

ПРОФЕССИОНАЛЬНО ОБ ORACLE

<IOUG>
independent oracle users group
share your experience

Oracle

ДЛЯ ПРОФЕССИОНАЛОВ

3-е издание

ТЕХНОЛОГИИ И РЕШЕНИЯ ДЛЯ ДОСТИЖЕНИЯ
ВЫСОКОЙ ПРОИЗВОДИТЕЛЬНОСТИ
И ЭФФЕКТИВНОСТИ

Томас Кайт, Дарл Кун



www.williamspublishing.com

Apress®
www.apress.com

Oracle

ДЛЯ ПРОФЕССИОНАЛОВ

3-е издание

Expert Oracle Database Architecture

Third Edition

Thomas Kyte
Darl Kuhn

Apress®

Oracle

ДЛЯ ПРОФЕССИОНАЛОВ

**Архитектура, методики
программирования и основные
особенности версий 9i, 10g, 11g и 12c**

3-е издание

**Томас Кайт
Дарл Кун**



**Москва • Санкт-Петербург • Киев
2016**

ББК 32.973.26-018.2.75

K15

УДК 681.3.07

Издательский дом “Вильямс”

Зав. редакцией *С.Н. Тригуб*

Перевод с английского *Ю.Н. Артеменко*

Под редакцией *Ю.Н. Артеменко*

По общим вопросам обращайтесь в Издательский дом “Вильямс” по адресу:
info@williamspublishing.com, <http://www.williamspublishing.com>

Кайт, Томас, Кун, Дарл.

K15 Oracle для профессионалов: архитектура и методики программирования, 3-е изд.
: Пер. с англ. — М. : ООО “И.Д. Вильямс”, 2016. — 960 с. : ил. — Парал. тит. англ.
ISBN 978-5-8459-2042-3 (рус.)

ББК 32.973.26-018.2.75

Все названия программных продуктов являются зарегистрированными торговыми марками соответствующих фирм.

Никакая часть настоящего издания ни в каких целях не может быть воспроизведена в какой бы то ни было форме и какими бы то ни было средствами, будь то электронные или механические, включая фотокопирование и запись на магнитный носитель, если на это нет письменного разрешения издательства APress, Berkeley, CA.

Authorized translation from the English language edition published by APress, Inc., Copyright © 2014 by Thomas Kyte and Darl Kuhn.

All rights reserved. No part of this work may be reproduced or transmitted in any form or by any means, electronic or mechanical, including photocopying, recording, or by any information storage or retrieval system, without the prior written permission of the copyright owner and the publisher.

Russian language edition is published by Williams Publishing House according to the Agreement with R&I Enterprises International, Copyright © 2016.

Научно-популярное издание

Томас Кайт, Дарл Кун

**Oracle для профессионалов:
архитектура и методики программирования,
3-е издание**

Верстка *Т.Н. Артеменко*

Художественный редактор *В.Г. Павлютин*

Подписано в печать 04.11.2015. Формат 70×100/16

Гарнитура Times

Усл. печ. л. 77,4. Уч.-изд. л. 62,5

Тираж 400 экз. Заказ № 6739

Отпечатано способом ролевой струйной печати

в АО «Первая Образцовая типография»

Филиал «Чеховский Печатный Двор»

142300, Московская область, г. Чехов, ул. Полиграфистов, д.1

ООО “И. Д. Вильямс”, 127055, г. Москва, ул. Лесная, д. 43, стр. 1

ISBN 978-5-8459-2042-3 (рус.)

ISBN 978-1-4302-6298-5 (англ.)

© Издательский дом “Вильямс”, 2016

© by Thomas Kyte and Darl Kuhn, 2014

Оглавление

| | |
|--|-----|
| Об авторах | 13 |
| Благодарности | 15 |
| Введение | 16 |
| Настройка среды | 24 |
| Глава 1. Разработка успешных приложений Oracle | 45 |
| Глава 2. Обзор архитектуры | 113 |
| Глава 3. Файлы | 137 |
| Глава 4. Структуры памяти | 205 |
| Глава 5. Процессы Oracle | 259 |
| Глава 6. Блокировка и защелкивание данных | 299 |
| Глава 7. Параллелизм и многоверсионность | 359 |
| Глава 8. Транзакции | 389 |
| Глава 9. Повтор и отмена | 427 |
| Глава 10. Таблицы базы данных | 487 |
| Глава 11. Индексы | 587 |
| Глава 12. Типы данных | 671 |
| Глава 13. Секционирование | 749 |
| Глава 14. Параллельное выполнение | 837 |
| Глава 15. Загрузка и выгрузка данных | 881 |
| Предметный указатель | 954 |

Содержание

| | |
|--|------------|
| Об авторах | 13 |
| Благодарности | 15 |
| Введение | 16 |
| О чем эта книга | 16 |
| Кому адресована эта книга | 17 |
| Как структурирована эта книга | 19 |
| Исходный код и обновления | 24 |
| От издательства | 24 |
| Настройка среды | 24 |
| Настройка учетной записи EODA | 25 |
| Настройка схемы SCOTT/TIGER | 26 |
| Выполнение сценария | 26 |
| Создание схемы без сценария | 27 |
| Настройка среды | 28 |
| Настройка средства AUTOTRACE в SQL*Plus | 29 |
| Настройка пакета Statspack | 30 |
| Специальные сценарии | 31 |
| Соглашения при написании кода | 44 |
| Глава 1. Разработка успешных приложений Oracle | 45 |
| Мой подход | 47 |
| Метод черного ящика | 48 |
| Как следует (и как не следует) разрабатывать приложения баз данных | 59 |
| Архитектура Oracle | 59 |
| Управление параллельной обработкой | 73 |
| Многоверсионность | 78 |
| Как заставить приложение выполняться быстрее? | 107 |
| Отношения между администратором базы данных и разработчиком | 110 |
| Резюме | 111 |
| Глава 2. Обзор архитектуры | 113 |
| Определение базы данных и экземпляра | 114 |
| Системная глобальная область и фоновые процессы | 123 |
| Подключение к Oracle | 125 |
| Выделенный сервер | 126 |
| Разделяемый сервер | 128 |
| Механизмы подключения через TCP/IP | 129 |
| Подключаемые базы данных | 132 |
| Снижение коэффициента использования ресурсов | 134 |
| Сокращение объема работ по обслуживанию | 134 |
| Отличия подключаемой базы данных | 135 |
| Резюме | 136 |
| Глава 3. Файлы | 137 |
| Файлы параметров | 138 |
| Что собой представляют параметры | 139 |
| Унаследованные файлы параметров init.ora | 143 |

| | |
|---|------------|
| Файлы параметров сервера (SPFILE) | 146 |
| Заключительные соображения по поводу файла параметров | 154 |
| Трассировочные файлы | 155 |
| Запрошенные трассировочные файлы | 157 |
| Трассировочные файлы, генерируемые в ответ на внутренние ошибки | 162 |
| Заключительные соображения по поводу трассировочных файлов | 168 |
| Сигнальный файл | 168 |
| Файлы данных | 172 |
| Краткий обзор механизмов файловой системы | 172 |
| Иерархия хранения в базе данных Oracle | 174 |
| Табличные пространства, управляемые словарем и управляемые локально | 179 |
| Временные файлы | 182 |
| Управляющие файлы | 184 |
| Файлы журналов повторения транзакций | 184 |
| Оперативный журнал повторения транзакций | 186 |
| Архивный журнал повторения транзакций | 189 |
| Файлы паролей | 191 |
| Файл отслеживания изменений | 194 |
| Ретроспективные журналы | 196 |
| Команда FLASHBACK DATABASE | 196 |
| Область для быстрого восстановления | 197 |
| Файлы DMP (файлы экспорта/импорта) | 198 |
| Файлы Data Pump | 200 |
| Плоские файлы | 202 |
| Резюме | 203 |
| Глава 4. Структуры памяти | 205 |
| Глобальная область процесса и глобальная область пользователя | 206 |
| Ручное управление памятью PGA | 208 |
| Автоматическое управление памятью PGA | 215 |
| Выбор между ручным и автоматическим управлением памятью | 228 |
| Заключительные соображения по поводу использования областей PGA и UGA | 230 |
| Системная глобальная область | 230 |
| Фиксированная область SGA | 236 |
| Буфер повторения | 237 |
| Кеш буферов блоков | 238 |
| Разделяемый пул | 246 |
| Большой пул | 249 |
| Пул Java | 251 |
| Пул Streams | 252 |
| Управление памятью SGA | 252 |
| Резюме | 258 |
| Глава 5. Процессы Oracle | 259 |
| Серверные процессы | 260 |
| Подключения посредством выделенного сервера | 260 |
| Подключения посредством разделяемого сервера | 262 |
| Резидентный пул соединений с базой данных | 264 |

| | |
|--|------------|
| Подключения и сеансы | 264 |
| Сравнение режимов выделенного сервера, разделяемого сервера и DRCP | 271 |
| DRCP | 276 |
| Заключительные соображения по поводу выделенного/разделяемого сервера | 276 |
| Фоновые процессы | 277 |
| Специализированные фоновые процессы | 278 |
| Служебные фоновые процессы | 292 |
| Подчиненные процессы | 296 |
| Подчиненные процессы ввода-вывода | 296 |
| Pnnn: серверы выполнения параллельного запроса | 297 |
| Резюме | 298 |
| Глава 6. Блокировка и зашелкивание данных | 299 |
| Понятие блокировки | 299 |
| Проблемы блокировки | 303 |
| Потерянные обновления | 303 |
| Пессимистическая блокировка | 305 |
| Оптимистическая блокировка | 307 |
| Выбор между оптимистической и пессимистической блокировкой | 315 |
| Взаимоблокировки | 320 |
| Эскалация блокировок | 326 |
| Типы блокировок | 326 |
| Блокировки DML | 327 |
| Блокировки DDL | 339 |
| Защелки | 344 |
| Блокировка вручную и блокировки, определенные пользователем | 356 |
| Резюме | 357 |
| Глава 7. Параллелизм и многоверсионность | 359 |
| Понятие управления параллелизмом | 359 |
| Уровни изоляции транзакций | 361 |
| Уровень изоляции READ UNCOMMITTED | 363 |
| Уровень изоляции READ COMMITTED | 364 |
| Уровень изоляции REPEATABLE READ | 366 |
| Уровень изоляции SERIALIZABLE | 369 |
| Уровень изоляции READ ONLY | 372 |
| Последствия многоверсионной согласованности чтения | 373 |
| Распространенный прием организации хранилищ данных, который не работает | 374 |
| Объяснение неожиданно высокой активности ввода-вывода в “горячих” таблицах | 375 |
| Согласованность записи | 378 |
| Согласованные чтения и текущие чтения | 378 |
| Наблюдение за перезапуском | 381 |
| Почему перезапуск важен для нас? | 385 |
| Резюме | 386 |
| Глава 8. Транзакции | 389 |
| Операторы управления транзакциями | 390 |

| | |
|--|-----|
| Атомарность | 392 |
| Атомарность на уровне оператора | 392 |
| Атомарность на уровне процедуры | 394 |
| Атомарность на уровне транзакции | 398 |
| Операторы DDL и атомарность | 398 |
| Постоянство | 399 |
| Расширения WRITE оператора COMMIT | 399 |
| Операторы COMMIT в нераспределенном блоке PL/SQL | 401 |
| Ограничения целостности и транзакции | 403 |
| Ограничения IMMEDIATE | 403 |
| Ограничения DEFERRABLE и каскадные обновления | 404 |
| Плохие привычки в отношении транзакций | 408 |
| Фиксация в цикле | 409 |
| Использование автоматической фиксации | 417 |
| Распределенные транзакции | 418 |
| Автономные транзакции | 420 |
| Как работают автономные транзакции | 421 |
| Когда использовать автономные транзакции | 423 |
| Резюме | 426 |
| Глава 9. Повтор и отмена | 427 |
| Что собой представляет redo | 428 |
| Что собой представляет undo | 429 |
| Совместная работа redo и undo | 433 |
| Пример сценария INSERT-UPDATE-DELETE-COMMIT | 433 |
| Обработка COMMIT и ROLLBACK | 438 |
| Что делает оператор COMMIT? | 439 |
| Что делает оператор ROLLBACK? | 446 |
| Исследование redo | 448 |
| Измерение redo | 448 |
| Можно ли отключить генерацию журналов redo? | 450 |
| Почему не удастся разместить новый журнал? | 454 |
| Очистка блоков | 456 |
| Конкуренция за журнал | 461 |
| Временные таблицы и redo/undo | 464 |
| Исследование undo | 470 |
| Что генерирует максимальный и минимальный объем информации undo? | 470 |
| Ошибка ORA-01555: snapshot too old | 473 |
| Резюме | 486 |
| Глава 10. Таблицы базы данных | 487 |
| Типы таблиц | 487 |
| Терминология | 490 |
| Сегмент | 490 |
| Управление пространством сегментов | 493 |
| Маркер максимального уровня заполнения | 494 |
| Списки свободных блоков | 496 |
| Параметры PCTFREE и PCTUSED | 501 |

| | |
|--|------------|
| Параметры LOGGING и NOLOGGING | 505 |
| Параметры INITRANS и MAXTRANS | 505 |
| Традиционные таблицы | 505 |
| Индекс-таблицы | 509 |
| Заключительные соображения по поводу индекс-таблиц | 526 |
| Кластеризованные индекс-таблицы | 527 |
| Заключительные соображения по поводу кластеризованных индекс-таблиц | 535 |
| Кластеризованные хеш-таблицы | 536 |
| Заключительные соображения по поводу кластеризованных хеш-таблиц | 545 |
| Отсортированные кластеризованные хеш-таблицы | 546 |
| Вложенные таблицы | 549 |
| Синтаксис вложенных таблиц | 550 |
| Хранение вложенных таблиц | 559 |
| Заключительные соображения по поводу вложенных таблиц | 562 |
| Временные таблицы | 563 |
| Сбор статистики до версии Oracle 12c | 568 |
| Сбор статистики, начиная с версии Oracle 12c | 572 |
| Заключительные соображения по поводу временных таблиц | 577 |
| Объектные таблицы | 577 |
| Заключительные соображения по поводу объектных таблиц | 585 |
| Резюме | 585 |
| Глава 11. Индексы | 587 |
| Обзор индексов Oracle | 588 |
| Индексы со структурой В-дерева | 590 |
| Сжатие ключей индекса | 593 |
| Индексы по реверсированным ключам | 596 |
| Индексы, упорядоченные по убыванию | 603 |
| Когда должен использоваться индекс со структурой В-дерева? | 605 |
| Заключительные соображения по поводу индексов со структурой В-дерева | 617 |
| Битовые индексы | 618 |
| Когда должен использоваться битовый индекс? | 619 |
| Битовые индексы соединений | 624 |
| Заключительные соображения по поводу битовых индексов | 627 |
| Индексы на основе функций | 627 |
| Простой пример индекса на основе функции | 628 |
| Индексация только некоторых строк | 638 |
| Реализация выборочной уникальности | 640 |
| Предостережение относительно ошибки ORA-01743 | 640 |
| Заключительные соображения по поводу индексов на основе функций | 642 |
| Индексы предметной области | 642 |
| Невидимые индексы | 644 |
| Множество индексов на одной и той же комбинации столбцов | 646 |
| Индексация расширенных столбцов | 646 |
| Решение с виртуальным столбцом | 648 |
| Решение с индексом на основе функции | 650 |
| Часто задаваемые вопросы и мифы об индексах | 651 |
| Работают ли индексы в представлениях? | 651 |

| | |
|---|-----|
| Могут ли значения NULL и индексы работать вместе? | 652 |
| Должны ли быть проиндексированы внешние ключи? | 655 |
| Почему индекс не используется? | 656 |
| Миф: пространство никогда не используется в индексе повторно | 663 |
| Миф: наиболее отличительные столбцы должны быть в индексе первыми | 667 |
| Резюме | 670 |
| Глава 12. Типы данных | 671 |
| Обзор типов данных Oracle | 671 |
| Символьные и двоичные строковые типы | 675 |
| Обзор NLS | 675 |
| Символьные строки | 679 |
| Двоичные строки: типы RAW | 686 |
| Расширенные типы данных | 689 |
| Числовые типы | 692 |
| Синтаксис и использование типа NUMBER | 694 |
| Синтаксис и использование типов BINARY_FLOAT/BINARY_DOUBLE | 698 |
| Несобственные числовые типы | 699 |
| Соображения по поводу производительности | 699 |
| Типы LONG | 701 |
| Ограничения типов LONG и LONG RAW | 701 |
| Копирование с участием унаследованных типов LONG | 702 |
| Типы DATE, TIMESTAMP и INTERVAL | 708 |
| Форматы | 709 |
| Тип DATE | 710 |
| Вариации типа TIMESTAMP | 716 |
| Тип INTERVAL | 724 |
| Типы LOB | 727 |
| Внутренние типы LOB | 728 |
| Тип BFILE | 744 |
| Типы ROWID и UROWID | 745 |
| Резюме | 747 |
| Глава 13. Секционирование | 749 |
| Обзор секционирования | 750 |
| Повышенная доступность | 750 |
| Облегчение задач администрирования | 753 |
| Улучшенная производительность операторов | 758 |
| Сокращение конкуренции в системе OLTP | 760 |
| Схемы секционирования таблиц | 760 |
| Секционирование по диапазонам ключей | 762 |
| Хеш-секционирование | 765 |
| Секционирование по списку значений ключа | 769 |
| Секционирование по интервалам ключей | 771 |
| Секционирование по ссылкам | 778 |
| Секционирование по интервалам ключей и по ссылкам | 783 |
| Секционирование по виртуальному столбцу | 785 |
| Составное секционирование | 787 |

12 Содержание

| | |
|--|------------|
| Перемещение строк | 789 |
| Заключительные соображения по поводу схем секционирования таблиц | 792 |
| Секционирование индексов | 793 |
| Сравнение локальных и глобальных индексов | 794 |
| Локальные индексы | 795 |
| Глобальные индексы | 801 |
| Системы OLTP и глобальные индексы | 813 |
| Частичные индексы | 818 |
| Еще раз о секционировании и производительности | 820 |
| Удобство средств обслуживания | 827 |
| Множественные операции обслуживания секций | 827 |
| Каскадное усечение | 830 |
| Каскадный обмен | 832 |
| Аудит и сжатие пространства сегментов | 834 |
| Резюме | 836 |
| Глава 14. Параллельное выполнение | 837 |
| Использование параллельного выполнения | 839 |
| Аналогия параллельной обработки | 840 |
| Oracle Exadata | 842 |
| Параллельный запрос | 843 |
| Параллельный DML | 849 |
| Параллельный DDL | 854 |
| Параллельный DDL и загрузка данных с использованием внешних таблиц | 855 |
| Параллельный DDL и усечение экстендов | 857 |
| Процедурный параллелизм | 867 |
| Параллельные конвейерные функции | 868 |
| Самодельный параллелизм | 871 |
| Самодельный параллелизм старой школы | 875 |
| Резюме | 879 |
| Глава 15. Загрузка и выгрузка данных | 881 |
| Внешние таблицы | 881 |
| Настройка внешних таблиц | 883 |
| Обработка ошибок | 893 |
| Использование внешней таблицы для загрузки разных файлов | 896 |
| Проблемы многопользовательского доступа | 896 |
| Предварительная обработка | 898 |
| Заключительные соображения по поводу внешних таблиц | 908 |
| Выгрузка Data Pump | 909 |
| Инструмент SQLLDR | 910 |
| Часто задаваемые вопросы по загрузке данных посредством SQLLDR | 915 |
| Предостережения относительно SQLLDR | 943 |
| Заключительные соображения по поводу SQLLDR | 943 |
| Выгрузка в плоский файл | 944 |
| Резюме | 953 |
| Предметный указатель | 954 |

Об авторах

Меня зовут **Том Кайт**. Я работал на компанию Oracle со времен версии 7.0.9 (или с 1993 года — для тех, кто не ведет отсчет времени по версиям Oracle). Однако с самим продуктом Oracle я имел дело примерно с версии 5.1.5c (однопользовательская версия для DOS, распространявшаяся на дискетах 360 Кбайт по цене \$99). До перехода в Oracle я более шести лет я проработал специалистом по интеграции систем, строя крупномасштабные гетерогенные базы данных и приложения, в основном предназначенные для военных и правительственных организаций. В то время мне приходилось тратить немало времени на СУБД Oracle, в частности, оказывая помощь тем, кто использовал базы данных Oracle. Я непосредственно контактировал с заказчиками, либо определяя требования и строя их системы, либо чаще помогая их перестраивать и настраивать (“настройка” нередко является синонимом “перестройки”). Кроме того, я тот самый Том, который ведет колонку “Ask Tom” (“Спросите у Тома”) в журнале *Oracle Magazine*, отвечая на вопросы по базам данных и инструментам Oracle. Обычно в течение дня на веб-сайте <http://asktom.oracle.com> я получаю и даю ответы на десятки вопросов. Каждые два месяца я публикую лучшие из них в журнале (все заданные вопросы доступны в Интернете; естественно, они хранятся в базе данных Oracle). Вдобавок я веду технические семинары по многим темам, которые вы найдете в настоящей книге. В целом очень много своего времени я трачу, помогая людям успешно работать с базами данных Oracle. И да, в свободное время я строю приложения и разрабатываю программное обеспечение в самой компании Oracle.

Эта книга отражает то, что я делаю каждый день. Представленный в ней материал раскрывает темы и вопросы, с которыми люди сталкиваются ежедневно. Эти вопросы освещены с точки зрения “Когда я использую это, то поступаю вот так”. Книга является кульминацией многолетнего опыта использования этого продукта в неслетном числе ситуаций.

Дарл Кун — администратор баз данных и разработчик, работающий на компанию Oracle. Он также преподает курсы по Oracle в Университете Реджис (Денвер, штат Колорадо) и является активным членом ассоциации пользователей Скалистых гор (Rocky Mountain Oracle Users Group). Дарл с удовольствием делится знаниями, что привело к реализации нескольких книжных проектов на протяжении последних лет.

О рецензентах

Мелани Кэффри — старший менеджер по разработке в Oracle Corporation, продвигающий клиентские и серверные решения Oracle для бизнес-потребностей разнообразных клиентов. Является соавтором нескольких технических публикаций, включая книгу *Expert Oracle Practice: Oracle Database Administration from the Oak Table*, вышедшую в издательстве Apress, а также *Oracle Web Application Programming for PL/SQL Developers*, *The Oracle DBA Interactive Workbook* и *Oracle Database Administration: The Complete Video Course*, опубликованные издательством Prentice Hall. Она преподает в Нью-Йорке студентам программу Колумбийского университета “Компьютерные технологии и приложения”, обучая их расширенному администрированию и разработке на PL/SQL. Кроме того, Мелани — частый докладчик на конференциях Oracle.

Кристофер Бек получил диплом Ратгерского университета и работал со многими СУБД более 19 лет. Последние 15 лет он является сотрудником компании Oracle, где занимает должность главного технолога, ориентируясь на основные технологии баз данных. Кристофер — соавтор двух патентов США в области программных методологий, которые стали основой того, что теперь известно как Oracle Application Express. Он рецензировал и другие книги по Oracle, включая первую книгу Тома — *Expert One-to-One* (Apress, 2003 г.), а также сам выступал в роли соавтора двух книг — *Beginning Oracle Programming* (Apress, 2003 г.) и *Mastering Oracle PL/SQL* (Apress, 2004 г.). Кристофер живет в Северной Вирджинии со своей женой Мартой и четырьмя детьми, и когда не тратит время на них, обычно играет в видеоигры или смотрит футбол.

Фриц Хугланд — профессионал в области информационных технологий, специализирующийся на производительности и внутреннем устройстве баз данных Oracle. Фриц часто презентует технические темы Oracle на конференциях по всему миру. В 2009 году он получил премию Oracle ACE от Oracle Technology Network и спустя год стал руководителем программы Oracle ACE. В 2010 году Фриц присоединился к OakTable Network. В дополнение к опыту разработки для Oracle он хорошо ориентируется в PostgreSQL и современных операционных системах. В настоящее время Фриц работает в компании Accenture из группы Accenture Enkitec.

Алекс Фаткулин — специалист по всему спектру технологий Oracle. Его мастерство было жизненно важным при решении ряда огромных проблем, с которыми столкнулись его заказчики.

Алекс опирается на годы опыта работы с крупнейшими мировыми компаниями, где ему приходилось иметь дело практически со всем, что относится к базам данных Oracle — от моделирования данных до построения архитектур в рамках решений высокой готовности и устранения проблем с производительностью исключительно больших производственных площадок.

Он получил степень бакалавра по вычислительной технике в Дальневосточном государственном университете (Владивосток, Россия), а также является членом Oracle ACE и OakTable.

Стефан Фаруль изучал язык SQL в 1983 году с помощью СУБД SQL/DS производства IBM. В 1986 году после полутора лет учебы в Оттавском университете (Канада) он присоединился к очень маленькой команде только что созданной компании Oracle France и затем проработал год в IBM France. В 1988 году Стефан оставил Oracle, занявшись разработкой программного обеспечения, и вернулся к базам данных в августе 1989 года. С тех пор он имел дело с базами данных, главным образом с Oracle, консультируя крупные французские компании на протяжении 25 лет. Стефан опубликовал несколько книг по SQL, вел семинары по Oracle в Азии, выкладывал учебные пособия по базам данных в YouTube и проводил для французских студентов учебные курсы по базам данных, одновременно консультируя один крупный банк. Получая истинное удовольствие от обучения, в настоящее время он преподает вычислительную технику в Университете штата Канзас.

Благодарности

Я хотел бы поблагодарить многих людей за оказанную ими помощь в завершении этой книги. Прежде всего, я хочу выразить благодарность вам, читателям этой книги. Если вы читаете данную книгу, то высока вероятность того, что вы тем или иным образом поучаствовали в работе веб-сайта <http://asktom.oracle.com/>, возможно, задав один или два вопроса. Такое действие — получение вопросов и исследование ответов — снабжает меня материалом для книги и знаниями, на которых основан этот материал. Без ваших вопросов я бы никогда не знал о СУБД Oracle столько, сколько знаю сейчас. Так что в конечном итоге именно вы сделали возможным появление настоящей книги.

Я хотел бы поблагодарить Тони Дэвиса за его усилия по приведению моей работы в действительно читабельный вид. Если вам нравится последовательность и разбиение книги на разделы, и вы находите представление материала четким, то в значительной степени это достигнуто стараниями Тони. Я работаю с Тони при написании технических материалов, начиная с 2000 года, и вижу, насколько выросли его знания Oracle за это время. Он способен выполнять не только литературное редактирование материала, но во многих случаях и научное редактирование. Многие из примеров, представленные в этой книге, появились благодаря именно ему (без них случайный читатель просто не смог бы разобратся в материале). Тони не редактировал второе издание книги, но содержимое во многом основано на первом издании. Без него настоящая книга не стала бы такой, как она есть.

Если бы не группа квалифицированных рецензентов, с которыми я сотрудничал во время написания этой книги и ее предыдущего издания, то мне пришлось бы переживать за ее содержимое. Первое издание рецензировали Джонатан Льюис, Родерик Маналак, Майкл Меллер и Гэйб Романеску. Они потратили массу времени на проверку правильности материала с технической точки зрения и оценку его полезности в реальной практике. Вторым изданием занималась команда отнюдь не меньшего калибра: Мелани Кэффри, Кристофер Бек и Джейсон Страуб. Я глубоко убежден, что техническая книга должна оцениваться не только теми, кто ее написал, но и сторонними рецензентами. Благодаря этим людям, я чувствую себя уверенным относительно изложенного материала.

В компании Oracle я работаю с самыми лучшими и одаренными людьми из всех, с кем мне доводилось встречаться, и каждый из них, так или иначе, внес свой вклад в создание этой книги. В частности, я хотел бы поблагодарить Кена Якобса за его поддержку и энтузиазм на протяжении многих лет. К сожалению, Кен больше не работает в Oracle, но его влияние сохранится надолго.

И последнее, но самое важное — я хочу выразить признательность за постоянную поддержку членам своей семьи. Вы знаете, насколько важно чувствовать себя кому-то нужным, когда вы пытаетесь сделать что-то, требующее траты многих часов “в нерабочее время”. Не представляю, как бы я смог завершить эту книгу без постоянной поддержки со стороны своей жены Мелани (которая также была техническим рецензентом книги), сына Алана и дочери Меган.

Том Кайт

Я хотел бы поблагодарить Тома за то, что он пригласил меня поработать с ним над этой книгой; это большая честь для меня. Также я хочу выразить благодарность Джонатану Геннику; его руководство (на протяжении многих лет и книг) заложило для меня основу, которая позволила быть в состоянии трудиться над книгой подобного уровня. И я также благодарен Хейди, Лизе и Бренди; без вашей поддержки я не смог бы успешно поучаствовать в настоящем проекте.

Дарл Кун

Введение

Вдохновение изложить собранный в этой книге материал появилось у меня вследствие приобретенного опыта разработки программного обеспечения Oracle, сотрудничества с другими разработчиками Oracle и оказания им помощи в построении надежных приложений на основе баз данных Oracle. В целом эта книга отражает то, чем мне приходится заниматься ежедневно, и в ней освещены вопросы, с которыми пользователи сталкиваются в своей повседневной работе.

Я раскрыл здесь то, что считаю наиболее важным, а именно — базы данных Oracle и их архитектуру. Я мог бы написать аналогично озаглавленную книгу, посвященную разработке приложений с использованием конкретного языка и архитектуры — например, приложений, которые применяют JavaServer Pages для взаимодействия с компонентами Enterprise JavaBean (EJB), которые, в свою очередь, используют JDBC для обмена данными с Oracle. Однако, по большому счету, для успешного построения таких приложений вы действительно должны хорошо разбираться во всех темах, рассмотренных в настоящей книге. В ней собран материал, который, по моему глубокому убеждению, должен быть всесторонне усвоен для успешной разработки приложений в среде Oracle, будь вы программистом на Visual Basic, применяющим ODBC, программистом на Java, использующим компоненты EJB и JDBC, или программистом на Perl, использующим DBI-интерфейс Perl. В книге не продвигается какая-то специфичная архитектура приложений; в ней не сравниваются трехуровневая и клиент-серверная архитектуры. Вместо этого в книге объясняется, что может делать база данных, и что вы должны понимать в плане особенностей ее функционирования. Поскольку база данных находится в основе архитектуры любого приложения, книга должна иметь широкую аудиторию.

Как следует из названия, эта книга сосредоточена на архитектуре базы данных и работе самой базы данных. Я подробно раскрываю архитектуру базы данных Oracle — файлы, структуры памяти и процессы, которые образуют базу данных и экземпляр. Затем я перехожу к обсуждению таких важных тем, как блокировка, управление параллелизмом, функционирование транзакций, повторение и отмена, и поясняю, почему эти вопросы настолько важны. Наконец, я рассматриваю физические структуры базы данных, такие как таблицы, индексы и типы данных, освещая приемы их оптимального применения.

О чем эта книга

Одна из проблем, обусловленная наличием многочисленных возможностей при разработке, заключается в том, что иногда трудно выбрать вариант, который лучше всего подходит в конкретных обстоятельствах. Каждый желает получить максимально возможную гибкость (столько вариантов, сколько вообще можно иметь), но при этом хочет, чтобы все оставалось очень ясным — другими словами, простым. СУБД Oracle предлагает разработчикам практически неограниченные возможности выбора. Никто никогда вам не скажет: “Вы не сможете сделать это в Oracle”. Наоборот, вы услышите: “Сколькими разными способами вы предпочитаете делать это в Oracle?”. Я надеюсь, что данная книга поможет вам производить правильный выбор.

Книга адресована тем, кто высоко ценит наличие выбора, но также предпочитает располагать какими-то руководящими принципами и сведениями о практических реализациях функциональных средств и возможностей Oracle. Например, в Oracle

имеется по-настоящему эффективное средство, которое называется параллельным выполнением. В документации Oracle рассказано, как его использовать и что оно делает. Однако в ней ничего не говорится о том, когда вы должны и — что, пожалуй, еще важнее — когда не должны применять это средство. В документации далеко не всегда сообщаются детали реализации средства, и если вы не осведомлены о них, это может еще не раз дать о себе знать. (Здесь я не имею в виду программные ошибки, а способ, которым средство предположительно функционирует, и для чего оно в действительности было предназначено.)

В этой книге я стремился не только описать, как работает та или иная функция, но также и объяснить, когда и почему вы должны обдумать возможность использования отдельного средства или реализации. Я твердо убежден, что в отношении всех сущностей важно понимать не только “как”, но также “когда” и “почему”, а заодно “когда нет” и “почему нет”.

Кому адресована эта книга

К целевой аудитории этой книги относится любой, кто занимается разработкой приложений на основе баз данных Oracle. Она предназначена для профессиональных разработчиков Oracle, которым необходимо знать способы выполнения действий в базе данных. Практическая природа книги означает, что многие разделы должны быть очень интересными также и администраторам баз данных. В большинстве примеров для демонстрации ключевых средств применяется SQL*Plus, так что не надейтесь обнаружить в них по-настоящему крутой графический пользовательский интерфейс, но вы узнаете, как работает база данных Oracle, что могут делать ее ключевые средства и когда они должны (и не должны) использоваться.

Книга ориентирована на любого, кто желает получить от Oracle максимум, прикладывая минимальные усилия. Она предназначена для любого, кто ищет новые методы применения существующих средств. Она рассчитана на тех, кто хочет увидеть, как эти средства можно задействовать в реальном мире (не просто ознакомиться с примерами использования функции, но в первую очередь понять, почему она оказывается подходящей в той или иной ситуации). Еще одной категорией людей, кого может заинтересовать эта книга, являются технические руководители, ответственные за разработчиков, которые занимаются проектами Oracle. В некоторых отношениях очень важно, чтобы технические руководители понимали, почему знание базы данных имеет решающее значение в достижении успеха. Эта книга может стать серьезным подспорьем руководителям, которые хотят, чтобы их персонал был обучен надлежащим технологиям, или желают убедиться, что персонал уже знает то, что должен знать.

Чтобы извлечь максимальную пользу из этой книги, читатель должен соответствовать следующим характеристикам.

- **Знать язык SQL.** Вовсе не обязательно быть лучшим кодировщиком SQL во все времена, но хорошие практические навыки помогут в освоении.
- **Понимать PL/SQL.** Это не является обязательным условием, но облегчит усвоение примеров. Например, данная книга не научит вас применению цикла FOR или объявлению типа записи; вопросы такого рода подробно освещены в документации по Oracle и многочисленных книгах. Тем не менее, нельзя сказать, что вы не изучите много нового о PL/SQL, когда будете читать эту книгу.

Вы освоите разнообразные возможности PL/SQL, узнаете о новых способах решения задач и ознакомитесь с пакетами/функциями, о существовании которых, возможно, даже не подозревали.

- **Иметь представление о каком-нибудь языке третьего поколения, таком как C или Java.** Я уверен, что любой, кто способен читать и писать код на языке третьего поколения, сможет успешно понять примеры, приведенные в этой книге.
- **Быть знакомым с руководством по концепциям базы данных Oracle (Oracle Database Concepts).**

Последний пункт требует дополнительной пары слов: из-за значительного объема набора документации по Oracle многие находят ее устрашающе большой. Если вы только приступили к изучению или пока еще не прочитали ни одного документа из набора, то могу вас заверить, что руководство *Oracle Database Concepts* — это именно то, с чего следует начинать. Оно занимает около 450 страниц (я хорошо это знаю, т.к. некоторые страницы написаны мною, а редактировать мне довелось целиком все руководство) и затрагивает многие из основных концепций Oracle, о которых необходимо знать. Оно может не предоставлять абсолютно все технические детали (им посвящены остальные 10 000–20 000 страниц документации), но обучит вас всем важнейшим концепциям. В руководстве *Oracle Database Concepts* раскрыты перечисленные ниже темы (здесь упомянуты далеко не все):

- структуры в базе данных и способ организации и хранения данных;
- распределенная обработка;
- архитектура памяти Oracle;
- архитектура процессов Oracle;
- объекты схемы, которые вы будете использовать (таблицы, индексы, кластеры и т.д.);
- встроенные и определяемые пользователем типы данных;
- хранимые процедуры SQL;
- работа транзакций;
- оптимизатор;
- целостность данных;
- управление параллелизмом.

Время от времени я и сам буду возвращаться к этим темам. Они являются фундаментальными. Без их знания вы будете создавать приложения Oracle, предрасположенные к отказам. Я настоятельно рекомендую прочитать это руководство и усвоить упомянутые вопросы.

Как структурирована эта книга

Чтобы помочь в работе с книгой, большинство глав организовано в виде четырех основных разделов (как описано ниже). Хотя это не жесткое разделение, оно позволяет быстрее переходить к области, в которой требуется дополнительная информация. Книга содержит 15 глав, причем каждая из них подобна “мини-книге” —

практически самостоятельному компоненту. Иногда я ссылаюсь на примеры или средства, описанные в других главах, но любую главу книги вполне можно читать независимо от остальных. Например, вовсе не обязательно сначала читать главу 10, посвященную таблицам базы данных, чтобы понять материал главы 14, в которой рассматривается параллелизм.

Формат и стиль многих глав практически идентичен и представлен ниже.

- Введение в функциональное средство или возможность.
- Причины, по которым может потребоваться применение (или отказ от использования) этого средства или возможности. Здесь даются пояснения относительно того, когда следует подумать о применении этого средства, а когда отказаться от него.
- Как использовать это средство. Приводимые здесь сведения не являются просто копией материала из справочника по SQL. Напротив, они представлены в пошаговом виде: то, что необходимо получить, то, что понадобится сделать, и вопросы, которые нужно выяснить, чтобы начать. В этом разделе рассматриваются следующие темы:
 - способ реализации средства;
 - примеры применения;
 - способ отладки средства;
 - предостережения относительно использования средства;
 - способ обработки ошибок (упреждающим образом);
 - резюме с кратким подведением итогов.

Книга содержит большое количество примеров и кода, которые доступны для загрузки на веб-сайте издательства. В последующих разделах приведены краткие описания содержимого всех глав.

Глава 1. Разработка успешных приложений Oracle

В этой главе я описываю свой подход к программированию для баз данных. Базы данных не создавались одинаковыми. Для того чтобы успешно и своевременно разрабатывать приложения, управляемые базами данных, необходимо четко понимать, что именно конкретная база данных способна делать, и каким образом она это делает. Без знания того, что может делать база данных, возникает риск постоянно заново изобретать колесо, создавая средства, которые база данных уже предоставляет. Без понимания того, каким образом работает база данных, скорее всего, будут получаться приложения, которые функционируют неэффективно и ведут себя непредсказуемо.

В главе предлагается эмпирический взгляд на некоторые приложения, где недостаток понимания базы данных привел к неудаче всего проекта. В соответствии с таким ориентированным на примеры подходом здесь обсуждаются основные средства и функции базы данных, которые должны хорошо пониматься разработчиком. Суть в том, что вы не должны позволять себе трактовать базу данных как черный ящик, который просто выдает ответы и самостоятельно заботится о масштабируемости и производительности.

Глава 2. Обзор архитектуры

В этой главе раскрываются основы архитектуры Oracle. Мы начнем с четких определений двух терминов, которые многими в мире Oracle понимаются совершенно неправильно — *экземпляр* и *база данных*. Затем мы рассмотрим два новых типа баз данных, появившиеся в версии Oracle 12c, в частности — *контейнерная база данных* и *подключаемая база данных*. Кроме того, мы кратко обсудим системную глобальную область (System Global Area — SGA) и процессы, лежащие в основе экземпляра Oracle, а посмотрим, каким образом выполняется простое действие “подключения к Oracle”.

Глава 3. Файлы

В этой главе подробно описаны восемь типов файлов, которые образуют базу данных и экземпляр Oracle. Начиная с простого файла параметров и заканчивая файлами данных и журнальными файлами повторения, мы посмотрим, что они собой представляют, каково их назначение и каким образом они используются.

Глава 4. Структуры памяти

В этой главе освещены вопросы использования памяти базой данных Oracle, как памяти в отдельных процессах (PGA), так и разделяемой памяти (SGA). Мы исследуем отличия между ручным и автоматическим управлением памятью PGA, автоматическим управлением разделяемой памятью в Oracle 10g и автоматическим управлением памятью в Oracle 11g, а также выясним, в каких случаях подходит каждый из этих методов. После прочтения этой главы вы получите полное представление о том, каким образом Oracle работает с памятью и управляет ею.

Глава 5. Процессы Oracle

В этой главе предлагается обзор типов процессов Oracle (серверных и фоновых процессов). Кроме того, в ней более подробно описаны отличия между подключением к базе данных с помощью процессов разделяемого и выделенного серверов. Мы также рассмотрим большинство фоновых процессов (таких как LGWR, DBWR, PMON, SMON и LREG), действующих во время запуска экземпляра Oracle, и обсудим их функции.

Глава 6. Блокировка и защелкивание данных

В разных базах данных задачи решаются по-разному (то, что хорошо работает в SQL Server, может не так хорошо работать в Oracle), и понимание особенностей реализации блокировки и управления параллелизмом абсолютно необходимо для успешного построения приложений. В этой главе рассматривается общий подход, используемый в Oracle для решения этих задач, типы блокировок, которые могут быть применены (DML, DDL и защелки), а также проблемы, которые могут возникнуть при неаккуратной реализации блокировки (взаимоблокировка, блокирование и эскалация).

Глава 7. Параллелизм и многоверсионность

В этой главе мы будем исследовать мое любимое средство Oracle — многоверсионность, а также его влияние на управление параллелизмом и всю структуру приложения. Здесь будет показано, что все базы данных отличаются друг от друга, а сама их реализация может влиять на структуру приложений. Мы начнем с обзора раз-

нообразных уровней изоляции транзакций, как они определены стандартом ANSI SQL, и посмотрим, каким образом они отображаются на реализацию Oracle (а также соответствие этому стандарту других баз данных). Затем мы проанализируем последствия, которые может иметь для нас многоверсионность — средство, позволяющее Oracle обеспечивать неблокирующее чтение базы данных.

Глава 8. Транзакции

Транзакции являются фундаментальным средством всех баз данных — это часть того, что отличает базу данных от файловой системы. Несмотря на это, их часто понимают неправильно, и многие разработчики даже не подозревают о том, что неумышленно отказываются от их применения. В этой главе показано, как должны использоваться транзакции в Oracle, а также раскрыты плохие привычки, которые могли быть приобретены при разработке приложений для других баз данных. В частности, мы взглянем на последствия атомарности и ее влияние на операторы в Oracle. Мы также обсудим операторы управления транзакциями (COMMIT, SAVEPOINT и ROLLBACK), ограничения целостности, распределенные транзакции (двухфазную фиксацию) и автономные транзакции.

Глава 9. Повтор и отмена

Можно было бы сказать, что разработчикам не обязательно разбираться в деталях реализации повтора (redo) и отмены (undo) в такой же степени, как администраторам баз данных, но разработчики должны знать, какую роль играют повтор и отмена в базе данных. После определения понятия повтора мы посмотрим, что в точности делает оператор COMMIT. Мы обсудим, каким образом выяснить объем сгенерированной информации redo и как с помощью конструкции NOLOGGING значительно сократить объем информации redo, генерируемой определенными операциями. Мы также исследуем генерацию данных redo в связи с такими проблемами, как очистка блоков и конкуренция за журналы.

В разделе главы, посвященном отмене, мы рассмотрим роль данных undo и операции, которые генерируют наибольший/наименьший объем информации undo. Наконец, мы обсудим печально известную ошибку ORA-01555: snapshot too old (ORA-01555: устаревший снимок), возможные причины ее возникновения и способы ее предотвращения.

Глава 10. Таблицы базы данных

В настоящее время Oracle поддерживает множество типов таблиц. В этой главе мы по очереди рассмотрим типы таблиц — традиционная таблица (т.е. стандартная, “нормальная” таблица), индекс-таблица, кластеризованная индекс-таблица, кластеризованная хеш-таблица, вложенная таблица, временная таблица и объектная таблица — и обсудим когда, как и почему они должны использоваться. Большую часть времени традиционной таблицы вполне достаточно, но данная глава поможет распознать ситуации, в которых лучше подходят другие типы.

Глава 11. Индексы

Индексы являются критически важным аспектом структуры приложения. Корректная реализация требует глубоких знаний данных, их распределения и планируемого использования. Слишком часто при разработке приложений к индексам относятся как к чему-то второстепенному, из-за чего в итоге страдает производительность.

В этой главе подробно рассматриваются различные типы индексов, включая индекс со структурой B-дерева (B*Tree), битовый индекс, индекс на основе функций и индекс предметной области, и обсуждаются обстоятельства, когда они должны, а когда не должны применяться. В разделе, посвященном часто задаваемым вопросам и мифам об индексах, я также отвечу на ряд распространенных вопросов вроде “Работают ли индексы на представлениях?” и “Почему индекс не используется?”.

Глава 12. Типы данных

На выбор доступно большое количество типов данных. В этой главе исследуются все 22 встроенных типа данных, приводятся объяснения их реализации и рассматриваются способы и случаи их применения. В начале главы приводится краткий обзор поддержки национальных языков (National Language Support — NLS), знание которой обязательно для полного понимания простых строковых типов в Oracle. Затем мы перейдем к обсуждению вездесущего типа NUMBER. Далее мы раскроем типы LONG и LONG RAW главным образом с исторической точки зрения. Основная цель здесь — показать, как работать в приложениях с унаследованными столбцами LONG и переводить их в тип LOB. После этого мы углубимся в разнообразные типы данных для хранения значений даты и времени и уделим внимание манипулированию различными типами данных для получения нужного результата. Также будут описаны основы поддержки часовых поясов.

Следующими мы обсудим типы данных LOB. Мы рассмотрим способы их хранения и смысл каждого из многочисленных настроек, таких как IN ROW, CHUNK, RETENTION, CACHE и т.д. При работе с типами данных LOB важно понимать, как они реализованы и каким образом хранятся по умолчанию — особенно, когда дело доходит до настройки их извлечения и хранения. Глава завершается описанием типов ROWID и UROWID. Это специальные патентованные типы Oracle, которые представляют адрес строки. Мы объясним, когда их использовать в качестве типа данных столбца в таблице (что почти никогда не случается).

Глава 13. Секционирование

Секционирование предназначено для облегчения управления очень большими таблицами и индексами за счет реализации концепции “разделяй и властвуй”; в основном оно сводится к разбиению таблицы или индекса на множество мелких и более управляемых частей. Это область, в которой администратор базы данных и разработчик должны работать вместе с целью обеспечения максимальной доступности и производительности приложения. Здесь также подробно раскрыты функциональные средства, появившиеся в версиях Oracle 11g и Oracle 12c.

В главе рассматривается секционирование как таблиц, так и индексов. Будет описано секционирование с использованием локальных индексов (распространено в хранилищах данных) и глобальных индексов (распространено в системах OLTP).

Глава 14. Параллельное выполнение

В этой главе представлена концепция и случаи применения параллельного выполнения в Oracle. Мы начнем с рассмотрения ситуаций, когда параллельная обработка полезна и должна использоваться, а также ситуаций, когда она планироваться не должна. После этого мы перейдем к анализу механики параллельного запроса — средства, которое большинство людей ассоциирует с параллельным выполнением.

Затем мы опишем язык PDML (Parallel DML), который позволяет проводить модификации с использованием параллельного выполнения. Мы посмотрим, как язык PDML физически реализован, и почему эта реализация приводит к ряду ограничений, касающихся PDML.

Далее мы займемся исследованием параллельного DDL. По моему мнению, именно здесь параллельное выполнение проявляется во всей своей красе. Обычно администраторы баз данных располагают небольшими окнами обслуживания, в рамках которых должны выполнять крупные операции. Параллельный DDL предоставляет администраторам баз данных возможность в полной мере задействовать все доступные ресурсы компьютера, позволяя им завершать большие и сложные операции за какую-то долю времени, которое они отняли бы при последовательном выполнении.

Глава завершается рассмотрением процедурного параллелизма, т.е. средства, с помощью которого мы можем выполнять код приложения параллельно. Здесь раскрываются две технологии. Первая из них — параллельные конвейерные функции, или способность Oracle параллельно выполнять хранимые функции динамическим образом. Вторая технология — это “самодельный” параллелизм (do-it-yourself — DIY), когда приложение проектируется так, чтобы оно могло выполняться параллельно.

Глава 15. Загрузка и выгрузка данных

Внимание в первой половине этой главы сосредоточено на внешних таблицах — высокоэффективном средстве, с помощью которого осуществляется массовая загрузка и выгрузка данных. Если вы выполняете большой объем загрузки данных, то должны всерьез подумать о применении внешних таблиц. Здесь также подробно обсуждается средство предварительной обработки внешних таблиц, которое позволяет командам операционной системы автоматически выполняться как часть выборки из внешней таблицы.

Вторая половина главы посвящена SQL*Loader (SQLLDR) и разнообразным способам использования этого инструмента для загрузки и модификации данных в базе. Рассматриваются такие вопросы, как загрузка данных с разделителями, обновление существующих строк и вставка новых строк, выгрузка данных и обращение к SQLLDR из хранимой процедуры. SQLLDR — это тщательно продуманный и очень важный инструмент, но с его практическим применением связано немало вопросов.

Исходный код и обновления

Усваивать материал этой книги эффективнее всего путем тщательной проработки и разбора практических примеров. Работая с примерами, вы можете предпочесть вводить весь код вручную. Многие читатели выбирают такой подход потому, что он является хорошим способом ознакомиться с используемыми приемами написания кода.

Решите вы набирать код самостоятельно или нет, исходный код всех примеров, приведенных в настоящей книге, доступен на веб-сайте издательства. Если вы предпочитаете вводить код вручную, то файлы исходного кода можно применять для проверки ожидаемых результатов — это следует делать в первую очередь, когда есть подозрение, что код был набран с ошибкой. Если вам не нравится вводить код вручную, то без загрузки исходного кода с веб-сайта издательства не обойтись. В любом случае файлы кода помогут вам с обновлением и отладкой.

От издательства

Вы, читатель этой книги, и есть главный ее критик и комментатор. Мы ценим ваше мнение и хотим знать, что было сделано нами правильно, что можно было сделать лучше и что еще вы хотели бы увидеть изданным нами. Нам интересно услышать и любые другие замечания, которые вам хотелось бы высказать в наш адрес.

Мы ждем ваших комментариев и надеемся на них. Вы можете прислать нам бумажное или электронное письмо, либо просто посетить наш веб-сайт и оставить свои замечания там. Одним словом, любым удобным для вас способом дайте нам знать, нравится или нет вам эта книга, а также выскажите свое мнение о том, как сделать наши книги более интересными для вас.

Посылая письмо или сообщение, не забудьте указать название книги и ее авторов, а также ваш обратный адрес. Мы внимательно ознакомимся с вашим мнением и обязательно учтем его при отборе и подготовке к изданию последующих книг.

Наши координаты:

E-mail: info@williamspublishing.com

WWW: <http://www.williamspublishing.com>

Информация для писем из:

России: 127055, г. Москва, ул. Лесная, д. 43, стр. 1

Украины: 03150, Киев, а/я 152

Настройка среды

В этом разделе я объясню, как настроить среду, чтобы в ней можно было выполнять приведенные в книге примеры кода. В частности, рассматриваются следующие вопросы.

- Настройка учетной записи EODA, используемой во многих примерах.
- Корректная настройка демонстрационной схемы SCOTT/TIGER.
- Настройка и запуск необходимой среды.
- Конфигурирование средства AUTOTRACE в SQL*Plus.
- Установка пакета Statspack.
- Установка и запуск Runstats и других специальных утилит, применяемых в книге повсеместно.
- Соглашения по кодированию, используемые в этой книге.

Все сценарии, которые не поставляются Oracle, доступны для загрузки на веб-сайте издательства. В загруженном архиве со сценариями имеются папки chNN, содержащие сценарии для каждой главы (NN — это номер главы). Сценарии, перечисленные в разделе “Настройка среды”, находятся в папке ch00.

Настройка учетной записи EODA

Пользователь EODA участвует в большинстве примеров, рассматриваемых в этой книге. Это просто учетная запись, которой была выдана роль администратора базы данных и привилегии EXECUTE и SELECT на определенных объектах, принадлежащих SYS:

```
connect / as sysdba
define username=eoda
define usernamepwd=foo
create user &username identified by &usernamepwd;
grant dba to &username;
grant execute on dbms_stats to &username;
grant select on V_$STATNAME to &username;
grant select on V_$MYSTAT to &username;
grant select on V_$LATCH to &username;
grant select on V_$TIMER to &username;
conn &username/&usernamepwd
```

Вы можете настроить любого желаемого пользователя для запуска примеров в этой книге. Я выбрал в качестве имени пользователя EODA просто потому, что оно является аббревиатурой англоязычного оригинала настоящей книги.

Настройка схемы SCOTT/TIGER

Схема SCOTT/TIGER часто уже будет существовать в базе данных. Обычно она включена в типовую установку, но не является обязательным компонентом базы данных. Пример схемы SCOTT можно установить для любой учетной записи базы данных; с применением учетной записи SCOTT не связано ничего магического. При желании можете установить таблицы EMP/DEPT прямо в собственную учетную запись базы данных.

Многие примеры в этой книге полагаются на таблицы в схеме SCOTT. Чтобы поработать с такими примерами, вам понадобятся эти таблицы. Если вы имеете дело с совместно используемой базой данных, целесообразно установить собственную копию этих таблиц от имени учетной записи, отличной от SCOTT, чтобы избежать побочных эффектов, возникающих из-за того, что другие пользователи затрагивают те же самые данные.

Выполнение сценария

Чтобы создать демонстрационные таблицы схемы SCOTT, выполните следующие действия:

- введите команду `cd $ORACLE_HOME/sqlplus/demo;`
- после подключения от имени любого пользователя запустите сценарий `demobld.sql`.

На заметку! В Oracle 10g и последующих версиях демонстрационные подкаталоги должны устанавливаться из установочного носителя. Ниже будут воспроизведены необходимые компоненты `demobld.sql`.

Сценарий `demobld.sql` создаст и заполнит данными пять таблиц. По завершении он автоматически завершит сеанс SQL*Plus, поэтому не удивляйтесь, когда после выполнения сценария окно SQL*Plus исчезнет с экрана — так было задумано.

В стандартных демонстрационных таблицах никаких ограничений ссылочной целостности не определено. Некоторые примеры рассчитывают на то, что в них поддерживается ссылочная целостность. Поэтому после запуска сценария `demobld.sql` рекомендуется выполнить также следующие операторы:

```
alter table emp add constraint emp_pk primary key(empno);
alter table dept add constraint dept_pk primary key(deptno);
alter table emp add constraint emp_fk_dept foreign key(deptno) references dept;
alter table emp add constraint emp_fk_emp foreign key(mgr) references emp;
```

На этом установка демонстрационной схемы завершена. Если в какой-то момент вы захотите удалить эту схему с целью освобождения места на диске, можете просто запустить сценарий `$ORACLE_HOME/sqlplus/demo/demodrop.sql`. Он удалит пять таблиц и завершит сеанс SQL*Plus.

Совет. Команды SQL для создания и удаления пользователя SCOTT можно также найти в сценарии `$ORACLE_HOME/rdbms/admin/utlsampl.sql`.

Создание схемы без сценария

Если доступ к файлу `demobld.sql` отсутствует, для работы с приведенными в книге примерами достаточно выполнить следующие команды:

```
CREATE TABLE EMP
(EMPNO NUMBER(4) NOT NULL,
 ENAME VARCHAR2(10),
 JOB VARCHAR2(9),
 MGR NUMBER(4),
 HIREDATE DATE,
 SAL NUMBER(7, 2),
 COMM NUMBER(7, 2),
 DEPTNO NUMBER(2)
);

INSERT INTO EMP VALUES (7369, 'SMITH', 'CLERK', 7902,
TO_DATE('17-DEC-1980', 'DD-MON-YYYY'), 800, NULL, 20);
INSERT INTO EMP VALUES (7499, 'ALLEN', 'SALESMAN', 7698,
TO_DATE('20-FEB-1981', 'DD-MON-YYYY'), 1600, 300, 30);
INSERT INTO EMP VALUES (7521, 'WARD', 'SALESMAN', 7698,
TO_DATE('22-FEB-1981', 'DD-MON-YYYY'), 1250, 500, 30);
INSERT INTO EMP VALUES (7566, 'JONES', 'MANAGER', 7839,
TO_DATE('2-APR-1981', 'DD-MON-YYYY'), 2975, NULL, 20);
INSERT INTO EMP VALUES (7654, 'MARTIN', 'SALESMAN', 7698,
TO_DATE('28-SEP-1981', 'DD-MON-YYYY'), 1250, 1400, 30);
INSERT INTO EMP VALUES (7698, 'BLAKE', 'MANAGER', 7839,
TO_DATE('1-MAY-1981', 'DD-MON-YYYY'), 2850, NULL, 30);
INSERT INTO EMP VALUES (7782, 'CLARK', 'MANAGER', 7839,
TO_DATE('9-JUN-1981', 'DD-MON-YYYY'), 2450, NULL, 10);
INSERT INTO EMP VALUES (7788, 'SCOTT', 'ANALYST', 7566,
TO_DATE('09-DEC-1982', 'DD-MON-YYYY'), 3000, NULL, 20);
INSERT INTO EMP VALUES (7839, 'KING', 'PRESIDENT', NULL,
TO_DATE('17-NOV-1981', 'DD-MON-YYYY'), 5000, NULL, 10);
INSERT INTO EMP VALUES (7844, 'TURNER', 'SALESMAN', 7698,
TO_DATE('8-SEP-1981', 'DD-MON-YYYY'), 1500, 0, 30);
INSERT INTO EMP VALUES (7876, 'ADAMS', 'CLERK', 7788,
TO_DATE('12-JAN-1983', 'DD-MON-YYYY'), 1100, NULL, 20);
INSERT INTO EMP VALUES (7900, 'JAMES', 'CLERK', 7698,
TO_DATE('3-DEC-1981', 'DD-MON-YYYY'), 950, NULL, 30);
INSERT INTO EMP VALUES (7902, 'FORD', 'ANALYST', 7566,
TO_DATE('3-DEC-1981', 'DD-MON-YYYY'), 3000, NULL, 20);
INSERT INTO EMP VALUES (7934, 'MILLER', 'CLERK', 7782,
TO_DATE('23-JAN-1982', 'DD-MON-YYYY'), 1300, NULL, 10);

CREATE TABLE DEPT
(DEPTNO NUMBER(2),
 DNAME VARCHAR2(14),
 LOC VARCHAR2(13)
);

INSERT INTO DEPT VALUES (10, 'ACCOUNTING', 'NEW YORK');
INSERT INTO DEPT VALUES (20, 'RESEARCH', 'DALLAS');
INSERT INTO DEPT VALUES (30, 'SALES', 'CHICAGO');
INSERT INTO DEPT VALUES (40, 'OPERATIONS', 'BOSTON');
```


Если вы создали схему посредством показанных выше команд, не забудьте выполнить команды создания ограничений предыдущего раздела.

Настройка среды

Большинство примеров в этой книге предназначено для выполнения полностью в среде SQL*Plus. Поэтому настройку и конфигурирование требует только SQL*Plus. Однако у меня есть совет относительно применения SQL*Plus. Почти во всех примерах тем или иным образом используется пакет DBMS_OUTPUT. Чтобы можно было работать с DBMS_OUTPUT, потребуется выдать следующую команду SQL*Plus:

```
SQL> set serveroutput on
```

Частый ввод этой команды довольно быстро становится утомительным. К счастью, SQL*Plus позволяет создать файл login.sql — сценарий, который выполняется при каждом запуске SQL*Plus. Более того, можно определить переменную среды SQLPATH, что позволит находить этот сценарий независимо от того, в каком каталоге он хранится.

Для всех примеров в этой книги применяется такой сценарий login.sql:

```
define _editor=vi
set serveroutput on size unlimited
set trimspool on
set long 5000
set linesize 100
set pagesize 9999
column plan_plus_exp format a80
set sqlprompt '&_user. @&_connect_identifier.> '
```

Ниже приведена аннотированная версия этого сценария.

1. `define _editor=vi`. Определяет текстовый редактор, который SQL*Plus будет использовать по умолчанию. Можете указать любой предпочитаемый текстовый редактор (не текстовый процессор), такой как Notepad или emacs.
2. `set serveroutput on size unlimited`. По умолчанию включает пакет DBMS_OUTPUT (следовательно, вводить каждый раз команду `set serveroutput on` не придется). Также устанавливает стандартный размер буфера в максимально возможное значение.
3. `set trimspool on`. При буферизации текста строки будут усекаться с отбрасыванием пробелов, поэтому они не будут иметь фиксированную длину. Если параметр `trimspool` установлен в `off` (по умолчанию), ширина буферизованных строк будет равна значению параметра `linesize`.
4. `set long 5000`. Устанавливает стандартное количество байтов, отображаемых при выборе столбцов LONG и CLOB.
5. `set linesize 100`. Устанавливает ширину строк, отображаемых SQL*Plus, в 100 символов.
6. `set pagesize 9999`. Устанавливает параметр `pagesize`, который управляет тем, насколько часто SQL*Plus выводит заголовки, в большое значение (мы будем получать один набор заголовков на страницу).

7. `column plan_plus_exp format a80`. Устанавливает стандартную ширину строки в выводе плана выполнения, получаемого с помощью AUTOTRACE. В общем случае значения `a80` вполне достаточно для отображения полного плана.

Последний фрагмент сценария `login.sql` настраивает приглашение на ввод команд SQL*Plus:

```
set sqlprompt '&_user.&_connect_identifier.> '
```

В результате приглашение интерфейса приобретает следующий вид, который позволяет видеть имя пользователя и идентификатор подключения:

```
EODA@ORA12CR1>
```

Настройка средства AUTOTRACE в SQL*Plus

AUTOTRACE — это средство SQL*Plus, которое выводит план выполнения запущенных запросов и сведения об использованных ими ресурсах. В этой книге AUTOTRACE применяется очень часто. Существует несколько способов конфигурирования этого средства.

Начальная установка

Средство AUTOTRACE полагается на доступность таблицы по имени `PLAN_TABLE`. Начиная с версии Oracle 10g, схема `SYS` содержит глобальную временную таблицу под названием `PLAN_TABLE$`. Все необходимые привилегии для работы с этой таблицей выданы пользователю `PUBLIC` и определен открытый синоним (с именем `PLAN_TABLE`, который указывает на `SYS.PLAN_TABLE$`). Это значит, что получать доступ к этой таблице может любой пользователь.

На заметку! Если вы имеете дело с очень старой версией Oracle, можете создать таблицу `PLAN_TABLE` вручную, выполнив сценарий `$ORACLE_HOME/rdbms/admin/utlxplan.sql`.

Вы должны также создать и назначить роль `PLUSTRACE`:

- введите команду `cd $ORACLE_HOME/sqlplus/admin;`
- войдите в SQL*Plus от имени `SYS` или пользователя, которому выданы привилегии `SYSDBA`;
- введите команду `@plustrce;`
- введите команду `GRANT PLUSTRACE TO PUBLIC;`

При желании можете заменить `PUBLIC` в команде `GRANT` другим именем пользователя.

Управление отчетом

Отчет о пути выполнения, который используется оптимизатором SQL, и статистику по выполнению операторов можно получать автоматически. Отчет генерируется после успешного выполнения операторов SQL DML (т.е. `SELECT`, `DELETE`, `UPDATE`, `MERGE` и `INSERT`). Он полезен для отслеживания и настройки производительности перечисленных операторов.

Отчетом можно управлять посредством настройки системной переменной AUTOTRACE.

- SET AUTOTRACE OFF. Отчет AUTOTRACE не генерируется. Это принято по умолчанию.
- SET AUTOTRACE ON EXPLAIN. Отчет AUTOTRACE будет отображать только путь выполнения, применяемый оптимизатором.
- SET AUTOTRACE ON STATISTICS. Отчет AUTOTRACE будет отображать только статистику по выполнению SQL-операторов.
- SET AUTOTRACE ON. Отчет AUTOTRACE будет содержать путь выполнения, используемый оптимизатором, и статистику по выполнению SQL-операторов.
- SET AUTOTRACE TRACEONLY. Похоже на SET AUTOTRACE ON, но подавляет вывод запроса пользователя, если он есть.
- SET AUTOTRACE TRACEONLY EXPLAIN: Похоже на SET AUTOTRACE ON, но подавляет вывод запроса пользователя (если он есть) и также статистику по выполнению.

Настройка пакета Statspack

Пакет Statspack должен устанавливаться при подключении к базе данных от имени пользователя SYS (CONNECT / AS SYSDBA) или пользователя, которому выданы привилегии SYSDBA. Во многих средах установка Statspack будет задачей, требующей участия администратора базы данных или системного администратора.

Установка Statspack не представляет сложности. Вы просто запускаете сценарий @spcreate.sql. Он находится в каталоге \$ORACLE_HOME/rdbms/admin и должен выполняться после подключения от имени SYS через SQL*Plus.

Перед запуском сценария spcreate.sql необходимо знать следующие три единицы информации:

- пароль, который вы хотите применять для создаваемой схемы PERFSTAT;
- стандартное табличное пространство для схемы PERFSTAT;
- временное табличное пространство, которое вы хотите использовать для схемы PERFSTAT.

Выполнение этого сценария дает примерно такой вывод:

```
$ sqlplus / as sysdba

SQL*Plus: Release 12.1.0.1.0 Production on Fri May 23 15:45:05 2014
Copyright (c) 1982, 2013, Oracle. All rights reserved.
Connected to:
Oracle Database 12c Enterprise Edition Release 12.1.0.1.0 - 64bit Production
With the Partitioning, OLAP, Advanced Analytics and Real Application
Testing options
SYS@ORA12CR1> @spcreate

Choose the PERFSTAT user's password
-----
Not specifying a password will result in the installation FAILING
Enter value for perfstat_password:
```

Выберите пароль для пользователя PERFSTAT

Не указание пароля приведет к неудачному завершению установки

Введите значение для `perfstat_password`:

... <для краткости остальной вывод не показан> ...

Во время выполнения сценарий запросит необходимую информацию. В случае опечатки или случайного прекращения установки перед попыткой установить Statspack повторно вы должны с помощью сценария `spdrop.sql`, находящегося в `$ORACLE_HOME/rdbms/admin`, удалить пользователя и установленные представления. Процесс установки Statspack создаст файл `spcpkg.lis`. Его необходимо просмотреть на предмет выяснения любых ошибок, которые могли произойти. Однако пользователь, представления и код PL/SQL должны установиться без каких-либо проблем при условии, что вы указали допустимые имена табличных пространств (и пользователь PERFSTAT пока еще не существует).

Совет. Документация по Statspack находится в текстовом файле `$ORACLE_HOME/rdbms/admin/spdoc.txt`.

Специальные сценарии

В этом разделе будут описаны требования (если есть) со стороны разнообразных сценариев, применяемых в книге. Кроме того, мы исследуем их код.

Runstats

Runstats — это инструмент, который я разработал для сравнения двух разных методов выполнения того же действия и выбора лучшего из них. Вы передаете инструменту Runstats два метода, а он делает все остальное. Runstats производит следующие три ключевых измерения.

- **Время в формате часов, минут и секунд или пройденное время.** Эту информацию полезно знать, но она не является самой важной.
- **Статистика системы.** Отображает рядом данные о том, сколько раз каждый из подходов делал что-либо (например, обращался к синтаксическому анализатору), и разницу между ними.
- **Сведения о зашелкивании.** Это основная часть отчета.

Как будет показано в книге, зашелка представляет собой легковесную блокировку. Блокировки — это механизмы сериализации, подавляющие параллелизм. Приложения, которые подавляют параллелизм, являются менее масштабируемыми, могут поддерживать меньшее число пользователей и требуют больший объем ресурсов. Наша цель всегда заключается в построении приложений, обладающих потенциалом масштабирования — приложений, которые могут обслуживать как одного, так и 1000 или 10 000 пользователей. Чем меньше происходит зашелкивания, тем лучше. Я могу выбрать подход, который приводит к увеличению времени выполнения, но снижает количество зашелок до 10% по отношению к другому подходу. Мне известно, что подход, использующий меньше зашелок, будет масштабироваться значительно лучше, чем подход с более высоким числом зашелок.

Утилиту Runstats лучше всего применять в изоляции, т.е. в однопользовательской базе данных. Мы будем измерять статистические показатели и активность по зашелкиванию (блокированию) данных, которые присущи анализируемым подходам. Мы не хотим, чтобы во время этих измерений другие сеансы вносили свой вклад в загрузку системы или зашелкивание. Для проведения тестов подобного рода удобно иметь небольшую тестовую базу данных. Например, я часто использую свой настольный ПК или ноутбук.

На заметку! Я уверен, что все разработчики должны располагать полностью контролируемой испытательной базой данных, в которой они смогут проверять свои идеи, не обращаясь постоянно за помощью к администратору базы данных. Разработчики определенно должны иметь такую базу данных на своих настольных ПК, даже если лицензия в персональной версии для разработчика предполагает просто “применение ее для разработки и тестирования, но не для развертывания”. В этом случае нечего терять! Кроме того, я проводил несколько неформальных опросов на конференциях и семинарах. Практически каждый администратор баз данных начинал свою деятельность как разработчик! Опыт и обучение, которое разработчики могут получить при наличии собственной базы данных — имея возможность видеть, как она в действительности работает — в долгосрочной перспективе сулит значительные преимущества.

Чтобы использовать Runstats, понадобится настроить доступ к нескольким представлениям V\$, создать таблицу для хранения статистики и создать пакет Runstats. Потребуется доступ к четырем таблицам V\$ (динамическим таблицам производительности): V\$STATNAME, V\$MYSTAT, V\$TIMER и V\$LATCH. Вот представление, с которым я работаю:

```
create or replace view stats
as select 'STAT...' || a.name name, b.value
   from v$statname a, v$mystat b
  where a.statistic# = b.statistic#
 union all
 select 'LATCH.' || name, gets
   from v$latch
 union all
 select 'STAT...Elapsed Time', hsecs from v$timer;
```

На заметку! Действительными именами объектов, к которым должен быть предоставлен доступ, являются V_\$STATNAME, V_\$MYSTAT и т.д. — имена объектов начинаются с V_\$, а не V\$. Имя V\$ — это синоним, который указывает на представление с именем, начинающимся с V_\$. Таким образом, V\$STATNAME — синоним, указывающий на представление V_\$STATNAME. У вас должен быть открыт доступ к этому представлению.

Вы можете иметь привилегию SELECT на V\$STATNAME, V\$MYSTAT и V\$LATCH, выданную непосредственно вам, так что у вас будет возможность создать представление самостоятельно. Кроме того, кто-то другой, располагающий привилегией SELECT на этих объектах, может создать представление и затем выдать вам привилегию SELECT на нем. Осталось создать небольшую таблицу для сбора статистики:

```
create global temporary table run_stats
( runid varchar2(15),
  name varchar2(80),
  value int )
on commit preserve rows;
```

Наконец, необходимо создать сам пакет для Runstats. Он содержит три простых API-вызова:

- RS_START (запуск Runstats) — процедура, которая должна вызываться в начале теста Runstats;
- RS_MIDDLE — процедура, которая должна вызываться в середине теста;
- RS_STOP — процедура, которая завершает выполнение и выводит отчет.

Спецификация этого пакета выглядит следующим образом:

```
EODA@ORA12CR1> create or replace package runstats_pkg
2 as
3     procedure rs_start;
4     procedure rs_middle;
5     procedure rs_stop( p_difference_threshold in number default 0 );
6 end;
7 /
```

Package created.

Пакет создан.

Параметр `p_difference_threshold` служит для управления объемом выводимых в конце данных. Пакет Runstats собирает статистические сведения и информацию о зашелках для каждого запуска, после чего выводит отчет о количестве используемых каждым тестом (каждым подходом) ресурсов и разнице между ними. Этот входной параметр можно применять для вывода только тех статистических сведений и информации о зашелках, которые отличаются больше, чем на это число. По умолчанию значение `p_difference_threshold` равно нулю, поэтому отображаются все результаты.

Теперь рассмотрим тело пакета, процедуру за процедурой. Пакет начинается с объявления ряда глобальных переменных. Они будут использоваться для хранения времени выполнения каждого теста:

```
EODA@ORA12CR1> create or replace package body runstats_pkg
2 as
3
4     g_start number;
5     g_run1 number;
6     g_run2 number;
7
```

Затем следует процедура RS_START. Она просто очищает таблицу со статистическими данными и заполняет ее статистическими сведениями “до запуска” и информацией о количестве зашелок. После этого она захватывает текущее значение таймера, которое можно применять для вычисления пройденного времени с точностью до сотых долей секунды:

```
8 procedure rs_start
9 is
10 begin
11     delete from run_stats;
12
13     insert into run_stats
14     select 'before', stats.* from stats;
```

```

15
16   g_start := dbms_utility.get_cpu_time;
17 end;
18

```

Следующая процедура — RS_MIDDLE. Она помещает в переменную G_RUN1 значение времени, затраченного на выполнение первого теста. Затем она вставляет текущий набор статистических данных и информацию о количестве защелок. Если вычесть эти значения из значений, которые ранее были сохранены во время выполнения процедуры RS_START, можно выяснить количество защелок, задействованных первым методом, количество использованных курсоров (статистический показатель) и т.д.

И, наконец, процедура записывает время начала следующего теста:

```

19 procedure rs_middle
20 is
21 begin
22   g_run1 := (dbms_utility.get_cpu_time-g_start);
23
24   insert into run_stats
25   select 'after 1', stats.* from stats;
26
27   g_start := dbms_utility.get_cpu_time;
28 end;
29

```

Последней процедурой в этом пакете является RS_STOP. Ее задача состоит в выводе совокупного времени процессора для каждого запуска, а также разницы между статистическими показателями по каждому из двух запусков (выводятся только те, что превышают пороговое значение):

```

30 procedure rs_stop(p_difference_threshold in number default 0)
31 is
32 begin
33   g_run2 := (dbms_utility.get_cpu_time-g_start);
34
35   dbms_output.put_line( 'Run1 ran in ' || g_run1 || ' cpu hsecs' );
36   dbms_output.put_line( 'Run2 ran in ' || g_run2 || ' cpu hsecs' );
37
38   if ( g_run2 <> 0 )
39   then
40     dbms_output.put_line
41     ( 'run 1 ran in ' || round(g_run1/g_run2*100,2) ||
42       '% of the time' );
43   end if;
44   dbms_output.put_line( chr(9) );
45
46   insert into run_stats
47   select 'after 2', stats.* from stats;
48
49   dbms_output.put_line
50   ( rpad( 'Name', 30 ) || lpad( 'Run1', 16 ) ||
51     lpad( 'Run2', 16 ) || lpad( 'Diff', 16 ) );
52

```

```

53  for x in
54  ( select rpad( a.name, 30 ) ||
55    to_char( b.value-a.value, '999,999,999,999' ) ||
56    to_char( c.value-b.value, '999,999,999,999' ) ||
57    to_char( ( c.value-b.value)-(b.value-a.value)),
58    '999,999,999,999' ) data
59    from run_stats a, run_stats b, run_stats c
60    where a.name = b.name
61    and b.name = c.name
62    and a.runid = 'before'
63    and b.runid = 'after 1'
64    and c.runid = 'after 2'
65
66    and abs( (c.value-b.value) - (b.value-a.value) )
67    > p_difference_threshold
68    order by abs( (c.value-b.value)-(b.value-a.value))
69  ) loop
70  dbms_output.put_line( x.data );
71  end loop;
72
73  dbms_output.put_line( chr(9) );
74  dbms_output.put_line
75  ( 'Run1 latches total versus runs -- difference and pct' );
76  dbms_output.put_line
77  ( lpad( 'Run1', 14 ) || lpad( 'Run2', 19 ) ||
78    lpad( 'Diff', 18 ) || lpad( 'Pct', 11 ) );
79
80  for x in
81  ( select to_char( run1, '9,999,999,999,999' ) ||
82    to_char( run2, '9,999,999,999,999' ) ||
83    to_char( diff, '9,999,999,999,999' ) ||
84    to_char( round( run1/decode( run2, 0, to_number(0), run2) *100,2 ),
'99,999.99' ) || '%' data
85    from ( select sum(b.value-a.value) run1, sum(c.value-b.value) run2,
86    sum( (c.value-b.value)-(b.value-a.value)) diff
87    from run_stats a, run_stats b, run_stats c
88    where a.name = b.name
89    and b.name = c.name
90    and a.runid = 'before'
91    and b.runid = 'after 1'
92    and c.runid = 'after 2'
93    and a.name like 'LATCH%'
94    )
95  ) loop
96  dbms_output.put_line( x.data );
97  end loop;
98  end;
99
100 end;
101 /

```

Package body created.

Тело пакета создано.

Итак, все готово для применения Runstats. В качестве примера мы покажем, как использовать Runstats, чтобы выяснить, какой метод эффективнее: одна групповая операция INSERT или построчная обработка. Начнем с определения двух таблиц, в которые будут вставляться 1 000 000 строк (сценарий создания таблицы BIG_TABLE представлен далее в этом разделе):

```
EODA@ORA12CR1> create table t1
  2 as
  3 select * from big_table
  4 where l=0;
Table created.
Таблица создана.
```

```
EODA@ORA12CR1> create table t2
  2 as
  3 select * from big_table
  4 where l=0;
Table created.
Таблица создана.
```

Теперь выполним первый метод вставки записей, предусматривающий применение одного SQL-оператора. Начнем с вызова процедуры RUNSTATS_PKG.RS_START:

```
EODA@ORA12CR1> exec runstats_pkg.rs_start;
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.
EODA@ORA12CR1> insert into t1
  2 select *
  3   from big_table
  4   where rownum <= 1000000;
1000000 rows created.
1000000 строк создано.
EODA@ORA12CR1> commit;
Commit complete.
Фиксация завершена.
```

Затем можно выполнить второй метод — построчную вставку данных:

```
EODA@ORA12CR1> exec runstats_pkg.rs_middle;
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.
EODA@ORA12CR1> begin
  2   for x in ( select *
  3               from big_table
  4               where rownum <= 1000000 )
  5   loop
  6       insert into t2 values X;
  7   end loop;
  8   commit;
  9 end;
10 /
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.
```

И, наконец, сгенерируем отчет:

```
EODA@ORA12CR1> exec runstats_pkg.rs_stop(1000000)
Run1 ran in 119 cpu hsecs
Run2 ran in 3376 cpu hsecs
run 1 ran in 3.52% of the time
```

| Name | Run1 | Run2 | Diff |
|--------------------------------|---------------|---------------|---------------|
| STAT...execute count | 29 | 1,000,032 | 1,000,003 |
| STAT...opened cursors cumulati | 29 | 1,000,035 | 1,000,006 |
| LATCH.shared pool | 582 | 1,001,466 | 1,000,884 |
| STAT...session logical reads | 148,818 | 1,158,009 | 1,009,191 |
| STAT...recursive calls | 183 | 1,010,218 | 1,010,035 |
| STAT...db block changes | 95,964 | 2,074,283 | 1,978,319 |
| LATCH.cache buffers chains | 443,882 | 5,462,356 | 5,018,474 |
| STAT...undo change vector size | 3,620,400 | 67,938,496 | 64,318,096 |
| STAT...KTFB alloc space (block | 109,051,904 | 176,160,768 | 67,108,864 |
| STAT...redo size | 105,698,540 | 384,717,388 | 279,018,848 |
| STAT...logical read bytes from | 1,114,251,264 | 9,300,803,584 | 8,186,552,320 |

```
Run1 latches total versus runs -- difference and pct
Run1          Run2          Diff          Pct
555,593       6,795,317       6,239,724       8.18%
```

PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.

Отчет подтверждает, что пакет RUNSTATS_PKG установлен, и показывает, почему при разработке приложений вы должны использовать одиночный оператор SQL вместо порции процедурного кода везде, где только возможно.

Mystat

Сценарий mystat.sql и сопровождающий его mystat2.sql применяются для того, чтобы показать, как изменяются некоторые “статистические показатели” Oracle до и после выполнения определенной операции. Сценарий mystat.sql собирает начальные значения показателей:

```
set echo off
set verify off
column value new_val V
define S="&1"

set autotrace off
select a.name, b.value
from v$statname a, v$mystat b
where a.statistic# = b.statistic#
and lower(a.name) = lower('&S')
/
set echo on
```

Сценарий mystat2.sql отображает разницу (&V заполняется в результате запуска сценария mystat.sql — он использует для этого средство NEW_VAL из SQL*Plus, которое содержит последнее значение VALUE, выбранное из предшествующего запроса):

```

set echo off
set verify off
select a.name, b.value V, to_char(b.value-&V, '999,999,999,999') diff
from v$statname a, v$mystat b
where a.statistic# = b.statistic#
and lower(a.name) = lower('&S')
/
set echo on

```

Например, чтобы посмотреть, сколько информации redo сгенерировано оператором UPDATE, можно поступить следующим образом:

```

EODA@ORA12CR1> @mystat "redo size"
EODA@ORA12CR1> set echo off

NAME                                VALUE
-----
redo size                          491167892

EODA@ORA12CR1> update big_table set owner = lower(owner)
2  where rownum <= 1000;

1000 rows updated.
1000 строк обновлено.

EODA@ORA12CR1> @mystat2
EODA@ORA12CR1> set echo off

```

| NAME | V | DIFF |
|-----------|-----------|--------|
| redo size | 491265640 | 97,748 |

Здесь видно, что оператор UPDATE для 1000 строк сгенерировал 97 748 байтов данных redo.

SHOW_SPACE

Процедура SHOW_SPACE выводит подробную информацию об утилизации пространства для сегментов базы данных. Ее интерфейс выглядит так:

```

EODA@ORA12CR1> desc show_space
PROCEDURE show_space
Argument Name          Type          In/Out  Default?
-----
P_SEGNAME              VARCHAR2      IN
P_OWNER                VARCHAR2      IN      DEFAULT
P_TYPE                 VARCHAR2      IN      DEFAULT
P_PARTITION            VARCHAR2      IN      DEFAULT

```

Ниже описаны аргументы этой процедуры.

- P_SEGNAME. Имя сегмента — например, имя таблицы или индекса.
- P_OWNER. По умолчанию принимается текущий пользователь, но эту процедуру можно применять для просмотра какой-то другой схемы.
- P_TYPE. По умолчанию этот аргумент принимает значение TABLE и представляет тип просматриваемого объекта. Например, оператор `select distinct segment_type from dba_segments` выводит список допустимых типов сегментов.

- **P_PARTITION.** Имя секции при просмотре пространства для секционированного объекта. Процедура `SHOW_SPACE` одновременно может отображать пространство только для одной секции.

Когда сегмент хранится в табличном пространстве **ASSM**, вывод из этой процедуры выглядит следующим образом:

```
EODA@ORA12CR1> exec show_space('BIG_TABLE');
Unformatted Blocks ..... 0
FS1 Blocks (0-25) ..... 0
FS2 Blocks (25-50) ..... 0
FS3 Blocks (50-75) ..... 0
FS4 Blocks (75-100) ..... 0
Full Blocks ..... 14,469
Total Blocks..... 15,360
Total Bytes..... 125,829,120
Total MBytes..... 120
Unused Blocks..... 728
Unused Bytes..... 5,963,776
Last Used Ext FileId..... 4
Last Used Ext BlockId..... 43,145
Last Used Block..... 296
```

PL/SQL procedure successfully completed.

Процедура PL/SQL успешно завершена.

Ниже приведены пояснения элементов вывода.

- **Unformatted Blocks (Неформатированных блоков).** Количество блоков, которые выделены для таблицы и хранятся ниже маркера максимального уровня заполнения (**high-water mark — HWM**), но не используются. Сумма неформатированных и неиспользуемых блоков равна общему количеству блоков, выделенных для таблицы, но не применяемых для хранения данных в объекте **ASSM**.
- **FS1 Blocks – FS4 Blocks (Блоков FS1 — Блоков FS4).** Количество форматированных блоков с данными. Диапазоны, указанные после их названий, представляют степень пустоты каждого блока. Например, диапазон (0–25) представляет количество блоков, пустых на 0–25%.
- **Full Blocks (Полных блоков).** Количество блоков, заполненных настолько, что они больше не являются кандидатами для выполнения последующих вставок данных.
- **Total Blocks (Всего блоков), Total Bytes (Всего байтов), Total Mbytes (Всего Мбайт).** Общий объем пространства, которое выделено сегменту, измеренное в блоках базы данных, байтах и мегабайтах.
- **Unused Blocks (Неиспользуемых блоков), Unused Bytes (Неиспользуемых байтов).** Эти значения представляют объем дискового пространства, которое никогда не использовалось. Эти блоки выделены для сегмента, но в настоящее время хранятся выше **HWM**-маркера сегмента.
- **Last Used Ext FileId (Идентификатор файла последнего использованного экстенда).** Идентификатор файла, который содержит последний экстенд с данными.

- Last Used Ext BlockId (Идентификатор блока последнего использованного экстенста). Идентификатор блока начала последнего экстенста; идентификатор блока внутри последнего использованного файла.
- Last Used Block (Последний использованный блок). Смещение идентификатора последнего блока, использованного в последнем экстенсте.

Когда процедура SHOW_SPACE применяется для просмотра объектов в табличных пространствах с ручным управлением пространством сегментов, вывод напоминает показанный далее:

```

EODA@ORA12CR1> exec show_space( 'BIG_TABLE' )
Free Blocks..... 1
Total Blocks..... 147,456
Total Bytes..... 1,207,959,552
Total MBytes..... 1,152
Unused Blocks..... 1,616
Unused Bytes..... 13,238,272
Last Used Ext FileId..... 7
Last Used Ext BlockId..... 139,273
Last Used Block..... 6,576

PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.

```

Единственное отличие от предыдущего отчета связано со строкой Free Blocks (Свободных блоков) в начале. Здесь указано количество блоков в первой группе списков свободных блоков сегмента. Сценарий выводит сведения только об этой группе списков свободных блоков. Чтобы сценарий мог отображать сведения о нескольких группах списков свободных блоков, его понадобится модифицировать.

Ниже приведен код с комментариями. Эта утилита представляет собой дополнительный уровень, действующий поверх API-интерфейса DBMS_SPACE в базе данных.

```

create or replace procedure show_space
( p_segname in varchar2,
  p_owner   in varchar2 default user,
  p_type    in varchar2 default 'TABLE',
  p_partition in varchar2 default NULL )
-- Данная процедура использует идентификатор authid текущего пользователя,
-- поэтому она может запрашивать представления DBA_* с использованием
-- привилегий из ROLE и допускает однократную установку для всей базы данных,
-- а не для каждого пользователя, который желает ее применять.
authid current_user
as
  l_free_blks          number;
  l_total_blocks       number;
  l_total_bytes        number;
  l_unused_blocks      number;
  l_unused_bytes       number;
  l_LastUsedExtFileId  number;
  l_LastUsedExtBlockId number;
  l_LAST_USED_BLOCK    number;
  l_segment_space_mgmt varchar2(255);
  l_unformatted_blocks number;
  l_unformatted_bytes  number;

```

```

l_fs1_blocks number; l_fs1_bytes number;
l_fs2_blocks number; l_fs2_bytes number;
l_fs3_blocks number; l_fs3_bytes number;
l_fs4_blocks number; l_fs4_bytes number;
l_full_blocks number; l_full_bytes number;

-- Встроенная процедура для вывода аккуратно
-- форматированных чисел с простыми метками.
procedure p( p_label in varchar2, p_num in number )
is
begin
    dbms_output.put_line( rpad(p_label,40,'.') ||
                          to_char(p_num,'999,999,999,999') );
end;
begin
-- Этот запрос выполняется динамически, чтобы позволить данной процедуре
-- быть созданной пользователем, имеющим доступ к представлениям
-- DBA_SEGMENTS/TABLESPACES посредством роли, как обычно принято.
-- ПРИМЕЧАНИЕ: во время выполнения вызывающая процедура
-- ДОЛЖНА иметь доступ к этим двум представлениям!
-- Этот запрос определяет, является ли данный объект объектом ASSM.
begin
    execute immediate
        'select ts.segment_space_management
          from dba_segments seg, dba_tablespaces ts
         where seg.segment_name      = :p_segname
           and (:p_partition is null or
                seg.partition_name = :p_partition)
           and seg.owner = :p_owner
           and seg.tablespace_name = ts.tablespace_name'
        into l_segment_space_mgmt
        using p_segname, p_partition, p_partition, p_owner;
exception
    when too_many_rows then
        dbms_output.put_line
            ( 'This must be a partitioned table, use p_partition => ');
-- Эта таблица должна быть секционированной, используйте p_partition =>
    return;
end;

-- Если объект является табличным пространством ASSM, то для получения
-- информации о пространстве мы должны применять этот API-вызов.
-- В противном случае мы используем API-интерфейс FREE_BLOCKS
-- для сегментов, управляемых пользователем.
if l_segment_space_mgmt = 'AUTO'
then
    dbms_space.space_usage
        ( p_owner, p_segname, p_type, l_unformatted_blocks,
          l_unformatted_bytes, l_fs1_blocks, l_fs1_bytes,
          l_fs2_blocks, l_fs2_bytes, l_fs3_blocks, l_fs3_bytes,
          l_fs4_blocks, l_fs4_bytes, l_full_blocks, l_full_bytes, p_partition);
    p( 'Unformatted Blocks ', l_unformatted_blocks );
    p( 'FS1 Blocks (0-25) ', l_fs1_blocks );

```

```

p( 'FS2 Blocks (25-50) ', l_fs2_blocks );
p( 'FS3 Blocks (50-75) ', l_fs3_blocks );
p( 'FS4 Blocks (75-100)', l_fs4_blocks );
p( 'Full Blocks      ', l_full_blocks );

```

```

else

```

```

  dbms_space.free_blocks(
    segment_owner   => p_owner,
    segment_name    => p_segname,
    segment_type    => p_type,
    freelist_group_id => 0,
    free_blks       => l_free_blks);

```

```

  p( 'Free Blocks', l_free_blks );

```

```

end if;

```

-- Затем с помощью API-вызова unused_space извлекается остальная информация.

```

dbms_space.unused_space
( segment_owner   => p_owner,
  segment_name    => p_segname,
  segment_type    => p_type,
  partition_name  => p_partition,
  total_blocks    => l_total_blocks,
  total_bytes     => l_total_bytes,
  unused_blocks   => l_unused_blocks,
  unused_bytes    => l_unused_bytes,
  LAST_USED_EXTENT_FILE_ID => l_LastUsedExtFileId,
  LAST_USED_EXTENT_BLOCK_ID => l_LastUsedExtBlockId,
  LAST_USED_BLOCK => l_LAST_USED_BLOCK );
p( 'Total Blocks', l_total_blocks );
p( 'Total Bytes', l_total_bytes );
p( 'Total MBytes', trunc(l_total_bytes/1024/1024) );
p( 'Unused Blocks', l_unused_blocks );
p( 'Unused Bytes', l_unused_bytes );
p( 'Last Used Ext FileId', l_LastUsedExtFileId );
p( 'Last Used Ext BlockId', l_LastUsedExtBlockId );
p( 'Last Used Block', l_LAST_USED_BLOCK );

```

```

end;
/

```

BIG_TABLE

Во многих примерах этой книги используется таблица BIG_TABLE. В зависимости от системы эта таблица может содержать от одной до 4 миллионов записей, а ее размер варьируется в пределах от 200 Мбайт до 800 Мбайт. Во всех случаях ее структура остается той же самой.

Для создания таблицы BIG_TABLE я написал сценарий, который выполняет перечисленные ниже действия.

- Создает пустую таблицу на основе представления ALL_OBJECTS. Это словарное представление применяется для заполнения таблицы BIG_TABLE.
- Устанавливает для этой таблицы режим NOLOGGING. Это не обязательно и было сделано ради производительности. Использовать режим NOLOGGING для тестовой таблицы совершенно безопасно. Однако не применяйте его в производственной системе, т.к. тогда будут отключены средства вроде Oracle Data Guard (Защита данных Oracle).

- Заполняет таблицу, занося в нее содержимое представления ALL_OBJECTS и затем выполняя на каждой итерации вставку таблицы в саму себя, что каждый раз увеличивает ее размер примерно вдвое.
- Создает на таблице ограничение первичного ключа.
- Собирает статистику.

Чтобы построить таблицу BIG_TABLE, можете запустить следующий сценарий в командной строке SQL*Plus и передать ему количество строк таблицы в качестве параметра. Сценарий остановится, когда достигнет указанного количества строк.

```
create table big_table
as
select rownum id, OWNER, OBJECT_NAME, SUBOBJECT_NAME, OBJECT_ID,
DATA_OBJECT_ID, OBJECT_TYPE, CREATED, LAST_DDL_TIME, TIMESTAMP,
STATUS, TEMPORARY, GENERATED, SECONDARY, NAMESPACE, EDITION_NAME
  from all_objects
 where 1=0
/

alter table big_table nologging;

declare
  l_cnt number;
  l_rows number := &numrows;
begin
  insert /*+ append */
  into big_table
  select rownum id, OWNER, OBJECT_NAME, SUBOBJECT_NAME, OBJECT_ID,
DATA_OBJECT_ID, OBJECT_TYPE, CREATED, LAST_DDL_TIME, TIMESTAMP,
STATUS, TEMPORARY, GENERATED, SECONDARY, NAMESPACE, EDITION_NAME
  from all_objects
 where rownum <= &numrows;
  --
  l_cnt := sql%rowcount;
  commit;
  while (l_cnt < l_rows)
  loop
    insert /*+ APPEND */ into big_table
    select rownum+l_cnt, OWNER, OBJECT_NAME, SUBOBJECT_NAME, OBJECT_ID,
DATA_OBJECT_ID, OBJECT_TYPE, CREATED, LAST_DDL_TIME, TIMESTAMP,
STATUS, TEMPORARY, GENERATED, SECONDARY, NAMESPACE, EDITION_NAME
    from big_table a
    where rownum <= l_rows-l_cnt;
    l_cnt := l_cnt + sql%rowcount;
    commit;
  end loop;
end;
/

alter table big_table add constraint
big_table_pk primary key(id);

exec dbms_stats.gather_table_stats( user, 'BIG_TABLE', estimate_percent=> 1);
```


Я произвел оценку базовой статистики о таблице. Для индекса, ассоциированного с первичным ключом, статистика вычисляется автоматически при его создании.

Соглашения при написании кода

Одно из используемых в этой книге соглашений при написании кода касается именования переменных в коде PL/SQL. Например, взгляните на следующее тело пакета:

```
create or replace package body my_pkg
as
    g_variable varchar2(25);
    procedure p( p_variable in varchar2 )
    is
        l_variable varchar2(25);
    begin
        null;
    end;
end;
/
```

Здесь присутствуют три переменные: глобальная переменная пакета `G_VARIABLE`, формальный параметр процедуры `P_VARIABLE` и локальная переменная `L_VARIABLE`. Переменные именуются согласно их области действия. Имена глобальных переменных начинаются с `G_`, имена параметров — с `P_`, а имена локальных переменных — с `L_`. Основная причина такого именования — провести различие между именами переменных PL/SQL и именами столбцов в таблице базы данных. Например, показанная ниже процедура будет всегда выводить все строки таблицы EMP, в которых поле ENAME не является пустым:

```
create procedure p( ENAME in varchar2 )
as
begin
    for x in ( select * from emp where ename = ENAME ) loop
        Dbms_output.put_line( x.empno );
    end loop;
end;
```

Встретив выражение `ename = ENAME`, механизм SQL сравнивает столбец ENAME с самим собой (разумеется). Можно было бы применить выражение `ename = P.ENAME`, т.е. уточнить ссылку на переменную PL/SQL с помощью имени процедуры, но об этом легко забыть и получить в итоге ошибку.

Я просто всегда именую свои переменные по области действия. В результате я легко могу отличать параметры от локальных и глобальных переменных, а также устраняю любую неоднозначность в отношении имен столбцов и имен переменных.

ГЛАВА 1

Разработка успешных приложений Oracle

Очень много времени я провожу, работая с программным обеспечением системы управления базами данных (СУБД) Oracle и, что еще важнее, общаясь с теми, кто использует это программное обеспечение. В течение последних 25 лет или около того мне довелось работать над многим проектами — как успешными, так и совершенно неудачными — и мой опыт можно было бы обобщить следующим образом.

- На успех или провал приложения, построенного на основе базы данных — и зависящего от базы данных — оказывает влияние то, как оно использует базу данных. Как следствие этого, все приложения построены на основе баз данных; я не в состоянии назвать ни одного полезного приложения, которое где-то не хранило бы данные.
- Приложения приходят и уходят, однако *данные* существуют всегда. Речь идет не о построении приложений; на самом деле речь идет о данных, лежащих в основе этих приложений.
- Команда разработки нуждается в нескольких разработчиках, разбирающихся в базах данных, которые будут отвечать за соблюдение логики базы данных и обеспечение работоспособности построенной системы с первого дня ее введения в эксплуатацию. Необходимость в проведении подстройки сразу после развертывания свидетельствует о том, что система не создавалась в подобной манере.

Эти соображения могут казаться совершенно очевидными, но я убедился, что слишком многие относятся к базе данных как к черному ящику — чему-то такому, о чем знать не обязательно. Возможно, они располагают генератором SQL-кода, который, по их мнению, избавит от бремени изучения языка SQL. Или же они полагают, что будут пользоваться базой данных только как двумерным файлом и выполнять “чтение по ключам”. Как бы там ни было, я могу утверждать, что в большинстве случаев такая точка зрения ошибочна: без понимания базы данных просто не обойтись. В этой главе будет показано, *почему* требуется знание базы данных, в частности, зачем необходимо разбираться в перечисленных ниже вопросах.

- Архитектура базы данных: как она работает и что собой представляет.
- Что такое средства управления параллельной обработкой, и каков в них смысл.

- Как выполнять повседневную настройку своего приложения.
- Каким образом в базе данных реализованы некоторые вещи — они не обязательно реализованы так, как вы могли подумать.
- Какие функциональные возможности база данных уже предлагает, и почему в общем случае лучше применять предоставленное средство, чем строить собственное.
- Почему может требоваться знание языка SQL, выходящее за рамки общего курса.
- Администраторы баз данных и разработчики находятся в одной команде, а не представляют два враждующих лагеря, которые постоянно пытаются перехитрить друг друга.

Поначалу может показаться, что перечень подлежащих изучению вопросов бесконечен, но примите во внимание следующую аналогию: что бы вы сделали в первую очередь, если бы пришлось разрабатывать в высшей степени масштабируемое производственное приложение для совершенно новой операционной системы (ОС)? Скорее всего, ваш ответ был бы следующим: “Сначала выяснили бы, как работает эта новая ОС, каким образом в ней будут выполняться приложения, и тому подобные вопросы”. Если ваш ответ отличается, то с большой вероятностью вы потерпите поражение.

Например, сравним ОС Windows и UNIX/Linux. Если вы являетесь опытным программистом в Windows, которого попросили разработать новое приложение для платформы UNIX/Linux, то вам придется заново изучить ряд вопросов. Управление памятью реализовано совершенно по-другому. Построение процессов сервера существенно отличается: в среде Windows должен быть разработан единственный процесс — единственный исполняемый модуль с множеством потоков. В среде UNIX/Linux разработка не сводится к построению одного автономного исполняемого модуля; вместо этого приходится иметь дело с множеством совместно работающих процессов. Это верно, что и Windows, и UNIX/Linux являются операционными системами. Обе они предлагают разработчикам множество одинаковых служб — управление файлами, управление памятью, управление процессами, обеспечение безопасности и т.п. Однако они серьезно отличаются по архитектуре, и большая часть того, что вы изучили в среде Windows, не применимо к UNIX/Linux (и, справедливости ради, наоборот). Чтобы добиться успеха на новой платформе, вы должны забыть все, что знали о старой платформе. То же самое касается и среды базы данных.

Все моменты, связанные с приложениями, которые выполняются непосредственно в среде ОС, относятся также и к приложениям, которые будут выполняться в среде базы данных: необходимо осознавать, что база данных является критически важной для успеха приложения. Если вы не понимаете, что делает конкретная база данных и как она это делает, то ваше приложение будет обречено на неудачу. Если вы считаете, что раз уж ваше приложение нормально работает в среде SQL Server, то оно обязательно будет успешно функционировать в среде Oracle, то это приложение, скорее всего, потерпит неудачу. Разумеется, корректно также и обратное утверждение: вовсе не обязательно, что масштабируемое, аккуратно разработанное приложение Oracle сможет работать в среде SQL Server без существенных архитектурных изменений. Подобно тому, как Windows и UNIX/Linux являются операцион-

ными системами, но фундаментально отличаются друг от друга, Oracle и SQL Server (здесь можно было бы упомянуть практически любую базу данных) представляя собой базы данных, существенно отличающиеся в архитектурном плане.

Мой подход

Прежде чем мы приступим, я должен пояснить вам свой подход к разработке. Я предпочитаю подходить к решению проблем, строя решение вокруг базы данных. Если какое-то действие может быть выполнено в базе данных, то так и делается. Это обусловлено несколькими причинами. Самая главная из них заключается в том, что если функциональность встроена в базу данных, то ее можно *развернуть* где угодно. Мне не известна ни одна популярная и коммерчески успешная серверная операционная система, для которой программное обеспечение Oracle было бы недоступным. Начиная с Windows и заканчивая десятками систем UNIX/Linux — везде доступно одно и то же программное обеспечение Oracle и те же самые варианты. Мне часто приходится строить и тестировать решения Oracle 12c, Oracle 11g или Oracle 10g на портативном компьютере под управлением ОС UNIX/Linux либо ОС Windows, которая выполняется на виртуальной машине. Полученные решения можно затем развертывать на разнообразных серверах, на которых функционирует одно и то же программное обеспечение баз данных, но в средах разных операционных систем. Если же функциональное средство приходится реализовать вне базы данных, то окажется, что его крайне трудно развертывать везде, где хочется. Одной из основных характеристик языка Java, которая делает его настолько привлекательным, является компиляция программ в одной и той же среде виртуальной машины Java (Java Virtual Machine — JVM), обеспечивающая высокую степень их переносимости; именно эта характеристика привлекает меня в базах данных. База данных — это *моя* виртуальная машина и это *моя* операционная система.

Итак, я стараюсь все выполнять в базе данных. Если же мои требования выходят за рамки того, что может предложить среда базы данных, я строю решение на языке Java вне этой среды. В результате почти все сложности, связанные с операционной системой, остаются скрытыми. Я по-прежнему должен понимать, каким образом работают *мои* “виртуальные машины” (Oracle и временами JVM) — ведь всегда нужно знать инструменты, с которыми приходится работать, — но они, в свою очередь, позаботятся о том, как лучше выполнять те или иные действия в среде заданной ОС.

Таким образом, простое знание особенностей работы одной “виртуальной ОС” позволяет строить приложения, которые будут успешно функционировать и масштабироваться под управлением многих операционных систем. Я вовсе не утверждаю, что можно находиться в полном неведении относительно используемой ОС, но как разработчик программного обеспечения, строящий приложения баз данных, вы довольно неплохо изолированы от нее и не будете иметь дела с большинством нюансов. Безусловно, администратору базы данных, отвечающему за работу программного обеспечения Oracle, придется намного больше взаимодействовать с ОС (если это не так, то поищите другого администратора). Если же вы разрабатываете клиент-серверное программное обеспечение, и большая часть кода выполняется за пределами базы данных и вне виртуальной машины (вероятно, наиболее популярной виртуальной машиной является JVM), то придется снова обращать внимание на ОС.

При разработке программного обеспечения базы данных я придерживаюсь достаточно простой философии, которая оставалась неизменной на протяжении многих лет.

- Все, что только возможно, должно делаться в одном операторе SQL. Верите или нет, но это возможно почти всегда. С течением времени это утверждение становится еще более справедливым. SQL — исключительно мощный язык.
- Если что-то нельзя выполнить в одном операторе SQL, то это необходимо реализовать на языке PL/SQL с помощью как можно более краткого кода. Следуйте принципу “больше кода = больше ошибок, меньше кода = меньше ошибок”.
- Если задачу нельзя решить средствами PL/SQL, попробуйте воспользоваться хранимой процедурой Java. Однако после выхода Oracle9i и последующих версий потребность в этом возникает очень редко. PL/SQL является полноценным и популярным языком третьего поколения (third-generation programming language — 3GL).
- Если задачу не удастся решить на языке Java, попробуйте написать внешнюю процедуру С. Именно такой подход применяют наиболее часто, когда нужно обеспечить высокую скорость работы приложения либо использовать API-интерфейс от независимых разработчиков, реализованный на языке С.
- Если вы не можете решить задачу с помощью внешней процедуры С, всерьез задумайтесь над тем, если в ней необходимость.

В этой книге вы будете повсеместно встречать примеры воплощения описанной выше философии. Мы будем использовать PL/SQL и его объектные типы для решения задач, которые в SQL сделать невозможно или неэффективно. Язык PL/SQL существует уже очень долгое время — на его отработку ушло более 27 лет (к 2015 году); в действительности, возвращаясь к версии Oracle 10g, был переписан сам компилятор PL/SQL, чтобы стать в первую очередь оптимизирующим компилятором. Никакой другой язык не связан настолько тесно с SQL и не является до такой степени оптимизированным для взаимодействия с SQL. Работа с SQL в среде PL/SQL происходит совершенно естественным образом, в то время как в любом другом языке, от Visual Basic до Java, применение SQL может оказаться довольно-таки обременительным. PL/SQL никогда не выглядел достаточно “естественным”, поскольку он не является расширением какого-то языка. Когда возможности PL/SQL исчерпываются, что в текущих выпусках базы данных случается очень редко, мы будем использовать Java. Иногда мы будем прибегать к языку С, но, как правило, только когда С является единственным выбором или когда требуется максимальная скорость, обеспечиваемая С. Часто последняя причина уходит на второй план благодаря низкоуровневой компиляции Java — возможности преобразования байт-кода Java в специфический для операционной системы объектный код платформы. Во многих случаях это позволяет коду Java работать почти так же быстро, как код С.

Метод черного ящика

Личный практический опыт (имеются в виду допущенные мною ошибки) позволил мне составить собственное мнение о причинах столь частых неудач при раз-

работке программного обеспечения, взаимодействующего с базами данных. Должен признаться, что я включил в эту книгу ряд проектов, которые хотя и не были официально признаны неудачными, однако потребовали значительно большего времени на разработку и развертывание, чем планировалось изначально. Причина в том, что в них приходилось вносить существенные изменения, изменять архитектуру или прилагать значительные усилия по настройке. Я называю такие затянувшиеся проекты неудачами, поскольку в большинстве случаев они должны были быть завершены в плановые сроки (или даже быстрее).

Единственная наиболее часто встречающаяся причина неудачи связана с недостаточными практическими знаниями самой базы данных, т.е. отсутствие общего представления об основном инструменте, который применяется. Метод черного ящика предполагает намеренное абстрагирование разработчиков от базы данных; фактически он поощряет отсутствие вообще каких-либо знаний о базе данных! Во многих случаях этот метод препятствует ее использованию. Похоже, что такой подход был порожден страхом, неуверенностью и сомнениями. Было принято считать, что базы данных, SQL, транзакции и целостность данных “трудны” для понимания. Итогом стало решение никого не заставлять делать что-то “трудное”. В результате исповедующие такую философию разработчики относятся к базе данных как к черному ящику и поручают генерацию всего необходимого кода какому-то программному средству. Они стараются отгородиться множеством защитных уровней, чтобы даже не касаться этой “трудной” базы данных.

Я никогда не мог понять такой подход к разработке базы данных, частично потому, что лично мне изучение Java и С далось значительно труднее, чем изучение концепций, лежащих в основе баз данных. Теперь я знаю языки Java и С довольно хорошо, но для овладения практическими навыками их применения мне понадобилось намного больше времени, чем это было при освоении баз данных. В случае базы данных нужно иметь представление о ее работе, но вовсе не обязательно знать абсолютно все ее внутренние и внешние особенности. При программировании на С или Java/J2EE необходимо знать все внутренние и внешние функций, количество которых в этих языках поистине *огромно*.

Если вы строите приложение базы данных, то *наиболее важной частью программного обеспечения является сама база данных*. Поэтому члены успешно работающей команды разработки будут учитывать это обстоятельство и стремиться к максимальному изучению базы данных, уделяя ей основное внимание. Много раз мне доводилось присоединяться к командам разработки проектов, в которых был принят почти диаметрально противоположный подход.

Типичный сценарий мог бы выглядеть следующим образом.

- Разработчики полностью овладели инструментом для построения графических пользовательских интерфейсов или языком, который они используют для создания такого интерфейса (таким как Java). Во многих случаях за их плечами были долгие недели, если не месяцы, обучения.
- В команде потратили ноль часов на обучение Oracle и ноль часов на обретение опыта работы с Oracle. Большинство членов команды вообще не имеют опыта взаимодействия с базами данных. Они также получили указание быть “независимыми от базы данных” — указание (выданное руководством или выведенное из теоретического свода правил), следовать которому у них мало шансов по

многим причинам. Самой очевидной причиной является отсутствие у членов команды достаточных знаний о том, что собой представляют базы данных или что они делают, которые позволили бы прийти к наименьшему общему знаменателю среди баз данных.

- Разработчики сталкиваются с многочисленными проблемами, которые касаются производительности, целостности данных, зависания приложения и т.п. (но зато им удастся создавать очень привлекательные окна).

В результате неизбежного возникновения проблем с производительностью ко мне обращались за помощью в их решении (в прошлом, как обучающийся разработчик, я часто сам был причиной таких проблем). В одном конкретном случае я не сумел вспомнить точный синтаксис новой команды, которой нужно было воспользоваться. Попросив принести *справочное руководство по SQL*, я получил в руки документ для версии Oracle 6.0. Сама разработка велась на основе версии 7.3 — по прошествии пяти лет после выхода версии 6.0! Этот документ был единственным, с чем имели дело разработчики, но было похоже на то, что данный факт их совершенно не смущал. Они даже не думали о том, что средства, которые они должны были знать для успешного проведения трассировки и настройки, в версии 6 в действительности не существовали. Их не волновало, что за те пять лет, которые прошли с момента написания имеющейся в их распоряжении документации, появились новые функциональные возможности, такие как триггеры, хранимые процедуры и многие сотни других. Понять, *почему* им потребовалась помощь, не составило труда — чего нельзя было сказать о решении возникших у них проблем.

На заметку! Даже в наши дни я часто обнаруживаю, что разработчики приложений баз данных совершенно не уделяют время чтению документации. На своем веб-сайте asktom.oracle.com я часто встречаю вопросы вроде “как выглядит синтаксис для...” в сочетании с заявлениями “у нас нет документации, поэтому, пожалуйста, просто расскажите нам”. Я отказываюсь давать прямые ответы на многие из таких вопросов, а вместо этого отправляю спрашивающих к сетевой документации, доступной любому пользователю где угодно в мире. За последние 15 лет оправдания типа “у нас нет документации” или “у нас нет доступа к ресурсам” практически утратили свою актуальность. Распространение Интернета и появление таких сайтов, как otn.oracle.com (Oracle Technology Network — сеть технологии Oracle), делает отсутствие под рукой полного набора документации совершенно непростительным! Сегодня каждый имеет доступ ко всей документации; нужно просто читать ее, или — что еще проще — искать в ней ответы на вопросы.

Сама идея о том, что разработчики, которые строят *приложение базы данных*, должны быть изолированы от базы данных, представляется мне странной, но такая позиция существует. Многие продолжают верить, что разработчики не могут тратить время на освоение базы данных и в целом вовсе не обязаны что-либо знать о базе данных. Почему? Не единожды мне доводилось слышать утверждение наподобие “...ведь Oracle является наиболее масштабируемой базой данных в мире, поэтому моему персоналу не нужно ее изучать; она просто работает”. Действительно, Oracle — самая масштабируемая база данных в мире. Однако в среде Oracle наряду с эффективным и масштабируемым кодом достаточно легко написать неудачный и не поддающийся масштабированию код. Замените слово “Oracle” названием любого другого программного обеспечения, и утверждение останется справедливым. Не подлежит

сомнению: приложение, которое работает плохо, написать значительно проще, чем приложение, которое работает хорошо. Даже в самой масштабируемой в мире базе данных иногда очень легко получить в итоге однопользовательскую систему, если вы не знаете, что делаете. База данных — это инструмент, а неправильное применение любого инструмента может приводить к проблемам. Стали бы вы разбивать грецкие орехи щипцами для колки орехов так, как если бы это был молоток? Конечно, щипцы можно было бы использовать и так, но подобное применение этого инструмента является неподходящим, и дало бы в результате месиво (а возможно и несколько ушибленных пальцев). Игнорирование особенностей базы данных может привести к аналогичным последствиям.

Как-то меня пригласили принять участие в одном проекте, который зашел в тупик. Разработчики столкнулись с крупными проблемами, касающимися производительности; создавалось впечатление, что их система выполняла многие транзакции последовательно. Вместо того чтобы множество пользователей могли работать параллельно, каждый из них помещался в длинную очередь и был вынужден дожидаться, пока завершит работу пользователь, стоящий перед ним. Архитекторы системы ознакомили меня со своим творением — классической трехзвенной моделью. В этой системе веб-браузер должен был обмениваться данными с сервером приложений среднего звена, на котором функционировали страницы JSP (JavaServer Pages). В свою очередь, JSP-страницы должны были использовать еще один уровень, бины Enterprise JavaBeans (EJB), который выполнял все SQL-запросы. Код SQL в бинах EJB генерировался инструментом третьей стороны в независимой от базы данных манере.

В системе было очень трудно диагностировать что-либо, т.к. ни один из фрагментов кода не был снабжен каким-нибудь инструментарием или средствами трассировки. Снабжение кода инструментарием — это тонкое искусство превращения каждой второй строки разработанного кода в отладочный код, который в случае возникновения проблем с производительностью, емкостью или даже логикой позволяет точно отследить место, где они произошли. В сложившейся ситуации можно было утверждать лишь то, что источник проблемы находился между браузером и базой данных — другими словами, под подозрением оказалась вся система. База данных Oracle тщательно инструментирована, но приложение должно быть в состоянии включать и отключать инструментарий в соответствующих точках, а это как раз то, что не было предусмотрено проектом.

Таким образом, мы столкнулись с попыткой диагностирования проблем производительности, располагая лишь теми сведениями, которые можно было почерпнуть из самой базы данных. К счастью, в этом случае решение было довольно простым. Просмотр таблиц `V$` (эти таблицы являются одним из способов предоставления доступа к инструментарии и статистической информации Oracle) показал, что основное соперничество велось за единственную таблицу — таблицу очереди на обработку. Приложение должно было помещать записи в эту таблицу, а другой набор процессов — извлекать записи из нее и обрабатывать. Исследовав эту таблицу более внимательно, мы обнаружили битовый индекс на одном из столбцов (более подробную информацию о битовых индексах можно найти в главе 11). Причина такого выбора состояла в том, что этот столбец — флаг обработки — мог содержать только два значения: `Y` и `N`. Записи, вставленные в таблицу, должны были иметь в этом столбце значение `N` (не обработана).

После считывания и обработки записи другими процессами эти процессы должны были изменять значение в столбце с N на Y , указывая на то, что запись обработана. Разработчикам нужно было быстро находить записи со значением N , поэтому столбец был проиндексирован. Они где-то прочитали, что битовые индексы предназначены для столбцов с низким кардинальным числом (столбцов, содержащих небольшое количество отличающихся значений), поэтому сочли такой вид индекса вполне естественным. (Попробуйте поискать в Google строку “when to use bitmap indexes” (когда использовать битовые индексы); словосочетание “low-cardinality” (низкое кардинальное число) будет встречаться повсеместно. К счастью, сейчас имеется также много статей, опровергающих эту излишне упрощенную концепцию.)

Именно битовый индекс был причиной всех проблем. В битовом индексе единственная запись ключа указывает на множество строк — сотни и более. Когда ключ битового индекса обновляется (и тем самым блокируется), то сотни записей, на которые он указывает, также блокируются. Таким образом, вставка новой записи со значением N приводит к блокированию записей N в битовом индексе, по существу блокируя также и сотни других записей со значением N . Между тем процесс, который читает эту таблицу и обрабатывает записи, не имеет возможности изменить значение какой-либо записи с N на Y , т.к. для этого ему пришлось бы заблокировать тот же самый ключ битового индекса. Фактически другие сеансы, которые пытаются всего лишь вставить новую запись в эту таблицу, также должны блокироваться, поскольку они предпринимают попытку заблокировать ту же самую запись битового ключа. Короче говоря, разработчики создали таблицу, вставлять и обновлять записи в которой мог только один пользователь за раз! Мы можем легко в этом убедиться, используя простой сценарий.

На заметку! Для демонстрации проблем блокировки и параллельного доступа на протяжении всей книги я буду использовать автономные транзакции. Я твердо убежден, что автономные транзакции являются средством, которое в Oracle не должны были открывать разработчикам по той простой причине, что большинство разработчиков не знают, когда и как применять их правильно. Некорректное использование автономных транзакций может — и будет — приводить к проблемам нарушения логической целостности данных. Помимо их применения в качестве демонстрационного инструмента, автономные транзакции имеют только еще одно предназначение — служить механизмом протоколирования ошибок. Чтобы запротолировать ошибку в блоке исключения, понадобится записать сведения о ней в таблицу и зафиксировать ее, не фиксируя ничего другого. Это было бы допустимым использованием автономной транзакции. Если вы применяете автономную транзакцию вне контекста протоколирования или демонстрации концепции, значит, вы почти наверняка делаете что-то совершенно неправильное.

Здесь я использую автономную транзакцию в базе данных, чтобы иметь две параллельные транзакции в одном сеансе. Автономная транзакция запускает “подтранзакцию” отдельно и независимо от любой уже установившейся транзакции в сеансе. Эта автономная транзакция ведет себя так, как если бы она находилась в совершенно отдельном сеансе — фактически родительская транзакция приостанавливается. Автономная транзакция может быть заблокирована родительской транзакцией (как мы увидим) и, более того, автономная транзакция не может видеть незафиксированных модификаций, произведенных родительской транзакцией. Например:

```
EODA@ORA12CR1> create table t
  2  ( processed_flag varchar2(1)
  3  );
Table created.
Таблица создана.

EODA@ORA12CR1> create bitmap index
  2  t_idx on t(processed_flag);
Index created.
Индекс создан.

EODA@ORA12CR1> insert into t values ( 'N' );
1 row created.
1 строка создана.

EODA@ORA12CR1> declare
  2      pragma autonomous_transaction;
  3  begin
  4      insert into t values ( 'N' );
  5      commit;
  6  end;
  7  /
declare
*
```

ERROR at line 1:
ORA-00060: deadlock detected while waiting for resource
ORA-06512: at line 4
ОШИБКА в строке 1:
ORA-00060: в процессе ожидания ресурса обнаружена взаимная блокировка
ORA-06512: в строке 4

Совет. За подробными сведениями о настройке приглашения командной строки SQL для отображения информации о среде, такой как имя пользователя и имя базы данных, обращайтесь в раздел “Настройка среды” в начале книги.

Поскольку использовалась автономная транзакция и создана подтранзакция, произошла взаимная блокировка, т.е. вторая вставка была заблокирована первой вставкой. Если бы имелось два отдельных сеанса, то никакой взаимной блокировки не возникло бы. Вместо этого вторая вставка просто была бы заблокирована и ожидала, пока не произойдет фиксация или откат первой вставки. Именно этот симптом проявился в проекте, с которым я имел дело — проблема с блокировкой и сериализацией.

Таким образом, в приведенном случае мы столкнулись с отсутствием понимания функциональной возможности базы данных (битовых индексов), в результате чего с самого начала база данных была обречена на плохую масштабируемость. Проблему усугубляло также и то, что не было никаких причин для написания кода работы с очередью. В базе данных доступны встроенные средства работы с очередями, которые появились в версии Oracle 8.0, вышедшей в 1997 году. Встроенное средство работы с очередями позволяет иметь много производителей (сеансов, которые вставляют необработанные записи), одновременно помещающих сообщения в очередь, и много потребителей (сеансов, которые получают записи со значением N для обработки), параллельно получающих эти сообщения. Это значит, что для реализации

очереди в базе данных писать какой-то специальный код не требуется. Разработчики должны были воспользоваться встроенным средством. И они могли это сделать, но им помешала полная неосведомленность о его существовании.

К счастью, как только причина проблемы была выявлена, ее устранение не составило особого труда. Столбец флага обработки нуждался в индексе, но не в битовом. Здесь необходим обычный индекс со структурой В-дерева. Это вызвало непродолжительные споры, поскольку никто не верил, что использование обычного индекса в столбце, содержащем два отличающихся значения, является удачной идеей. Однако после выполнения моделирования (я большой поклонник моделирования, тестирования и экспериментирования) мы смогли удостовериться в правильности этого подхода.

На заметку! Мы создаем индексы — любого типа — обычно для поиска небольшого количества строк в крупном наборе данных. В этом случае количество строк, которые мы хотели искать через индекс, было равно *одному*. Мы нуждались в поиске одной необработанной записи. Одна строка — это очень маленькое количество, следовательно, индекс является подходящим решением. Подошел бы индекс любого типа. Индекс со структурой В-дерева весьма удобен при поиске одиночной записи в большом наборе записей.

При создании индекса необходимо выбрать один из следующих подходов.

- Просто создать индекс на столбце флага обработки.
- Создать индекс на столбце флага обработки только для тех записей, в которых флагом обработки является N — т.е. индексировать только интересующие значения. Как правило, мы не хотим применять индекс, когда флаг обработки равен Y, из-за того, что подавляющее большинство записей в таблице имеют флаг обработки Y. Обратите внимание, что здесь не было сформулировано “мы никогда не хотим...”. Если по какой-то причине необходимо часто подсчитывать количество обработанных записей, то наличие индекса для таких записей может быть очень полезно.

В главе, посвященной индексации, мы более детально рассмотрим оба подхода. В конце концов, мы создали очень маленький индекс только на записях, в которых флаг обработки был равен N. Доступ к этим записям был исключительно быстрым, а подавляющее большинство записей со значением Y вообще не участвовали в этом индексе. Мы использовали индекс, основанный на функции `decode(processed_flag, 'N', 'N')`, которая возвращает либо N, либо NULL; учитывая, что ключ NULL в обычный индекс В-дерева не помещается, в результате мы проиндексировали только записи со значением N.

На заметку! В главе 11 вы найдете дополнительные сведения о NULL и индексации.

Закончилась ли на этом история? Вовсе нет. Решение, находящееся в руках моего клиента, все еще оставалось неоптимальным. По-прежнему приходилось выполнять сериализацию “извлечения из очереди” необработанных записей. Мы могли быстро найти первую необработанную запись, используя `select * from queue_table where decode(processed_flag, 'N', 'N')='N' FOR UPDATE`, но выполнять эту операцию мог только один сеанс за раз. В проекте применялась версия Oracle 10g,

и потому нельзя было воспользоваться относительно новым средством SKIP LOCKED, которое появилось только в Oracle 11g Release 1. Конструкция SKIP LOCKED позволила бы многим сеансам параллельно находить первую неблокированную и необработанную запись, после чего заблокировать ее и обработать. Вместо этого мы должны были реализовать код для нахождения первой неблокированной записи и заблокировать ее вручную. Такой код в Oracle 10g и предшествующих версиях в целом выглядел бы так, как показано ниже. Мы начинаем с создания таблицы с требуемым индексом, который был описан ранее, и наполнения ее некоторыми данными:

```
EODA@ORA12CR1> create table t
  2 ( id          number primary key,
  3   processed_flag varchar2(1),
  4   payload      varchar2(20)
  5 );
Table created.
Таблица создана.
EODA@ORA12CR1> create index
  2 t_idx on
  3 t( decode( processed_flag, 'N', 'N' ) );
Index created.
Индекс создан.
EODA@ORA12CR1> insert into t
  2 select r,
  3        case when mod(r,2) = 0 then 'N' else 'Y' end,
  4        'payload ' || r
  5   from (select level r
  6         from dual
  7         connect by level <= 5)
  8 /
5 rows created.
5 строк создано.
EODA@ORA12CR1> select * from t;
```

```
      ID P PAYLOAD
-----
1 Y payload 1
2 N payload 2
3 Y payload 3
4 N payload 4
5 Y payload 5
```

Затем понадобится найти все необработанные записи. Для каждой записи мы отправляем базе данных вопрос вида: “Заблокирована ли уже эта запись? Если нет, то заблокировать ее и предоставить”. Код выглядит примерно так:

```
EODA@ORA12CR1> create or replace
  2 function get_first_unlocked_row
  3 return t%rowtype
  4 as
  5     resource_busy exception;
  6     pragma exception_init( resource_busy, -54 );
  7     l_rec t%rowtype;
```

```

8  begin
9      for x in ( select rowid rid
10                 from t
11                 where decode(processed_flag, 'N', 'N') = 'N')
12      loop
13          begin
14              select * into l_rec
15                  from t
16                  where rowid = x.rid and processed_flag='N'
17                  for update nowait;
18              return l_rec;
19          exception
20              when resource_busy then null;
21              when no_data_found then null;
22          end;
23      end loop;
24  return null;
25  end;
/
Function created.
Функция создана.

```

На заметку! В приведенном выше коде был выполнен код DDL — CREATE OR REPLACE FUNCTION. Непосредственно перед запуском оператора DDL происходит автоматическая фиксация, поэтому здесь присутствует неявный вызов COMMIT. Вставляемые строки фиксируются в базе данных — и этот факт необходим для корректной работы последующих примеров. Вообще говоря, факт автоматической фиксации используется повсеместно в оставшейся части книги. Если вы запускаете эти примеры без выполнения CREATE OR REPLACE, не забудьте сначала сделать COMMIT!

Если теперь применить две разные транзакции, можно заметить, что они получают разные записи. Кроме того, они получают их одновременно (снова с использованием автономных транзакций для демонстрации проблем параллельного доступа):

```

EODA@ORA12CR1> declare
2  l_rec t%rowtype;
3  begin
4  l_rec := get_first_unlocked_row;
5  dbms_output.put_line( 'I got row ' || l_rec.id || ', ' || l_rec.payload );
6  end;
7  /
I got row 2, payload 2
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно выполнена.

EODA@ORA12CR1> declare
2  pragma autonomous_transaction;
3  l_rec t%rowtype;
4  begin
5  l_rec := get_first_unlocked_row;
6  dbms_output.put_line( 'I got row ' || l_rec.id || ', ' || l_rec.payload );
7  commit;
8  end;
9  /

```

I got row 4, payload 4
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно выполнена.

Теперь в Oracle 11g Release 1 и последующих версиях показанную выше логику можно реализовать с применением конструкции **SKIP LOCKED**. В приведенном далее примере мы снова выполним две параллельные транзакции и обнаружим, что каждая находит и блокирует записи одновременно.

```
EODA@ORA12CR1> declare
  2   l_rec t%rowtype;
  3   cursor c
  4   is
  5   select *
  6   from t
  7   where decode(processed_flag, 'N', 'N') = 'N'
  8   FOR UPDATE
  9   SKIP LOCKED;
10 begin
11   open c;
12   fetch c into l_rec;
13   if ( c%found )
14   then
15       dbms_output.put_line( 'I got row ' || l_rec.id || ', '
                             || l_rec.payload );
16   end if;
17   close c;
18 end;
19 /
```

I got row 2, payload 2
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно выполнена.

```
EODA@ORA12CR1> declare
  2   pragma autonomous_transaction;
  3   l_rec t%rowtype;
  4   cursor c
  5   is
  6   select *
  7   from t
  8   where decode(processed_flag, 'N', 'N') = 'N'
  9   FOR UPDATE
10   SKIP LOCKED;
11 begin
12   open c;
13   fetch c into l_rec;
14   if ( c%found )
15   then
16       dbms_output.put_line( 'I got row ' || l_rec.id || ', '
                             || l_rec.payload );
17   end if;
18   close c;
19   commit;
20 end;
21 /
```

I got row 4, payload 4
 PL/SQL procedure successfully completed.
Процедура PL/SQL успешно выполнена.

Оба продемонстрированных выше “решения” должны помочь справиться со второй проблемой сериализации, с которой столкнулся мой клиент при обработке сообщений. Но намного ли проще получилось бы решение, если бы клиент просто применял средство *Advaced Queuing* и вызывал процедуру *DBMS_AQ.DEQUEUE*? Чтобы решить проблему сериализации поставщика сообщений, нам пришлось реализовать индекс, основанный на функции. Чтобы решить проблему сериализации для потребителя, пришлось использовать этот основанный на функции индекс для извлечения записей и написания кода. Таким образом, мы устранили основную проблему, порожденную неполным пониманием применяемых средств, и лишь после тщательного изучения выяснили, что система была недостаточно инструментирована. Однако перечисленные ниже проблемы остались нерешенными.

- Приложение было построено без единого соображения относительно масштабирования на уровне базы данных.
- В приложении была реализована функциональность (таблица очереди), которую база данных уже предоставляла, причем с высокой степенью параллелизма и масштабируемости. Речь идет о программном обеспечении *AQ (Advaced Queuing — расширенная организация очередей)*, встроенном в базу данных. Разработчики пытались построить эту функциональность заново.
- Опыт показал, что от 80 до 90 (и более!) процентов всех действий по настройке должны быть выполнены на уровне приложения (обычно в интерфейсном коде, который читает и записывает в базу данных), а не на уровне базы данных.
- Разработчики не имели ни малейшего представления о том, что делают бины в базе данных, или где искать потенциальные проблемы.

На этом проблемы данного проекта не исчерпывались. Нам также необходимо было выяснить следующие моменты.

- Как настроить SQL-код, не изменяя его. Вообще говоря, это очень трудно. Владеть до некоторой степени этой магией в *Oracle 10g* и последующих версиях можно с помощью *SQL Profiles* (эта опция требует лицензии на пакет настройки *Oracle (Oracle Tuning Pack)*), в *Oracle 11g* и выше — посредством расширенной статистики, а в *Oracle 12c* и далее — с использованием адаптивной оптимизации запросов. Однако неэффективный SQL-код остается неэффективным SQL-кодом.
- Как измерить производительность.
- Как выявить имеющиеся узкие места.
- Как и что индексировать. И так далее.

К концу недели разработчики, ранее сторонившиеся базы данных, были поражены теми возможностями, которые она могла предоставить, и простотой получения информации. Самое главное, они увидели, насколько ее применение может повысить производительность приложения. В итоге разработчики добились успеха, отстав от графика всего на несколько недель.

Мой рассказ о мощи возможностей базы данных приведен здесь не для того, чтобы критиковать средства или технологии, подобные Hibernate, EJB и контейнерному подходу к реализации постоянства. Критика касается намеренного игнорирования базы данных, особенностей ее функционирования и способов ее использования. Примененные в рассмотренном случае технологии вполне себя оправдали — разумеется, после того, как разработчики получили некоторое представление о самой базе данных.

Подводя итог, можно сказать, что база данных, как правило, служит краеугольным камнем приложения. Если она работает недостаточно хорошо, то все остальное не имеет особого значения. Что делать, если используемый черный ящик плохо работает? Единственное что можно предпринять в этой ситуации — смотреть на него и гадать, почему он работает не так, как нужно. Вы не можете ни починить его, ни настроить. Вы просто не понимаете, каким образом он работает — и вы сознательно поставили себя в такое положение. Альтернативой является отстаиваемый мною подход: понимать базу данных, знать, как она работает, что она может предложить, и в полной мере задействовать ее потенциал.

Как следует (и как не следует) разрабатывать приложения баз данных

Но достаточно строить гипотезы, сейчас, по крайней мере. В оставшейся части этой главы будет применяться более практический подход с обсуждением того, почему знание базы данных и ее элементов определенно позволит значительно приблизиться к успешной реализации (не переписывая приложение по нескольку раз). Некоторые проблемы устранять проще, если понимать, каким образом их искать. Устранение других проблем может потребовать радикального переписывания кода. Одна из целей настоящей книги — в первую очередь помочь вам избежать этих проблем.

На заметку! В приведенных далее разделах некоторые основные функциональные средства Oracle обсуждаются без подробного освещения их сути и всех особенностей применения. При необходимости я буду отправлять за деталями либо к одной из последующих глав книги, либо к соответствующей документации Oracle.

Архитектура Oracle

Мне приходилось сотрудничать со многими клиентами, использующими крупные производственные приложения, которые были “перенесены” в Oracle из другой платформы баз данных (например, SQL Server). Слово “перенесено” взято в кавычки потому, что большинство встречаемых мною адаптаций сводились к точке зрения “найти минимальные изменения, которые обеспечили бы успешную компиляцию и выполнение кода SQL Server на платформе Oracle”. Откровенно говоря, приложения, построенные в результате такого подхода к делу, попадались мне чаще всего, поскольку именно они требовали наибольшей помощи. Я вовсе не критикую SQL Server в этом отношении — ведь справедливо и обратное! Перенос приложения Oracle и его помещение с минимальными изменениями в среду SQL Server приведет к получению столь же плохо работающего кода, как и наоборот; проблема имеет обоюдный характер.

Однако в одном конкретном случае архитектура SQL Server и способ применения SQL Server действительно были навязаны реализацией Oracle. В качестве конечной цели ставилось масштабирование, но обратившиеся ко мне разработчики на самом деле не хотели переходить на другую базу данных. Они хотели провести перенос с минимальными усилиями со своей стороны, и потому оставили архитектуру в основном прежней — на уровне клиента и базы данных. Это решение имело два важных последствия.

- Архитектура подключений в Oracle осталась той же самой, что и использованная в SQL Server.
- Разработчики применяли литеральный (неограниченный) SQL-код.

Эти два обстоятельства повлекли за собой построение системы, которая не могла поддерживать необходимую пользовательскую нагрузку (серверу базы данных просто не хватало доступной памяти) и обладала крайне низкой производительностью.

Используйте единственное подключение в Oracle

В среде SQL Server открытие подключения к базе данных для каждого параллельно выполняющегося оператора является весьма распространенной практикой. При выполнении пяти запросов в среде SQL Server вполне можно встретить пять подключений. С другой стороны, в Oracle вне зависимости от того, сколько запросов нужно выполнить, пять или даже пятьсот, максимальное количество открываемых подключений равно единице. В результате то, что принято в SQL Server, совершенно не рекомендуется применять в среде Oracle; наличие множества подключений к базе данных просто нежелательно.

Однако как раз это и было сделано. Простое веб-приложение для каждой веб-страницы может открывать 5, 10, 15 и более подключений, а это значит, что сервер мог поддерживать только $\frac{1}{5}$, $\frac{1}{10}$, $\frac{1}{15}$ и менее параллельно работающих пользователей от того числа, которое должен. Кроме того, была предпринята попытка использования базы данных на обычной платформе Windows — в среде простого сервера Windows без доступа к Datacenter-версии Windows Server. В результате архитектура Windows с единственным процессом ограничила общий объем оперативной памяти, доступной серверу баз данных Oracle, до приблизительно 1,75 Гбайт. Поскольку каждое подключение Oracle занимает, как минимум, определенный фиксированный объем памяти, возможности масштабирования количества пользователей, работающих с приложением, были существенно ограничены. Объем оперативной памяти сервера составлял 8 Гбайт, но из них можно было использовать только около 2 Гбайт.

На заметку! В среде 32-разрядной ОС Windows доступны способы использования большего объема оперативной памяти, такие как ключ `/AWE`, но для этого требуются версии ОС, которые в описанной ситуации не применялись.

Существовало три подхода к решению этой проблемы, причем все три были достаточно трудоемкими — и это после завершения “переноса”!

Были доступны следующие варианты.

- Изменить архитектуру приложения, чтобы оно могло получить преимущества выполнения в среде Oracle, и во время генерирования страницы применять одно подключение, а не от 5 до 15. Это единственное решение, которое действительно устранило бы проблему.
- Провести модернизацию ОС (отнюдь не простая задача) и задействовать модель с большим объемом доступной памяти, предлагаемую версией Windows Server Datacenter (что само по себе совсем не просто, т.к. сопряжено со сложной настройкой базы данных, с определением косвенных буферов данных и других нестандартных параметров).
- Перенести базу данных из Windows в среду какой-то другой ОС, которая поддерживает множество процессов. Это фактически позволило бы базе данных задействовать всю установленную оперативную память. На 32-разрядной платформе Windows вы ограничены примерно 2 Гбайт памяти для комбинированных областей PGA/SGA (2 Гбайт для обеих вместе), поскольку они выделяются единственным процессом. При использовании платформы с множеством процессов, которая также является 32-разрядной, вы будете ограничены примерно 2 Гбайт для SGA и 2 Гбайт на процесс для PGA, что существенно больше, чем обеспечивает 32-разрядная платформа Windows.

Как видите, ни одна из этих возможностей не относилась к решениям, о которых можно было бы сказать: “Хорошо, мы сделаем это до обеда”. Каждая из них представляла собой сложное решение проблемы, которую проще всего было решить на этапе переноса базы данных, пока вносились изменения в код и в наиболее важные элементы системы. Более того, простое тестирование масштабирования, проведенное до развертывания системы, позволило бы выявить проблемы подобного рода еще до того, как конечные пользователи начали бы испытывать какие-то неудобства.

Используйте переменные привязки

Если бы я писал книгу о том, как строить *немасштабируемые* приложения Oracle, то первая и последняя глава называлась бы “Не используйте переменные привязки”. Отказ от применения переменных привязки является основной причиной возникновения проблем с производительностью и главным ограничителем возможностей масштабирования, не говоря уже об огромной степени риска в плане безопасности. Способ работы Oracle с разделяемым пулом (очень важная структура данных в разделяемой памяти) в большинстве случаев определяется разработчиками, использующими переменные привязки. Если вы хотите заставить транзакционную реализацию Oracle работать медленно, вплоть до полной ее остановки, просто откажитесь от применения переменных привязки.

Переменная привязки — это метка-заполнитель в запросе. Например, для извлечения записи сотрудника 123 можно выполнить следующий запрос:

```
select * from emp where empno = 123;
```

В качестве альтернативы можно запустить такой запрос:

```
select * from emp where empno = :empno;
```

В типичной системе запрос информации о сотруднике 123 вполне может быть выполнен один или два раза и больше никогда на протяжении длительного периода времени. Позже может требоваться информация о сотруднике 456, затем — о сотруднике 789 и т.д. Или, как в предшествующих операторах SELECT, если вы не указываете в своих операторах вставки переменные привязки, то значения первичного ключа будут жестко закодированы в них, и мне известен тот факт, что такие операторы вставки никогда не смогут использоваться повторно! Если в запросе применяются литералы (константы), то каждый запрос оказывается совершенно новым, никогда ранее не выполнявшимся в базе данных. Он должен быть синтаксически разобран, определен (произведено распознавание имен), проверен на соблюдение правил безопасности, оптимизирован и т.п.

Короче говоря, каждый запускаемый уникальный оператор будет требовать компиляции при каждом своем выполнении.

Во втором запросе использовалась переменная привязки :empno, значение которой передается во время выполнения запроса. Этот запрос компилируется только один раз, а затем план запроса сохраняется в разделяемом пуле (библиотечном кеше), из которого он может быть извлечен и применен повторно. Разница между этим двумя методами с точки зрения производительности и масштабируемости не просто велика — она огромна.

Из предыдущего объяснения должно быть совершенно понятно, что синтаксический разбор оператора с жестко закодированными переменными (называемый *полным* разбором) будет проходить дольше и потреблять намного больше ресурсов, чем повторное использование уже разобранного плана запроса (которое называется *частичным* разбором). Однако степень снижения количества пользователей, которых может поддерживать система, при первом методе может оказаться не настолько ясной. Очевидно, что отчасти это объясняется увеличением потребления системных ресурсов, но более значительный фактор связан с влиянием механизмов защелок, применяемых к библиотечному кешу.

При выполнении полного разбора запроса база данных будет дольше хранить определенные низкоуровневые устройства последовательной обработки, называемые *защелками* (или *внутренними блокировками*); за дополнительной информацией о них обращайтесь в главу 6. Эти защелки защищают структуры данных в разделяемой памяти Oracle от одновременных изменений двумя сеансами (в противном случае структуры данных были бы повреждены) и от считывания структуры данных кем-либо во время ее изменения. Чем дольше и чаще приходится “защелкивать” эти структуры данных, тем более длинной будет становиться очередь для получения таких защелок. Это приведет к монополизации ограниченных ресурсов. Временами компьютер может выглядеть недогруженным, тем не менее, все действия в базе данных будут выполняться очень медленно. Внешне все выглядит так, будто кто-то владеет одним из механизмов последовательной обработки, создавая очередь — достичь максимальной производительности не удастся. Достаточно наличия в базе данных одного неправильно ведущего себя приложения, чтобы производительность всех приложений значительно снизилась. Единственное небольшое приложение без переменных привязки со временем приведет к удалению из разделяемого пула всех SQL-запросов, принадлежащих остальным хорошо настроенным приложениям. Одной ложки дегтя хватит, чтобы испортить бочку меда.

На заметку! Чтобы увидеть отличие между полным и частичным разбором в действии, рекомендуется пересмотреть демонстрационный видеоролик, доступный по ссылке <http://tinyurl.com/RWP-OLTP-PARSING>. Он был смонтирован командой, с которой я работал — командой Real World Performance (Производительность в реальном мире) из Oracle. В нем наглядно показана разница между полным и частичным разбором — она близка к отличию на порядок! В транзакционной системе, архитектура которой ориентирована на использование переменных привязки, можно добиться десятикратного увеличения скорости выполнения в случае их применения. Вы можете использовать эту короткую визуальную презентацию, чтобы убедить других разработчиков о высоком влиянии переменных привязки (либо их отсутствия) на производительность!

Если же вы применяете переменные привязки, то любой, кто отправляет точно такой же запрос, ссылающийся на тот же самый объект, будет использовать скомпилированный план из пула. Компиляция процедуры производится только один раз, после чего ею можно пользоваться снова и снова. Это очень эффективно, и именно на такую работу рассчитана база данных. При этом не только уменьшается объем задействованных ресурсов (частичный разбор требует намного меньшего количества ресурсов), но также сокращается время действия зашелок, да и потребность в них будет возникать реже. В результате производительность и масштабируемость приложений значительно возрастают.

Чтобы получить общее представление о том, насколько огромна разница в производительности в случае применения переменных привязки, можно запустить небольшой тест. В этом тесте мы просто вставим несколько строк в следующую таблицу:

```
EODA@ORA12CR1> create table t ( x int );  
Table created.  
Таблица создана.
```

Теперь создадим две очень простые хранимые процедуры. Обе они будут вставлять в таблицу числа от 1 до 10 000, но первая процедура использует единственный SQL-оператор с переменной привязки:

```
EODA@ORA12CR1> create or replace procedure proc1  
2 as  
3 begin  
4     for i in 1 .. 10000  
5     loop  
6         execute immediate  
7             'insert into t values ( :x )' using i;  
8     end loop;  
9 end;  
10 /  
Procedure created.  
Процедура создана.
```

Вторая процедура конструирует уникальный SQL-оператор для каждой вставляемой строки:

```
EODA@ORA12CR1> create or replace procedure proc2  
2 as  
3 begin
```

```

4      for i in 1 .. 10000
5      loop
6          execute immediate
7              'insert into t values ( '||i||')';
8      end loop;
9  end;
10 /

```

Procedure created.

Процедура создана.

Единственное отличие между этими двумя процедурами состоит в том, что в одной применяется переменная привязки, а в другой — нет. Обе процедуры используют динамический SQL-код с идентичной логикой. Разница заключается только в применении переменной привязки в первой процедуре. Мы можем оценить два подхода с использованием простого разработанного мною инструмента runstats для детального сравнения:

```

EODA@ORA12CR1> exec runstats_pkg.rs_start
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно выполнена.

EODA@ORA12CR1> exec proc1
PL/SQL procedure successfully completed.

EODA@ORA12CR1> exec runstats_pkg.rs_middle
PL/SQL procedure successfully completed.

EODA@ORA12CR1> exec proc2
PL/SQL procedure successfully completed.

EODA@ORA12CR1> exec runstats_pkg.rs_stop(9500)
Run1 ran in 34 cpu hsecs
Run2 ran in 432 cpu hsecs
run 1 ran in 7.87% of the time

```

На заметку! За подробными сведениями о runstats и других утилитах обращайтесь в раздел “Настройка среды” в начале этой книги. Вы можете получить несколько отличающиеся значения для процессорного времени или любой другой метрики. На это влияет версия Oracle, операционная система и аппаратная платформа. Общая идея будет той же, но точные значения наверняка окажутся другими.

Предыдущий результат четко показывает, что в плане процессорного времени вставка 10 000 строк заняла значительно больше времени и ресурсов в случае без переменных привязки по сравнению с ситуацией, когда они применялись. Фактически вставка строк без участия переменных привязки потребовала более чем на порядок больше процессорного времени. Для каждой вставки без переменных привязки практически все время, затраченное на выполнение оператора, уходило просто на разбор этого оператора! Но есть кое-что похуже. Если мы посмотрим другую информацию, то обнаружим существенное отличие в утилизации ресурсов при каждом подходе:

| Name | Run1 | Run2 | Diff |
|--------------------------------|------------|-------------|-------------|
| STAT...CCursor + sql area evic | 2 | 9,965 | 9,963 |
| STAT...enqueue requests | 35 | 10,012 | 9,977 |
| STAT...enqueue releases | 34 | 10,012 | 9,978 |
| STAT...execute count | 10,020 | 20,005 | 9,985 |
| STAT...opened cursors cumulati | 10,019 | 20,005 | 9,986 |
| STAT...table scans (short tabl | 3 | 10,000 | 9,997 |
| STAT...sorts (memory) | 3 | 10,000 | 9,997 |
| STAT...parse count (hard) | 2 | 10,000 | 9,998 |
| LATCH.session allocation | 5 | 10,007 | 10,002 |
| LATCH.session idle bit | 17 | 10,025 | 10,008 |
| STAT...db block gets | 10,447 | 30,376 | 19,929 |
| STAT...db block gets from cach | 10,447 | 30,376 | 19,929 |
| STAT...db block gets from cach | 79 | 20,037 | 19,958 |
| LATCH.shared pool simulator | 8 | 19,980 | 19,972 |
| STAT...calls to get snapshot s | 22 | 20,003 | 19,981 |
| STAT...parse count (total) | 18 | 20,005 | 19,987 |
| LATCH.call allocation | 4 | 20,016 | 20,012 |
| LATCH.enqueue hash chains | 70 | 20,211 | 20,141 |
| STAT...consistent gets | 266 | 40,093 | 39,827 |
| STAT...consistent gets from ca | 266 | 40,093 | 39,827 |
| STAT...consistent gets pin (fa | 219 | 40,067 | 39,848 |
| STAT...consistent gets pin | 219 | 40,067 | 39,848 |
| STAT...calls to kcmgcs | 117 | 40,085 | 39,968 |
| STAT...session logical reads | 10,713 | 70,469 | 59,756 |
| STAT...recursive calls | 10,058 | 70,005 | 59,947 |
| STAT...KTFB alloc space (block | 196,608 | 131,072 | -65,536 |
| LATCH.cache buffers chains | 51,835 | 171,570 | 119,735 |
| LATCH.row cache objects | 206 | 240,686 | 240,480 |
| LATCH.shared pool | 20,090 | 289,899 | 269,809 |
| STAT...session pga memory | 65,536 | -262,144 | -327,680 |
| STAT...logical read bytes from | 87,760,896 | 577,282,048 | 489,521,152 |

Run1 latches total versus runs -- difference and pct

| Run1 | Run2 | Diff | Pct |
|--------|---------|---------|-------|
| 73,620 | 784,913 | 711,293 | 9.38% |

PL/SQL procedure successfully completed.

Процедура PL/SQL успешно выполнена.

Утилита runstats генерирует отчет, который демонстрирует разницу в использовании защепок, а также отличия в статистике. Здесь runstats было указано напечатать все отличия, превышающие значение 9 500. Можно заметить, что *полный разбор* выполнялся 2 раза при первом подходе с применением переменных привязки и 10 000 раз при подходе без переменных привязки (по одному разу для каждой вставки). Однако разница в количестве операций разбора — лишь вершина айсберга. Здесь также видно, что в случае подхода без переменных привязки использовалось на порядок больше “защепок”, чем с переменными привязки. Эта разница может породить вопрос о том, что собой представляет защелка?

Давайте ответим на этот вопрос. Защелка (latch) — это тип блокировки, применяемый для сериализации доступа к разделяемым структурам данных, которые находятся в глобальной системной области (System Global Area — SGA), где Oracle хранит разобранный скомпилированный SQL-код. При модификации чего-либо в

этой разделяемой структуре вы должны позаботиться о том, чтобы в каждый момент времени только одному процессу был разрешен доступ к ней. (Очень плохо, если два процесса или потока попытаются обновить одну и ту же структуру в памяти одновременно — ее повреждение неизбежно.) Поэтому в Oracle используется механизм *защелкивания*, который представляет собой облегченный метод блокировки для обеспечения последовательного доступа. Пусть слово “облегченный” не вводит в заблуждение. Защелки являются устройствами сериализации, которые в каждый момент времени разрешают доступ к структуре данных в памяти только одному процессу. Защелки, применяемые реализацией полного разбора, являются одними из наиболее часто используемых. Сюда входят защелки для разделяемого пула и библиотечного кеша. Это “долговременные” защелки, за которые часто соперничают пользователи. Все это означает, что при увеличении количества пользователей, пытающихся одновременно выполнить полный разбор SQL-кода, с течением времени производительность будет ухудшаться в геометрической прогрессии. Чем больше пользователей выполняет разбор, тем их большее количество ждет возможности защелкнуть разделяемый пул, и тем длиннее очереди и дольше ожидание.

Выполнение SQL-операторов без переменных привязки во многом подобно компиляции подпрограммы перед каждым вызовом метода. Представьте себе, что клиентам поставлен исходный код Java, который требует перед вызовом метода любого класса запускать компилятор Java, компилировать класс, выполнять метод, а затем отбрасывать байт-код. Для повторного вызова этого же метода приходится предпринимать те же самые действия: компилировать его, выполнить и затем отбросить байт-код. Такой подход совершенно неприемлем в приложении, неприемлем он и в базе данных.

Другое следствие отказа от использования переменных привязки для разработчиков, применяющих конкатенацию строк, касается безопасности, а именно — уязвимости к *внедрению SQL* (SQL injection). Если вы еще не знакомы с этим термином, я советую на какое-то время отвлечься от книги и поискать в Интернете упоминания “SQL injection”. По состоянию на момент написания этой книги вы получите свыше пяти миллионов ссылок. Проблема внедрения SQL хорошо документирована.

На заметку! Внедрение SQL — это брешь в системе безопасности, которая возникает, когда разработчик принимает данные, введенные конечным пользователем, объединяет их с запросом, затем компилирует и выполняет получившийся запрос. На самом деле разработчик принимает от конечного пользователя фрагменты SQL-кода, компилирует их и выполняет. Такой прием потенциально позволяет конечному пользователю модифицировать SQL-оператор так, чтобы он делал то, что разработчик оператора не намеревался. Это очень похоже на ситуацию, когда оставляют терминал с открытым сеансом SQL*Plus, который был подключен с привилегиями пользователя SYSDBA. Тем самым вы приглашаете кого угодно сесть и ввести любую команду, скомпилировать ее и затем выполнить. Последствия могут оказаться катастрофическими.

Фактически если вы не используете переменные привязки и применяете прием с конкатенацией строк, как в показанной ранее процедуре PROC2, то код становится уязвимым для атак внедрением SQL и потому должен быть тщательно пересмотрен. Причем это должен сделать кто-то другой, а не разработчик, написавший этот код, потому что код должен быть пересмотрен критически и объективно. Если проверя-

ющий является партнером автора кода или хуже того — другом или подчиненным, то проверка будет не столь критической, какой должна быть. К коду, разработанному без использования переменных привязки, следует относиться с подозрением — он должен быть скорее исключением, нежели нормой.

Чтобы продемонстрировать, насколько коварным может быть внедрение SQL, рассмотрим следующую короткую процедуру:

```
EODA@ORA12CR1> create or replace procedure inj( p_date in date )
2  as
3      l_username  all_users.username%type;
4      c           sys_refcursor;
5      l_query      varchar2(4000);
6  begin
7      l_query := '
8      select username
9      from all_users
10     where created = ''' || p_date || ''';
11
12     dbms_output.put_line( l_query );
13     open c for l_query;
14
15     for i in 1 .. 5
16     loop
17         fetch c into l_username;
18         exit when c%notfound;
19         dbms_output.put_line( l_username || '.....' );
20     end loop;
21     close c;
22 end;
23 /
```

Procedure created.

Процедура создана.

На заметку! Этот код выводит только максимум пять записей. Он был разработан для выполнения в “пустой” схеме. Схема с множеством существующих таблиц может приводить к разнообразным эффектам, которые изменяют результат, показанный ниже. Один из эффектов может выражаться в том, что вы не увидите таблицу, которую я пытаюсь показать в примере — это может быть связано с выводом всего лишь пяти записей. Другой эффект может проявиться как ошибка в числе или значении — по причине длинного имени таблицы. Ни один из этих фактов не делает пример недействительным; желающий похитить данные все это может обойти.

Сейчас большинство известных мне разработчиков посмотрят на этот код и скажут, что он безопасен с точки зрения внедрения SQL. Они обоснуют это тем, что входные данные для процедуры должны быть переменной Oracle типа DATE — 7-байтовым двоичным форматом, представляющим век, год, месяц, день, час, минуту и секунду. Не существует способа, которым переменная DATE могла бы изменить смысл приведенного оператора SQL. Как выяснилось, это очень большая ошибка. Код уязвим в отношении внедрения — модификации во время выполнения — со стороны любого, кто знает, каким образом это сделать (вполне очевидно, что та-

кие люди есть). Если выполнить процедуру так, как разработчик “ожидает”, вот что можно увидеть:

```
EODA@ORA12CR1> exec inj( sysdate )
      select *
      from all_users
      where created = '12-MAR-14'

PL/SQL procedure successfully completed.
Процедура PL/SQL успешно выполнена.
```

Полученный результат показывает, что оператор SQL сконструирован безопасным — вполне ожидаемо. Тогда как можно использовать эту процедуру ненадлежащим образом? Давайте представим, что вы подключили к проекту нового разработчика — разработчика-злоумышленника. Все разработчики имеют доступ для выполнения этой процедуры с целью просмотра пользователей, зарегистрированных в базе данных сегодня, но не имеют доступа к любой таблице схемы, к которой принадлежит данная процедура. Они не знают, какие таблицы имеются в схеме — специалисты из команды обеспечения безопасности решили, что подход вида “безопасность через незнание” вполне эффективен, и потому решили запретить всем публикацию имен таблиц. Таким образом, разработчики не знают о существовании следующей конкретной таблицы:

```
EODA@ORA12CR1> create table user_pw
2  ( uname varchar2(30) primary key,
3    pw  varchar2(30)
4  );
Table created.
Таблица создана.
```

```
EODA@ORA12CR1> insert into user_pw
2  ( uname, pw )
3  values ( 'TKYTE', 'TOP SECRET' );
1 row created.
1 строка создана.
```

```
EODA@ORA12CR1> commit;
Commit complete.
Фиксация выполнена.
```

Таблица USER_PW выглядит довольно важной, но вспомните, что пользователи не знают о ее существовании. Однако они (пользователи с минимальными привилегиями) имеют доступ к процедуре INJ:

```
EODA@ORA12CR1> create user devacct identified by foobar;
User created.
Пользователь создан.

EODA@ORA12CR1> grant create session to devacct;
Grant succeeded.
Выдано успешно.

EODA@ORA12CR1> grant execute on inj to devacct;
Grant succeeded.
```

Злоумышленник из числа разработчиков или пользователей может просто запустить следующий код:

```
EODA@ORA12CR1> connect devacct/foobar;
Connected.
Подключено.
```

```
DEVACCT@ORA12CR1> alter session set
  2  nls_date_format = ''''union select tname from tab--''';
Session altered.
```

Сеанс изменен.

```
DEVACCT@ORA12CR1> exec eoda.inj( sysdate )
```

```
select username
  from all_users
 where created =
''union select tname from tab--'
USER_PW.....
```

PL/SQL procedure successfully completed.

Процедура PL/SQL успешно выполнена.

В приведенном выше коде оператор select выполняет оператор (не возвращающий строк):

```
select username from all_users where created = ''
```

И объединяет его с оператором:

```
select tname from tab
```

Взгляните на последнюю конструкцию, т.е. -- '. В SQL*Plus двумя дефисами, идущими подряд, обозначается комментарий; в итоге последняя одиночная кавычка попадает в комментарий, который необходим для того, чтобы оператор был синтаксически правильным.

Здесь нас интересует NLS_DATE_FORMAT — большинство даже не подозревает о возможности включения символьных строковых литералов с типом NLS_DATE_FORMAT. (К сожалению, многие даже не знают, что формат даты можно изменять и без этого “трюка”. Им также не известна возможность изменять сеанс (для установки NLS_DATE_FORMAT), не располагая привилегией ALTER SESSION!) Злоумышленник здесь может воспользоваться *вашим* кодом для опроса таблицы, которую вы не намеревались запрашивать, *с применением вашего набора привилегий*. Представление словаря TAB ограничивает его доступом к набору таблиц, которые может видеть текущая схема. Когда пользователи запускают эту процедуру, владельцем данной процедуры является текущая схема, используемая для авторизации (т.е. вы, а не пользователи). Теперь они могут видеть, какие таблицы находятся в схеме. Злоумышленники видят таблицу USER_PW и проявляют к ней вполне естественный интерес. Затем они пытаются обратиться к этой таблице:

```
DEVACCT@ORA12CR1> select * from eoda.user_pw;
select * from eoda.user_pw
*
```

ERROR at line 1:

ORA-00942: table or view does not exist

ОШИБКА в строке 1:

ORA-00942: таблица или представление не существует

Злоумышленник не может получить доступ к таблице непосредственно; ему не хватает привилегии SELECT на этой таблице. Но не переживайте — есть и другой путь. Пользователь желает узнать о столбцах таблицы. Вот один способ посмотреть структуру таблицы:

```
DEVACCT@ORA12CR1> alter session set
  2 nls_date_format = ''''union select tname||'/'||cname from col--'';
Session altered.
Сеанс изменен.

DEVACCT@ORA12CR1> exec eoda.inj( sysdate )

select username
  from all_users
 where created =
''union select tname||'/'||cname from col--'
USER_PW/PW.....
USER_PW/UNAME.....

PL/SQL procedure successfully completed.
Процедура PL/SQL успешно выполнена.
```

Итак, теперь известны также имена столбцов. Зная имена таблиц и столбцов таблиц в схеме, можно еще раз изменить NLS_DATE_FORMAT для опроса этой таблицы, а не таблиц словаря данных. Вот что злоумышленник может предпринять далее:

```
DEVACCT@ORA12CR1> alter session set
  2 nls_date_format = ''''union select uname||'/'||pw from user_pw--'';
Session altered.
Сеанс изменен.

DEVACCT@ORA12CR1> exec eoda.inj( sysdate )

select username
  from all_users
 where created =
''union select uname||'/'||pw from user_pw--'
TKYTE/TOP SECRET.....

PL/SQL procedure successfully completed.
Процедура PL/SQL успешно выполнена.
```

И вот, пожалуйста — злоумышленник из числа разработчиков или пользователей теперь обладает вашей ценной информацией об именах пользователей и паролях. Продолжим наши рассуждения: а что, если данный разработчик располагает привилегией CREATE PROCEDURE? Вполне обосновано предположить, что это так (в конце концов, он ведь разработчик). Мог ли он продвинуться еще дальше в рассматриваемом примере? Безусловно. Эта невинно выглядящая хранимая процедура, как минимум, предоставляет гарантированный доступ по чтению ко всему, к чему имеет доступ схема EODA; и если учетная запись, эксплуатирующая такую ошибку, обладает привилегией CREATE PROCEDURE, то хранимая процедура позволит злоумышленнику выполнить любую команду, которую могла бы выполнить схема EODA! Чтобы удостовериться в этом, мы выдадим схеме привилегию CREATE PROCEDURE:

```
DEVACCT@ORA12CR1> connect eoda/foo
Connected.
Подключено.
```

```
EODA@ORA12CR1> grant create procedure to devacct;  
Grant succeeded.  
Выдано успешно.  
  
EODA@ORA12CR1> connect devacct/foobar;  
Connected.  
Подключено.
```

На заметку! В настоящем примере предполагается, что пользователю EODA была назначена роль администратора базы данных (DBA) с конструкцией WITH ADMIN OPTION.

Затем мы от имени разработчика создадим функцию, которая выдает привилегии DBA. С этой функцией связаны два важных факта: она является подпрограммой с правами вызывающего, т.е. будет выполняться с привилегиями, выданными запустившему ее пользователю, и она же является подпрограммой с прагмой `autonomous_transaction`, что означает создание подтранзакции, которая фиксируется или откатывается перед возвратом из подпрограммы, тем самым делая возможным ее вызов из SQL-кода. Вот эта функция:

```
DEVACCT@ORA12CR1> create or replace function foo  
2   return varchar2  
3   authid CURRENT_USER  
4   as  
5       pragma autonomous_transaction;  
6   begin  
7       execute immediate 'grant dba to devacct';  
8       return null;  
9   end;  
10  /
```

Function created.

Функция создана.

Теперь осталось лишь обманом заставить EODA (администратора базы данных, который может предоставить роль DBA другим пользователям) выполнить показанную функцию. Учитывая то, что уже было сделано для использования дефекта внедрения SQL, это легко. Мы установим `NLS_DATE_FORMAT`, чтобы включить ссылку на нашу функцию, и выдадим пользователю EODA привилегию на ее выполнение:

```
DEVACCT@ORA12CR1> alter session set  
2   nls_date_format = ''''union select devacct.foo from dual--''';  
Session altered.  
Сеанс изменен.
```

```
DEVACCT@ORA12CR1> grant execute on foo to eoda;  
Grant succeeded.  
Выдано успешно.
```

Цель достигнута! У нас есть администратор базы данных:

```
DEVACCT@ORA12CR1> select * from session_roles;  
no rows selected  
строки не выбраны  
  
DEVACCT@ORA12CR1> exec eoda.inj( sysdate )
```

```

        select username
        from all_users
        where created =
'union select devacct.foo from dual--'
.....
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно выполнена.
DEVACCT@ORA12CR1> connect devacct/foobar
Connected.
Подключено.
DEVACCT@ORA12CR1> select * from session_roles;
ROLE
-----
DBA
SELECT_CATALOG_ROLE
...
XS_RESOURCE
OLAP_DBA
24 rows selected.
24 строки выбрано.

```

Совет. Чтобы просмотреть, какие роли выданы другим ролям, запросите представление `ROLE_ROLE_PRIVS`.

Как же защититься от этого? Нужно использовать переменные привязки, например:

```

EODA@ORA12CR1> create or replace procedure NOT_inj( p_date in date )
2 as
3     l_username  all_users.username%type;
4     c           sys_refcursor;
5     l_query     varchar2(4000);
6 begin
7     l_query := '
8     select username
9     from all_users
10    where created = :x';
11
12    dbms_output.put_line( l_query );
13    open c for l_query USING P_DATE;
14
15    for i in 1 .. 5
16    loop
17        fetch c into l_username;
18        exit when c%notfound;
19        dbms_output.put_line( l_username || '.....' );
20    end loop;
21    close c;
22 end;
23 /

```

Procedure created.
Процедура создана.

```
EOADA@ORA12CR1> exec NOT_inj(sysdate)

select username
  from all_users
 where created = :x

PL/SQL procedure successfully completed.
Процедура PL/SQL успешно выполнена.
```

Налицо простой и очевидный факт: в случае применения переменных привязки вы не станете жертвой внедрения SQL. Если вы *не* используете переменные привязки, то должны тщательно проинспектировать каждую строку кода с позиции “злого гения” (знающего *абсолютно все* тонкости Oracle) и оценить возможность совершения атаки на код. Если бы я был уверен в том, что 99,9999% моего кода безопасны к атакам внедрением SQL, и пришлось бы переживать только об оставшихся 0,0001% (где по каким-то причинам применение переменных привязки невозможно), то я спал бы гораздо спокойнее, чем в ситуации, когда пришлось бы волноваться по поводу того, что все 100% кода незащищены к внедрению SQL.

Во всяком случае, в конкретном проекте, с описания которого был начат данный раздел, переписывание существующего кода для использования переменных привязки было единственно возможным выходом. Полученный в результате код функционирует на порядок быстрее и во много раз увеличивает количество параллельно работающих пользователей, которые система в состоянии поддерживать. К тому же код стал более безопасным — не пришлось пересматривать всю кодовую базу на предмет возможных проблем с внедрением SQL. Однако эта безопасность досталась дорогой ценой в смысле затрат времени и усилий, потому что моему клиенту пришлось писать код системы, после чего *переписывать код заново*. Дело не в том, что применение переменных привязки трудно или чревато ошибками, а в том, что они не использовали такие переменные изначально, а потому вынуждены были пересмотреть почти весь код и модифицировать его. Моему клиенту не пришлось бы платить эту цену, если бы разработчики понимали жизненную важность переменных привязки с самого первого дня.

Управление параллельной обработкой

Управление параллельной обработкой является одной из тех областей, в которых базы данных различаются. Именно в этой области базы данных отличаются от файловой системы и друг от друга. Исключительно важно, чтобы приложение базы данных корректно работало в условиях параллельного доступа, однако о необходимости проверки выполнения этого требования часто забывают. Технологии, которые прекрасно работают, когда все происходит последовательно, не обязательно будут работать хорошо, когда все пользователи делают что-то одновременно. Отсутствие четкого представления о том, как механизмы управления параллельной обработкой реализованы в конкретной базе данных, приведет к следующим проблемам:

- нарушение целостности данных;
- снижение скорости работы приложений даже при небольшом количестве пользователей;
- уменьшение возможностей масштабирования для поддержки большего числа пользователей.

Обратите внимание, что здесь не говорится “может...” или “вы рискуете...”. Напротив — все это неизбежно *произойдет* при отсутствии соответствующей реализации управления параллельной обработкой, причем даже без вашего ведома. Отсутствие правильно реализованного управления параллельной обработкой приведет к нарушению целостности базы данных, поскольку какие-то функции, работающие в условиях изоляции, будут работать не так, как ожидалось, в ситуации с несколькими пользователями. Приложение будет работать медленнее, чем должно, потому что ему придется ожидать доступа к данным. Возможность масштабирования будет утеряна из-за проблем, связанных с блокировкой и соперничеством. По мере того, как очереди для доступа к ресурсам будут удлиняться, время ожидания будет становиться все больше и больше.

В качестве аналогии можно привести пробку на пункте взимания дорожного сбора. Если автомобили подъезжают по очереди, друг за другом, движение не стопорится. Если же множество автомобилей подъезжает одновременно, начинается образовываться пробка. Более того, время ожидания увеличивается совсем не пропорционально количеству автомобилей на контрольном пункте. С определенного момента значительное дополнительное время начинает тратиться на “разбирательства” между водителями, ожидающими своей очереди, и на их обслуживание (в базе данных этому соответствует переключение контекста).

Проблемы параллельной обработки наиболее трудны для отслеживания. Эта задача аналогична выполнению отладки в многопоточной программе. Программа может нормально работать в управляемой искусственной среде отладчика, но оказаться совершенно неработоспособной в реальной среде. Например, в условиях состязаний может случиться так, что два потока одновременно изменяют одну и ту же структуру данных. Подобные ошибки крайне трудно отслеживать и исправлять. Если проводить тестирование приложения только в изоляции, а затем развернуть его для десятков параллельно работающих пользователей, скорее всего, это приведет к проявлению множества скрытых проблем параллельной обработки.

В следующих двух разделах будут продемонстрированы два небольших примера того, как недостаточное понимание вопросов параллельной обработки может разрушить данные или привести к снижению производительности и масштабируемости.

Реализация блокировки

База данных использует *блокировки* (lock) для обеспечения того, что изменение указанного фрагмента данных в каждый конкретный момент времени может совершать максимум одна транзакция. По существу блокировки — это механизмы, которые разрешают параллельную обработку; без применения какой-либо модели блокировки для предотвращения, например, одновременных обновлений одной и той же строки, многопользовательский доступ к базе данных был бы невозможным. Однако при злоупотреблении или неправильном использовании блокировки могут в действительности только мешать параллельной обработке. Если программист или сама база данных блокирует данные без необходимости, то меньшее число людей сможет параллельно выполнять операции. Поэтому понимание сущности блокировки и особенностей ее работы в базе данных исключительно важно при построении масштабируемого и корректно функционирующего приложения.

Очень важно также понимать, что каждая база данных реализует блокировку по-своему. В одних блокировка осуществляется на уровне страниц, в других — на

уровне строк. Одни базы данных распространяют блокировки с уровня строк на уровень страниц, другие этого не делают. В одних базах данных блокировки чтения применяются, в других — нет. В некоторых базах данных последовательные транзакции реализованы через блокирование, а в других — посредством согласованных по чтению представлений данных (без блокировок). Все эти мелкие различия могут приводить к огромным проблемам с производительностью или явным программным ошибкам, если не понимать особенностей их работы.

Ниже описаны принципы политики блокировки Oracle.

- Oracle блокирует данные на уровне строк при модификации данных. Блокировка не распространяется на уровни блока или таблицы.
- Oracle никогда не блокирует данные, чтобы всего лишь прочитать их. Простые операции чтения не помещают какие-либо блокировки на строки данных.
- Процесс, записывающий данные, не блокирует процесс, который читает данные. Позвольте повторить: *операции чтения не блокируются операциями записи*. Это фундаментальным образом отличается от многих других баз данных, в которых операции чтения блокируются операциями записи. Хотя такая характеристика кажется исключительно положительной (в основном, так оно и есть), попытка наложения ограничений целостности посредством логики приложения при недостаточном понимании этой концепции, *скорее всего, приведет к получению неправильной реализации*.
- Процесс, записывающий данные, блокируется только в том случае, если другой процесс, выполняющий запись, уже заблокировал строку, к которой обращается этот процесс. Процесс, читающий данные, никогда не блокирует процесс, который выполняет запись.

Указанные обстоятельства должны учитываться при разворачивании приложения. При этом необходимо осознавать, что описанная политика уникальна для Oracle — каждая база данных обладает собственными особенностями реализации блокировки. Даже приводя все свои приложения к “наименьшему общему знаменателю SQL”, следует учитывать, что модели блокировки и параллельной обработки, используемые каждым поставщиком баз данных, определяют различия в поведении приложения. Разработчик, который не понимает, каким образом база данных реализует параллельную обработку, неизбежно столкнется с проблемами целостности данных. (Это особенно часто происходит при переходе от другой базы данных к Oracle и наоборот, а также из-за пренебрежения учета в приложениях отличий между механизмами параллельной обработки.)

Предотвращение потерянных обновлений

Один из побочных эффектов применения неблокирующего подхода Oracle состоит в том, что в действительности необходимо обеспечить, чтобы во время выполнения каких-либо действий самим разработчиком доступ к строке предоставлялся не более чем одному пользователю.

Разработчик демонстрировал мне созданную им недавно программу резервирования ресурсов (конференц-залов, проекторов и т.п.), которая находилась в процессе разворачивания. В этом приложении было реализовано бизнес-правило для предотвращения выделения ресурса более чем одному лицу в любой заданный период

времени. То есть приложение содержало код, который специально проверял, не зарезервирован ли данный временной интервал за другим пользователем (по крайней мере, разработчик полагал, что это так). Код запрашивал таблицу SCHEDULES и при отсутствии в ней строк с резервированием, перекрывающим интересующий временной интервал, вставлял новую строку. Таким образом, в основном разработчик имел дело с двумя таблицами:

```

EODA@ORA12CR1> create table resources
 2 ( resource_name varchar2(25) primary key,
 3   other_data    varchar2(25)
 4 );

Table created.
Таблица создана.

EODA@ORA12CR1> create table schedules
 2 ( resource_name varchar2(25) references resources,
 3   start_time    date,
 4   end_time      date
 5 );

Table created.
Таблица создана.

```

Сразу после вставки записи о бронировании помещения в таблицу SCHEDULES, но перед фиксацией приложение выполняло запрос:

```

EODA@ORA12CR1> select count(*)
 2   from schedules
 3  where resource_name = :resource_name
 4         and (start_time < :new_end_time)
 5         AND (end_time > :new_start_time)
 6 /

```

Все казалось простым и “пуленепробиваемым” (во всяком случае, разработчику); если результатом подсчета был 0, то помещение поступало в ваше распоряжение. Если результат подсчета больше 0, то помещение не может быть зарезервировано на заявленный период времени. Выяснив, в чем состояла логика разработчика, я выполнил очень простой тест, чтобы показать ему ошибку, которая должна была проявиться в процессе эксплуатации приложения — ошибку, которую было бы невероятно трудно отследить и диагностировать. Многие были бы убеждены, что это *должна* быть программная ошибка базы данных.

Все что я сделал — это попросил кого-то воспользоваться соседним терминалом. Оба пользователя (разработчик и второе лицо) открыли один и тот же экран приложения, и “на счет три” каждый из них щелкнул на кнопке выполнения и попытался зарезервировать одно и то же помещение на почти перекрывающийся период времени. Обоим пользователям удалось выполнить резервирование. Логика, которая прекрасно работала в условиях изоляции, дала сбой в многопользовательской среде. В этом случае проблема частично была вызвана неблокирующими операциями чтения Oracle. Один из сеансов никак не блокировал другой. Оба сеанса просто запустили запрос и выполнили действия по резервированию помещения. Оба сеанса смогли выдать запрос для просмотра информации о резервировании, даже если один из них уже начал модифицировать таблицу SCHEDULES (изменение не будет видно другому сеансу вплоть до выполнения операции фиксации, ко времени которой уже слишком

поздно). Поскольку казалось, что пользователи никогда не будут пытаться изменять одну и ту же строку в таблице `SCHEDULES`, то их блокировка не производилась и, следовательно, бизнес-правило не могло вступить в силу, как было задумано.

Это вызвало удивление у разработчика, написавшего к этому моменту немало приложений баз данных, т.к. его опыт был связан с базой данных, поддерживающей блокировки чтений. То есть сеанс, читающий данные, блокировался бы сеансом, записывающим данные, а сеанс, записывающий данные, был бы заблокирован параллельным чтением этих данных. В его мире одна из двух описанных выше транзакций заблокировала бы другую — или, возможно, в приложении возникла бы взаимоблокировка. Однако в итоге транзакция все равно потерпела бы неудачу.

Таким образом, разработчику требовался метод принудительного применения бизнес-правила в многопользовательской среде — способ обеспечения того, чтобы в каждый момент времени только одно лицо могло делать резервирование конкретного ресурса. В этом случае решение заключалось в небольшой сериализации. Кроме выполнения приведенного ранее запроса `count (*)` разработчик сначала запускал следующий запрос:

```
select * from resources where resource_name = :resource_name FOR UPDATE;
```

Здесь производится блокирование ресурса (помещения), предназначенного для резервирования, непосредственно *перед* его резервированием — иначе говоря, до запроса этого ресурса из таблицы `SCHEDULES`. За счет блокирования ресурса, который пользователь пытается зарезервировать, разработчик гарантирует, что никто другой в то же самое время не сможет изменить график эксплуатации этого ресурса. Любым пользователям, желающим выполнить оператор `SELECT FOR UPDATE` для данного ресурса, придется ожидать, пока не будет осуществлена фиксация транзакции, после чего они получают возможность увидеть график использования ресурса. Шансы перекрытия графиков исчезали. Разработчики обязаны понимать, что в многопользовательской среде они должны иногда применять приемы, аналогичные используемым при многопоточном программировании. В рассматриваемом случае конструкция `FOR UPDATE` действует подобно семафору. Она обеспечивает последовательный доступ к конкретной строке таблицы `RESOURCES`, не позволяя каким-то двум пользователям зарезервировать ресурс одновременно.

Применение подхода `FOR UPDATE` по-прежнему обеспечивает высокую степень параллелизма, т.к. потенциально могут резервироваться тысячи ресурсов. Это тот редкий случай, когда устанавливается ручная блокировка данных, которые в действительности мы обновлять не собираемся. Вы должны уметь распознавать ситуации, когда необходимо блокировать ручную, и что вероятно еще важнее — когда этого делать не следует (позже будет приведен соответствующий пример). Более того, конструкция `FOR UPDATE` не блокирует чтение данных о ресурсе другими пользователями, как это бывает в других базах данных. В результате такой подход будет очень хорошо масштабироваться.

Проблемы, подобные описанной в этом разделе, имеют крупные последствия при переносе приложения из одной базы данных в другую (я вернусь к этой теме далее в главе), время от времени создавая разработчикам разнообразные препятствия. Например, если вы обладаете опытом работы с другими базами данных, в которых процессы записи данных блокируют процессы чтения и наоборот, то можете привыкнуть полагаться на этот факт для предотвращения проблем с целостностью

данных. *Отсутствие* параллелизма является одним из способов защиты от проблем подобного рода, и именно так работают многие базы данных, отличные от Oracle. В базе данных Oracle правила параллельной обработки обладают наивысшим приоритетом, и следует помнить, что это обуславливает совершенно иное поведение приложений (в противном случае неприятные последствия неизбежны).

Мне приходилось сталкиваться с ситуациями, когда разработчики, даже ознакомившись с подобным примером, поднимали на смех идею о том, что им необходимо понимать, каким образом все это работает. Они говорили: “Мы просто помечаем флажок ‘транзакционное’ для нашего приложения Hibernate и оно само заботится обо всем, что связано с транзакциями; нам незачем знать всю эту кухню”. Я спрашивал их: “Так что, Hibernate будет генерировать отличающийся код для SQL Server, DB2 и Oracle — совершенно разный код, разное количество операторов SQL, разную логику?”. Они отвечали — нет, но приложение будет транзакционным. Возникла путаница. Транзакционное в данном контексте означает просто поддержку фиксации и отката, а не то, что код транзакционно согласован (читай: что код корректен). Независимо от инструмента или платформы, используемой для доступа к базе данных, знание средств управления параллельным доступом жизненно важно, если вы не хотите повредить свои данные.

На протяжении 99% времени блокировка совершенно прозрачна, и о ней можно не беспокоиться. Но именно оставшийся один процент случаев нужно уметь распознавать. Для решения этой проблемы не существует какого-то простого контрольного перечня типа “чтобы сделать это, необходимо выполнить то”. Все зависит от понимания поведения приложения в многопользовательской среде и в базе данных.

Мы рассмотрим эту тему значительно подробнее в главах, посвященных блокированию и управлению параллельной обработкой. В них вы узнаете, что применение ограничений целостности данных такого типа, как описанный в этом разделе, когда вы должны реализовать правило, распространяющееся на несколько строк в одной таблице или на две и более таблиц (вроде ограничения ссылочной целостности), всегда требует особого внимания. Кроме того, в многопользовательской среде, скорее всего, придется использовать ручное блокирование или какую-то другую технологию для поддержания целостности.

Многоверсионность

Эта тема очень тесно связана с управлением параллелизмом, поскольку многоверсионность служит основой механизма управления параллельной обработкой Oracle. База данных Oracle действует на основе многоверсионной, согласованной по чтению модели параллелизма. В главе 7 мы раскроем технические аспекты этой модели более подробно, но по существу она представляет собой механизм, посредством которого Oracle обеспечивает следующие концепции.

- **Согласованные по чтению запросы.** Запросы, которые создают непротиворечивые результаты для конкретного момента времени.
- **Неблокирующие запросы.** Запросы никогда не блокируются процессом, записывающим данные, как это имеет место в других базах данных.

В базе данных Oracle эти две концепции очень важны. В основном термин *многоверсионность* описывает способность Oracle одновременно поддерживать несколько версий данных, извлеченных из базы (начиная с версии 3.0, вышедшей в 1983 году!).

Понятие *согласованность чтения* отражает тот факт, что запрос в Oracle будет возвращать результаты из согласованного момента времени. Каждый блок, используемый запросом, будет в точности соответствовать этому моменту времени, даже если данные были модифицированы или заблокированы, пока выполнялся запрос (это было действительным, начиная с версии 4.0, которая появилась еще в 1984 году!). Если вы понимаете, как совместно работают многоверсионность и согласованность чтения, то всегда поймете ответы, полученные из базы данных. Прежде чем приступить к более подробному рассмотрению реализации многоверсионности в Oracle, давайте ознакомимся с ее работой на простом примере:

```
EODA@ORA12CR1> create table t
2 as
3 select username, created
4 from all_users
5 /
```

Table created.

Таблица создана.

```
EODA@ORA12CR1> set autoprint off
EODA@ORA12CR1> variable x refcursor;
EODA@ORA12CR1> begin
2 open :x for select * from t;
3 end;
4 /
```

PL/SQL procedure successfully completed.

Процедура PL/SQL успешно выполнена.

```
EODA@ORA12CR1> declare
2 pragma autonomous_transaction;
3 -- вы могли бы также сделать это
4 -- в другом сеансе sqlplus,
5 -- и результат был бы идентичным
6 begin
7 delete from t;
8 commit;
9 end;
10 /
```

PL/SQL procedure successfully completed.

Процедура PL/SQL успешно выполнена.

```
EODA@ORA12CR1> print x
```

| USERNAME | CREATED |
|------------------|-----------|
| DEVACCT | 02-SEP-13 |
| OP\$MELANIE | 17-JUL-13 |
| SCOTT | 03-JUL-13 |
| OP\$TKYTE | 02-SEP-13 |
| APEX_040200 | 28-JUN-13 |
| APEX_PUBLIC_USER | 28-JUN-13 |
| ... | |
| AUDSYS | 28-JUN-13 |
| SYS | 28-JUN-13 |

36 rows selected.

36 строк выбрано.

В этом примере мы создаем тестовую таблицу `T` и загружаем в нее данные из таблицы `ALL_USERS`. Затем мы открываем на этой таблице курсор. *Никакие данные* из курсора не извлекаются: он только открывается и остается в таком состоянии.

На заметку! Имейте в виду, что СУБД Oracle не “отвечает предварительно” на запрос. В случае открытия курсора данные никуда не копируются — представьте себе, сколько времени потребовалось бы для открытия курсора в таблице, содержащей 1 миллиард строк, если это вообще удалось бы сделать. Курсор открывается немедленно и отвечает на запрос по мере перемещения. Другими словами, чтение данных из таблицы выполняется только при их извлечении из курсора.

Затем в этом же сеансе (или, возможно, в другом — подход все равно сработал бы) мы перешли к удалению всех данных из таблицы. Мы даже выполнили фиксацию (`COMMIT`) этого удаления. Строки исчезли — но так ли это на самом деле? Фактически их можно извлечь через курсор (или посредством запроса `FLASHBACK` с применением конструкции `AS OF`). Дело в том, что результирующий набор, возвращенный командой `OPEN`, был предопределен на момент его открытия. Во время открытия мы затронули не единственный блок данных в этой таблице, но ответ уже был “высечен в камне”. У нас нет возможности узнать, каким будет этот ответ, до тех пор, пока не извлечем данные; однако с точки зрения нашего курсора результат является неизменным. Это не значит, что при открытии курсора Oracle копирует предыдущие данные в какое-то другое место; на самом деле данные сохранила для нас команда `DELETE`, поместив их (начальные копии образов строк, существующие до выполнения `DELETE`) в область данных, которая называется *сегментом отмены* или *сегментом отката*.

Ретроспектива

В прошлом решение о моменте времени, когда выполняемые запросы должны быть согласованными, всегда принимала СУБД Oracle. То есть СУБД Oracle делала это так, чтобы любой открытый результирующий набор был текущим относительно одного из следующих двух моментов времени.

- *Момент времени, когда курсор был открыт.* Это стандартное поведение для режима изоляции `READ COMMITTED` (отличия между уровнями транзакций `READ COMMITTED`, `READ ONLY` и `SERIALIZABLE` рассматриваются в главе 7).
- *Момент времени начала транзакции, к которой относится данный запрос.* Это стандартное поведение для уровней транзакций `READ ONLY` и `SERIALIZABLE`.

Тем не менее, начиная с версии Oracle9i, мы получили в свое распоряжение средство *ретроспективного запроса*, с помощью которого можно предложить базе данных Oracle выполнить запрос “на момент” (разумеется, с определенными разумными ограничениями в отношении промежутка времени, на который можно возвратиться в прошлое). Таким образом, вы можете “видеть” согласованность чтения и многоверсионность еще более наглядно.

На заметку! Архив ретроспективных данных, применяемый для обслуживания ретроспективных запросов (отстоящих на месяцы и годы в прошлое), который доступен в Oracle 11g Release 1 и последующих версиях, для производства версии данных, существовавших в

базе на определенный момент времени в прошлом, не использует согласованность чтения и многоверсионность. Вместо этого он применяет копии прежних образов записей, которые были помещены в архив. Мы возвратимся к теме архивных ретроспективных данных в одной из последующих глав. Также обратите внимание, что архив ретроспективных данных является *функциональным средством* базы данных, начиная с версии 11.2.0.4 и далее. Ранее это была возможность для базы данных, доступная за отдельную плату; теперь это средство для всех, используемое без дополнительных лицензионных затрат.

Рассмотрим следующий пример. Начнем с получения значения SCN (System Change Number (номер системного изменения) или System Commit Number (номер системной фиксации); эти два термина взаимозаменяемы). Номер SCN представляет собой значение внутренних часов Oracle: каждый раз, когда происходит фиксация, показание этих часов увеличивается (инкрементируется). Можно было бы также применять дату или отметку времени, но значение SCN является легко доступным и очень точным:

```
SCOTT@ORA12CR1> variable scn number
SCOTT@ORA12CR1> exec :scn := dbms_flashback.get_system_change_number;
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно выполнена.
SCOTT@ORA12CR1> print scn

      SCN
-----
13646156
```

На заметку! В вашей системе доступ к пакету DBMS_FLASHBACK может быть ограничен. Я выдал право на выполнение этого пакета в своей базе данных пользователю SCOTT; вы можете сделать то же самое.

Имея номер SCN, мы можем потребовать от Oracle предоставить данные на момент времени, соответствующий значению SCN; вместо SCN можно было бы указать дату и отметку времени. Мы хотим иметь возможность запросить базу данных Oracle позже и выяснить, что находилось в таблице в этот конкретный момент времени. Но давайте посмотрим, какую информацию содержит таблица EMP прямо сейчас:

```
SCOTT@ORA12CR1> select count(*) from emp;

COUNT(*)
-----
      14
```

Теперь удалим всю информацию и удостоверимся, что она “исчезла”:

```
SCOTT@ORA12CR1> delete from emp;
14 rows deleted.
14 строк удалено.
SCOTT@ORA12CR1> select count(*) from emp;

COUNT(*)
-----
       0
```

```
SCOTT@ORA12CR1> commit;
Commit complete.
Фиксация выполнена.
```

Однако, используя ретроспективный запрос с конструкцией AS OF SCN или AS OF TIMESTAMP, можно также потребовать от Oracle отобразить содержимое таблицы на желаемый момент времени:

```
SCOTT@ORA12CR1> select count(*),
2         :scn then_scn,
3         dbms_flashback.get_system_change_number now_scn
4   from emp as of scn :scn;

COUNT(*)  THEN_SCN  NOW_SCN
-----
14         13646156  13646157
```

Наконец, в Oracle 10g и последующих версиях доступна команда flashback, которая позволяет посредством лежащей в основе технологии многоверсионности возвращать объекты в состояние, в котором они пребывали в какой-то момент времени в прошлом. В рассматриваемом случае мы можем привести таблицу EMP к виду, который она имела перед удалением всей информации (в качестве части процесса нам понадобится разрешить перемещение строк (row movement), что позволит изменять назначенные идентификаторы строк (rowid) — необходимое условие для выполнения ретроспективного отката таблицы):

```
SCOTT@ORA12CR1> alter table emp enable row movement;
Table altered.
Таблица изменена.

SCOTT@ORA12CR1> flashback table emp to scn :scn;
Flashback complete.
Фиксация завершена.

SCOTT@ORA12CR1> select cnt_now, cnt_then,
2         :scn then_scn,
3         dbms_flashback.get_system_change_number now_scn
4   from (select count(*) cnt_now from emp),
5        (select count(*) cnt_then from emp as of scn :scn)
6   /

CNT_NOW  CNT_THEN  THEN_SCN  NOW_SCN
-----
14       14     13646156  13646786
```

Вот что собой представляют согласованность чтения и многоверсионность. Если вы не понимаете, как работает многоверсионность Oracle и что она подразумевает, то не сможете воспользоваться всеми преимуществами Oracle или писать корректные приложения в Oracle (гарантирующие целостность данных).

На заметку! Ретроспективный откат таблицы требует наличия версии Enterprise Edition базы данных Oracle.

Согласованность чтения и неблокирующие чтения

А теперь рассмотрим последствия многоверсионности: согласованные по чтению запросы и неблокирующие чтения. Приведенный ниже код может вызвать удивление у тех, кто не знаком с многоверсионностью. Для простоты предположим, что считываемая таблица содержит по одной строке для каждого блока базы данных (наименьшей единицы хранения в базе данных), а в этом примере выполняется полный просмотр таблицы.

Мы будем запрашивать простую таблицу `ACCOUNTS`, которая хранит балансы банковских счетов и имеет очень простую структуру:

```
create table accounts
( account_number number primary key,
  account_balance number
);
```

В реальном приложении таблица `ACCOUNTS` должна была бы содержать сотни тысяч строк, но для простоты мы рассмотрим таблицу с четырьмя строками, как показано в табл. 1.1 (в главе 7 мы проанализируем этот пример более подробно).

Таблица 1.1. Содержимое таблицы `ACCOUNTS`

| Строка | Номер счета (<code>account_number</code>) | Баланс счета (<code>account_balance</code>) |
|--------|--|--|
| 1 | 123 | \$500.00 |
| 2 | 234 | \$250.00 |
| 3 | 345 | \$400.00 |
| 4 | 456 | \$100.00 |

В конце рабочего дня мы хотим получить отчет о наличии денежных средств в банке. Это предельно простой запрос:

```
select sum(account_balance) from accounts;
```

Ответ в этом примере, конечно же, очевиден: \$1250. Однако что произойдет, если мы успешно прочитали строку 1, а во время чтения строк 2 и 3 банкомат генерирует транзакции в отношении этой таблицы и переводит сумму \$400 со счета 123 на счет 456? Запрос обнаружит \$500 в строке 4 и возвратит сумму, равную \$1650, не так ли? Естественно, подобного следует избегать, потому что ответ был бы ошибочным — ни в один из моментов времени в столбце баланса счетов такая денежная сумма не присутствует. В Oracle описанные ситуации не возникают благодаря согласованности чтения. Методы, используемые в Oracle, отличаются от применяемых в большинстве остальных баз данных, и вы должны понимать эти отличия.

Во многих других базах данных для получения “согласованного” и “корректного” ответа на этот запрос придется либо заблокировать целую таблицу, пока сумма вычисляется, либо заблокировать строки при их чтении. Это препятствует изменению ответа во время его получения. Если вы заблокируете таблицу заранее, то ответ будет соответствовать состоянию базы данных на момент, когда запрос начал выполняться. Если вы будете блокировать данные по мере их считывания (метод, обычно называемый *блокировкой разделяемого чтения*, который предотвращает обновления,

но не препятствует доступу к данным со стороны других процессов, выполняющих чтение), то получите ответ, соответствующий состоянию базы данных на момент завершения запроса. Оба метода в значительной степени снижают степень параллелизма. Блокировка таблицы препятствует любым обновлениям целой таблицы на протяжении выполнения запроса (для таблицы, состоящей из четырех строк, этот промежуток времени оказывается очень коротким, но для таблиц, которые содержат сотни тысяч строк, он мог бы занять несколько минут). Метод “блокировки по ходу” предотвращает обновления данных, которые были прочитаны и уже обработаны, и может приводить к возникновению взаимоблокировок между этим запросом и другими операциями обновления.

Как уже упоминалось ранее, без полного понимания концепции многоверсионности невозможно в полной мере воспользоваться преимуществами, предлагаемыми Oracle. Разберем одну из причин этого. Многоверсионность в Oracle применяется для получения ответа, соответствующего моменту начала запроса, и запрос будет выполнен *без блокирования чего-либо*. (Пока транзакция перевода денежных сумм с одного счета на другой обновляет строки 1 и 4, эти строки будут заблокированы для других процессов записи, но не для процессов чтения данных, таких как запрос `SELECT SUM. . .`) Фактически в Oracle отсутствует блокировка “разделяемого чтения” (распространенный тип блокировки в других базах данных), т.к. необходимость в ней отсутствует. Все, что снижает степень параллелизма, из Oracle было удалено.

Мне на практике встречались случаи, когда отчет, созданный разработчиком, который не понимал возможности многоверсионности Oracle, полностью блокировал всю систему. Причиной было то, что разработчик стремился получить согласованные по чтению (т.е. правильные) результаты запросов. В любой другой базе данных, с которой разработчик имел дело, это требовало блокирования таблиц либо использования оператора `SELECT . . . WITH HOLDLOCK` (механизм SQL Server, предназначенный для блокировки строк в разделяемом режиме по мере их считывания). В итоге разработчик должен был либо блокировать таблицы перед тем, как приступить к запуску построения отчета, либо применять запрос `SELECT . . . FOR UPDATE` (ближайший аналог `WITH HOLDLOCK`). Это приводило к практически полной остановке обработки транзакций — причем совершенно неоправданно.

Каким же образом СУБД Oracle получает правильный непротиворечивый ответ (\$1250) во время чтения без блокировки каких-либо данных — другими словами, без снижения степени параллелизма? Секрет кроется в используемых Oracle транзакционных механизмах. При каждом изменении данных Oracle создает записи в двух разных местах (большинство других баз данных помещают обе записи в одно место; для них сегмент отката (undo) и журнальный сегмент (redo) — просто “транзакционные данные”). Одна запись поступает в журналы повторного выполнения, где Oracle хранит информацию для *повторного выполнения* (redo), или “наката” (roll forward) транзакции. Для операции вставки это будет вставляемая строка. Для операции удаления концептуально это будет сообщение об удалении строки из файла X, блока Y, строкового слота Z и т.д. Еще одной записью является *запись отмены* (undo), помещаемая в сегмент отмены. Если транзакция отказывает и должна быть отменена, то Oracle читает образ “до того” из сегмента отмены и восстанавливает данные. В дополнение к применению этих данных сегмента отмены для отмены транзакций, Oracle использует их для отмены изменений в блоках данных и их восстановления в том состоянии, которое соответствовало моменту начала выполнения запроса.

Это позволяет производить чтение, минуя блокировку, и получать согласованные, правильные ответы, не блокируя какие-либо данные самостоятельно.

Итак, применительно к рассматриваемому примеру, Oracle приходит к правильному ответу, как описано в табл. 1.2.

Таблица 1.2. Многоверсионность в действии

| Момент времени | Запрос | Транзакция перевода денежных сумм с одного счета на другой |
|----------------|--|--|
| T1 | Считывает строку 1; на данный момент баланс = \$500; сумма = \$500 | |
| T2 | | Обновляет строку 1; помещает монопольную блокировку на строку 1, предотвращая другие обновления (но не чтения). Теперь строка 1 содержит \$100 |
| T3 | Считывает строку 2; на данный момент баланс = \$250; сумма = \$750 | |
| T4 | Считывает строку 3; на данный момент баланс = \$400; сумма = \$1150 | |
| T5 | | Обновляет строку 4; помещает монопольную блокировку на строку 4, предотвращая другие обновления (но не чтения). Теперь строка 4 содержит \$500 |
| T6 | Считывает строку 4; обнаруживает, что строка 4 была изменена. В действительности будет выполнен откат блока, чтобы его вид соответствовал моменту времени T1. Из этого блока запрос прочитает значение \$100 | |
| T7 | | Фиксирует транзакцию |
| T8 | Представляет в качестве ответа значение \$1250 | |

В момент времени T6 по существу Oracle выполняет “чтение через” блокировку, которую наша транзакция поместила на строку 4. Именно так реализованы неблокирующие чтения: Oracle только проверяет, изменились ли данные, не беспокоясь о том, заблокированы ли данные в текущий момент (что подразумевает возможность изменения данных). СУБД Oracle просто извлекает старые данные из сегмента отката и переходит к следующему блоку данных.

Рассмотрим еще одну простую демонстрацию многоверсионности. В базе данных доступно несколько версий одного и того же фрагмента информации для различных моментов времени. СУБД Oracle способна использовать эти снимки данных, соответствующие разным моментам времени, для выполнения согласованных по чтению запросов и неблокирующих чтений.

Такое согласованное по чтению представление данных всегда создается на уровне SQL-оператора. Результаты любого одиночного SQL-оператора являются непро-

тиворечивыми относительно момента начала его выполнения. Именно это качество делает результат выполнения оператора, подобного следующей вставке, предсказуемым набором данных:

```
for x in (select * from t)
loop
    insert into t values (x.username, x.user_id, x.created);
end loop;
```

Результат запроса `SELECT * FROM` предопределен на момент начала его выполнения. Оператор `SELECT` не увидит никаких новых данных, сгенерированных оператором `INSERT`. Вообразите себе, что если бы он делал это — оператор мог бы превратиться в бесконечный цикл. Если бы оператор `SELECT` мог “видеть” новые строки, вставленные в таблицу `T` оператором `INSERT`, то приведенный выше код создавал бы неизвестное количество строк. Скажем, если бы в начале таблица `T` имела 10 строк, то в конце она могла бы содержать 20, 21, 23 или бесконечное число строк. Поведение было бы совершенно непредсказуемым. Согласованное чтение обеспечивается для всех операторов, поэтому оператор `INSERT` вроде показанного ниже также является предсказуемым:

```
insert into t select * from t;
```

Оператор `INSERT` будет работать с согласованным по чтению представлением таблицы `T`. Он не увидит строки, которые только что вставил; вместо этого он вставит только те строки, которые существовали на момент начала выполнения оператора `SELECT`. Определенные базы данных вообще не допускают применения рекурсивных операторов, подобных приведенному, т.к. они не в состоянии определить, сколько строк нужно на самом деле вставить.

Таким образом, если вы привыкли к способу поддержания целостности и параллелизма запросов, принятому в других базах данных, или вам никогда не приходилось иметь дело с такими концепциями (например, из-за полного отсутствия опыта использования баз данных), то теперь важность понимания их работы должна быть ясна. Чтобы максимально задействовать потенциал Oracle и реализовать корректно функционирующий код, *необходимо* хорошо понимать реализацию этих концепций именно в Oracle, а не в других базах данных.

Нужна ли независимость от базы данных?

К настоящему моменту вы уже можете догадаться, о чем пойдет речь в этом разделе. Я ссылаюсь на другие базы данных и на то, что в каждой из них функциональные средства реализованы по-своему. Я убежден, что за исключением некоторых приложений, предназначенных только для чтения данных, построение полностью независимого от базы данных и при этом хорошо масштабируемого приложения — исключительно трудная задача. Фактически это почти невозможно, если только не знать совершенно точно, каким образом работает каждая база данных. А если вы очень хорошо знаете, как именно работает каждая база данных, то должны понимать, что независимость от базы данных — вовсе не то, к чему следует действительно стремиться.

В целях иллюстрации вернемся к ранее рассмотренному примеру резервирования ресурсов (до добавления конструкции `FOR UPDATE`). Предположим, что это прило-

жение было разработано в базе данных, модель блокировки/параллелизма в которой совершенно отличается от модели, принятой в Oracle. Я собираюсь показать, что при переносе приложения из одной базы данных в другую придется удостовериться в том, что она продолжает правильно работать в этих разных средах, и существенно его изменить!

Пусть первоначальное приложение резервирования ресурсов было развернуто в базе данных, в которой применяются блокирующие чтения (т.е. процессы чтения блокируются процессами записи). Учтем также, что бизнес-правило было реализовано с помощью триггера базы данных (*после* выполнения операции INSERT, но перед фиксацией транзакции необходимо убедиться, что таблица содержит только одну строку, соответствующую данному временному промежутку). В системе с блокирующими чтениями из-за наличия недавно вставленных данных справедливо отметить, что операции вставки в эту таблицу должны выполняться последовательно. Первый пользователь вставил бы свой запрос на резервирование “помещения А” с 14:00 до 15:00 в пятницу, а затем выполнил запрос для выявления возможных накладок. Если бы следующий пользователь попытался вставить запрос, пересекающийся с первым, то во время просмотра на предмет наличия накладок этот запрос оказался бы заблокированным (ожидая, пока недавно вставленные данные станут доступными для чтения). Вполне очевидно, что в такой базе данных с блокирующими чтениями наше приложение должно работать успешно, хотя легко могла бы возникнуть *взаимоблокировка* (эта концепция раскрывается в главе 6), если бы мы одновременно вставляли строки и затем пытались читать данные друг друга. Проверка наличия перекрывающихся интервалов для выделенных ресурсов будет выполняться последовательно — и никогда параллельно.

Если мы перенесем это приложение в Oracle и просто предположим, что оно будет вести себя аналогичным образом, то испытаем немалое потрясение. Окажется, что в среде Oracle, в которой производится блокировка на уровне строк и обеспечиваются неблокирующие чтения, приложение ведет себя некорректно. Как мы уже выяснили, для обеспечения последовательного доступа необходимо использовать конструкцию FOR UPDATE. Без этой конструкции два пользователя могли бы зарезервировать один ресурс на то же самое время. Это прямое следствие непонимания особенностей работы применяемой базы данных в многопользовательской среде.

С подобного рода проблемами мне приходилось многократно встречаться при переносе приложения из базы данных А в базу данных Б. Когда приложение, которое безошибочно функционировало в базе данных А, не работает или работает странно в базе данных Б, прежде всего, первой приходит мысль об “ущербности” базы данных Б. Истинная правда заключается в том, что база данных Б всего лишь работает *по-другому*. Ни одна из этих баз данных не является неправильной или “ущербной”; они просто различны. Знание и четкое понимание особенностей функционирования каждой из них чрезвычайно помогает в решении таких проблем. Перенос приложения из Oracle в SQL Server делает актуальными вопросы блокирующих чтений и взаимоблокировки — иначе говоря, проблемы возникают в обоих направлениях.

Например, как-то ко мне обратились за помощью в преобразовании кода Transact-SQL (язык хранимых процедур, разработанный для SQL Server) в код PL/SQL. Разработчик, занимающийся преобразованием, жаловался, что в среде Oracle запросы SQL возвращали “ошибочные” ответы. Запросы выглядели следующим образом:

```

declare
  l_some_variable  varchar2(25);
begin
  if ( some_condition )
  then
    l_some_variable := f( ... );
  end if;

  for x in ( select * from T where x = l_some_variable )
  loop
    ...
  end loop;
end;

```

Цель здесь состояла в нахождении в таблице T всех строк, в которых значение x было равно NULL при несоблюдении определенного условия или конкретному значению, если это условие было соблюдено.

Разработчик утверждал, что в Oracle этот запрос не возвращал никаких данных, если переменная L_SOME_VARIABLE не была установлена в конкретное значение (когда она оставалась равной NULL). В Sybase или SQL Server все было в порядке — запрос находил строки, в которых значение переменной x было равно NULL. Подобная ситуация встречалась мне практически во всех переносах приложений из Sybase или SQL Server в Oracle. Язык SQL рассчитан на работу в контексте трехзначной логики, а в Oracle сравнения с NULL реализованы в соответствии с требованиями стандарта ANSI SQL (где NULL обозначает состояние неизвестной величины и не значение). Согласно этим правилам результат сравнения x с NULL будет ни истинным, ни ложным — он в действительности является *неизвестным*. Это утверждение иллюстрирует следующий фрагмент кода:

```

EODA@ORA12CR1> select * from dual where null=null;
no rows selected
строки не выбраны

```

```

EODA@ORA12CR1> select * from dual where null <> null;
no rows selected
строки не выбраны

```

```

EODA@ORA12CR1> select * from dual where null is null;

```

```

D
-
X

```

Поначалу этот фрагмент может показаться запутанным. Он доказывает, что в Oracle значение NULL не является ни равным, ни не равным NULL. В SQL Server по умолчанию ситуация совершенно иная: в SQL Server и в Sybase значение NULL равно NULL (так принято по умолчанию; в текущих выпусках SQL Server стандартное поведение может быть изменено, чтобы отразить стандарт ANSI). Нельзя считать, что в упомянутых базах данных обработка выполняется *неправильно* — она просто *отличается*. В действительности все эти базы данных совместимы с ANSI (совместимость с ANSI отнюдь не означает стопроцентную поддержку стандарта; за деталями обращайтесь в раздел “Влияние стандартов” далее в главе), но все же они функционируют по-разному. Существуют неоднозначности, проблемы обратной совместимости и другие сложности, которые приходится преодолевать. Например, SQL Server поддерживает

ANSI-метод сравнения с NULL, но не по умолчанию (иначе нарушилась бы работа тысяч унаследованных приложений, построенных на основе этой базы данных).

В рассматриваемом случае одно из возможных решений проблемы предусматривает написание такого запроса:

```
select *  
  from t  
 where ( x = l_some_variable OR (x is null and l_some_variable is NULL  
))
```

Однако это приводит к другой проблеме. В SQL Server такой запрос использовал бы индекс на *x*. В Oracle это может быть не так, поскольку индекс со структурой В-дерева вообще не будет индексировать запись NULL (технологии индексации описаны в главе 11). Следовательно, если требуется искать значения NULL, то индексы со структурой В-дерева не особенно полезны.

На заметку! До тех пор пока хотя бы один столбец индекса со структурой В-дерева определен как NOT NULL, все строки таблицы будут в действительности присутствовать в индексе, так что предикат *x is null* в *where* при извлечении строк может — и будет — использовать индекс.

В этом случае для минимизации влияния на код мы присвоили *x* значение, которое в реальности никогда не может встретиться. Здесь по определению значение *x* было положительным числом, поэтому мы выбрали число -1. Таким образом, запрос приобрел следующий вид:

```
select * from t where nvl(x,-1) = nvl(l_some_variable,-1)
```

Вдобавок мы создали индекс на основе функции:

```
create index t_idx on t( nvl(x,-1) );
```

За счет внесения минимальных изменений мы добились того же самого конечного результата. Ниже перечислены важные выводы, которые должны быть сделаны на основе этого примера.

- Базы данных отличаются друг от друга. Опыт работы с одной из них частично применим в другой, но при этом необходимо быть готовым к ряду как фундаментальных, так и совсем мелких различий.
- Незначительные отличия (вроде обработки значений NULL) могут оказывать столь же крупное влияние, как и фундаментальные отличия (такие как механизмы управления параллельной обработкой).
- Единственным способом предотвращения таких проблем является знание базы данных, особенностей ее работы и реализации функциональных возможностей.

Разработчики часто (обычно несколько раз в день) меня спрашивают, как сделать в базе данных специфичную работу, например, создать временную таблицу в хранимой процедуре. Я никогда не даю прямого ответа на подобные вопросы. Вместо этого я задаю встречный вопрос: *почему* вы хотите сделать это? В ответ я часто слышу: в SQL Server мы создавали временные таблицы в хранимых процедурах, и нам нужно делать это в Oracle. Нечто подобное я и ожидал услышать. В таких ситуациях мой

ответ прост: вам не нужно создавать временные таблицы в хранимых процедурах — вы только думаете, что должны это делать. И действительно, создание таких таблиц в Oracle — очень неудачная затея. Если вы создадите таблицы в хранимой процедуре Oracle, то обнаружите следующие моменты.

- Выполнение DDL-кода снижает возможности масштабирования.
- Постоянное выполнение DDL-кода снижает скорость работы.
- Выполнение DDL-кода приводит к фиксации транзакции.
- Для доступа к таким таблицам во всех хранимых процедурах придется использовать динамический, а не статический SQL-код (из-за того, что эти таблицы на этапе компиляции не существуют).
- Динамический SQL-код в PL/SQL выполняется медленнее и оптимизирован в меньшей степени, чем статический SQL-код.

Итог всего сказанного состоит в том, что вовсе не обязательно поступать в точности так, как в SQL Server (если временная таблица вообще необходима в Oracle). Все задачи должны решаться наиболее эффективным способом, принятым в Oracle. Точно так же при переносе приложения из Oracle в SQL Server не следует создавать для всех пользователей единственную таблицу для совместного использования временных данных (как делается в Oracle). Это ограничило бы масштабируемость и параллелизм в SQL Server. Все базы данных отличаются друг от друга.

Речь не идет о том, что применять временные таблицы в Oracle невозможно. Вы можете и, скорее всего, будете использовать их. Просто вы будете *применять их в Oracle не так, как в SQL Server* (и наоборот).

Влияние стандартов

Раз уж все базы данных совместимы со стандартом SQL99, то они должны быть одинаковыми. Во всяком случае, так часто полагают. В этом разделе я собираюсь развеять этот миф.

SQL99 — это стандарт баз данных, принятый ANSI/ISO. Он был наследником ANSI/ISO-стандарта SQL92, который, в свою очередь, заменил собой ANSI/ISO-стандарт SQL89. В настоящее время сам SQL99 вытесняется обновлениями стандартов SQL2003, SQL2008 и SQL2011. Стандарт определяет язык (SQL) и поведение (транзакции, уровни изоляции и т.д.) базы данных. Известно ли вам, что многие коммерческие базы данных совместимы с SQL99 хотя бы до определенной степени? А знаете ли вы, что это обстоятельство означает очень мало, когда речь заходит о переносимости запросов и приложений?

Начиная с SQL92, стандарты определяют четыре уровня совместимости.

- **Начальный уровень.** Совместимость на этом уровне обеспечивают большинство поставщиков баз данных. Он является незначительным расширением предшествующего стандарта SQL89. Ни один из поставщиков баз данных не был сертифицирован на более высоком уровне. Более того, в действительности Национальный институт стандартов и технологий (National Institute of Standards and Technology — NIST), который обычно занимается выдачей сертификатов совместимости с SQL, вообще не сертифицирует совместимость на более высоких уровнях. Я был членом группы, которая в 1993 году получила у

NIST сертификат совместимости с SQL92 на начальном уровне для Oracle 7.0. Набор функциональных средств базы данных, совместимых со стандартом на начальном уровне — это подмножество возможностей Oracle 7.0.

- **Переходный уровень.** По объему набора функциональных средств этот уровень находится приблизительно посередине между начальным и промежуточным уровнями.
- **Промежуточный уровень.** Этот уровень включает многие дополнительные возможности, в том числе (но не ограничиваясь ими):
 - динамический SQL;
 - каскадное удаление для обеспечения ссылочной целостности;
 - типы данных DATE и TIME;
 - области;
 - символьные строки переменной длины;
 - выражение CASE;
 - функции приведения для преобразований между типами данных.
- **Уровень полной совместимости.** Обеспечивает поддержку перечисленных ниже возможностей (как и в предыдущем случае, список далеко не полон):
 - управление подключениями;
 - строковый тип данных BIT;
 - откладываемые ограничения целостности;
 - производные таблицы в конструкции FROM;
 - подзапросы в конструкциях проверочного ограничения целостности;
 - временные таблицы.

Стандарт совместимости начального уровня не включает такие функциональные особенности, как внешние соединения, новый синтаксис внутреннего соединения и т.п. Переходный уровень определяет синтаксис внешних и внутренних соединений. Набор средств промежуточного уровня еще шире, а уровень полной совместимости, естественно, включает все, что определено в стандарте SQL92. В большинстве книг, посвященных SQL92, различия между уровнями совместимости никак не описаны, что создает путаницу. В этих книгах демонстрируется, каким образом выглядела бы теоретическая реализация базы данных, полностью совместимая со стандартом SQL92. В результате невозможно взять какую-то книгу по SQL92 и применить изложенный в ней материал к сколько-нибудь реальной базе данных, совместимой с SQL92. Это обусловлено тем, что стандарт определяет не слишком много функциональных средств на начальном уровне, а использование любых средств промежуточного или более высокого уровня увеличивает риск невозможности переноса приложения.

Стандарт SQL99 определяет только два уровня совместимости: основной и расширенный. При разработке этого стандарта была предпринята попытка выхода далеко за рамки традиционного SQL и в нем были введены объектно-реляционные конструкции (массивы, коллекции и т.д.). Стандарт охватывает тип MM (мультимедиа) SQL, объектно-реляционные типы и т.п.

Ни один поставщик не подтверждает сертификатом “совместимость” своих баз данных со стандартом SQL99 на основном или расширенном уровне. Более того, я не знаю ни одного поставщика, который заявлял бы, что его продукт полностью совместим со стандартом на любом из упомянутых уровней.

Вы не должны бояться применять специфичные для поставщика функциональные средства — в конце концов, вы заплатили за них немалые деньги. Каждая база данных обладает собственным набором ловких приемов и в каждой из них всегда можно найти способ выполнения требуемой операции. Поэтому используйте то, что лучше всего подходит для текущей базы данных, а при переходе на другие базы данных заново реализуйте соответствующие компоненты. Для минимизации объема таких изменений применяйте подходящие технологии программирования. Те же самые приемы используются при написании приложений, переносимых между операционными системами; к примеру, их применяют разработчики ядра Oracle.

Удостоверьтесь, что можете адаптироваться

Целью должно быть полноценное использование доступных средств, но убедитесь, что можете изменять реализацию от случая к случаю. В качестве аналогии: Oracle является переносимым приложением, которое функционирует под управлением многочисленных операционных систем. Тем не менее, в среде Windows база данных Oracle работает в стиле Windows: применяет потоки и другие особенности, характерные для Windows. В противоположность этому под управлением UNIX/Linux база данных Oracle запускается как сервер с множеством процессов, используя отдельные процессы для выполнения тех действий, для которых в Windows действуют потоки — т.е. работает в стиле UNIX/Linux. Функциональность “ядра Oracle” доступна на обеих платформах, но внутренне она реализована совершенно по-разному. То же самое относится к приложениям баз данных, которые должны выполняться в средах нескольких СУБД.

Например, распространенной функцией многих приложений баз данных является генерация уникального ключа для каждой строки. Когда вы вставляете строку, система должна автоматически сгенерировать ключ. Для этого в Oracle реализован объект базы данных, который называется `SEQUENCE`, кроме того, функция `SYS_GUID()` также предоставляет уникальные ключи. В Informix имеется тип данных `SERIAL`, а в Sybase и SQL Server — тип `IDENTITY`. В каждой базе данных существует какой-то способ решения этой задачи. Однако конкретные методы различаются как способом выполнения, так и возможным результатом. Таким образом, перед хорошо осведомленным разработчиком открываются две возможности:

- разработка полностью независимого от базы данных метода генерации уникального ключа;
- согласование отличающихся реализаций и применение разных приемов при реализации ключей в каждой базе данных.

На заметку! Теперь в Oracle также имеется тип `IDENTITY`, начиная с версии Oracle 12c. Внутренне он создает последовательность и по умолчанию присваивает очередное значение столбцу, что по поведению очень похоже на тип `IDENTITY` в SQL Server.

Теоретическое преимущество первого подхода заключается в том, что для перехода от одной базы данных к другой ничего изменять не придется. Я назвал это преимущество “теоретическим”, поскольку недостаток такой реализации настолько значителен, что делает все решение совершенно непригодным. Чтобы разработать полностью независимый от базы данных процесс, пришлось бы создать таблицу вроде приведенной ниже:

```
EODA@ORA12CR1> create table id_table
  2  ( id_name varchar2(30) primary key,
  3    id_value number );
Table created.
Таблица создана.

EODA@ORA12CR1> insert into id_table values ( 'MY_KEY', 0 );
1 row created.
1 строка создана.

EODA@ORA12CR1> commit;
Commit complete.
Фиксация завершена.
```

Затем для получения нового ключа понадобилось бы выполнить следующий код:

```
EODA@ORA12CR1> update id_table
  2    set id_value = id_value+1
  3  where id_name = 'MY_KEY';
1 row updated.
1 строка обновлена.

EODA@ORA12CR1> select id_value
  2    from id_table
  3  where id_name = 'MY_KEY';

ID_VALUE
-----
1
```

Все выглядит достаточно просто, но приводит к перечисленным далее последствиям.

- В каждый момент времени обрабатывать строку транзакции может только один пользователь. Указанную строку необходимо обновить, чтобы инкрементировать значение счетчика, а это вынудит программу выполнять операцию последовательно. В лучшем случае генерировать новое значение для этого ключа одновременно сможет только один пользователь.
- В Oracle (и это поведение в других базах данных может отличаться) все пользователи кроме первого, которые попытаются параллельно выполнить эту операцию при уровне изоляции `SERIALIZABLE`, должны будут получить ошибку “ORA-08177: can’t serialize access for this transaction” (“ORA-08177: невозможно обеспечить последовательный доступ для выполнения этой транзакции”).

Например, использование последовательных транзакций (что более распространено в среде J2EE, где многие инструменты автоматически применяют этот режим изоляции как стандартный, часто без ведома разработчиков) обусловило бы нижеследующее поведение.

Обратите внимание, что приглашение командной строки SQL содержит информацию о том, какой сеанс активен в данном примере:

```
ops$tkyte session(419,269)> set transaction isolation level serializable;
Transaction set.
Режим транзакций установлен.

ops$tkyte session(419,269)> update id_table
  2      set id_value = id_value+1
  3      where id_name = 'MY_KEY';
1 row updated.
1 строка обновлена.

ops$tkyte session(419,269)> select id_value
  2      from id_table
  3      where id_name = 'MY_KEY';

  ID_VALUE
-----
        7
```

Теперь мы перейдем в другой сеанс SQL*Plus и выполним ту же самую операцию — параллельный запрос уникального идентификатора:

```
ops$tkyte session(6,479)> set transaction isolation level serializable;
Transaction set.
Режим транзакций установлен.

ops$tkyte session(6,479)> update id_table
  2      set id_value = id_value+1
  3      where id_name = 'MY_KEY';
```

Этот запрос будет заблокирован, поскольку в каждый данный момент времени только одна транзакция может обновлять строку. Приведенный пример демонстрирует первый возможный результат — блокирование и ожидание строки. Но поскольку в Oracle мы используем режим `SERIALIZABLE`, при фиксации транзакции первого сеанса будем наблюдать следующее поведение:

```
ops$tkyte session(419,269)> commit;
Commit complete.
Фиксация завершена.
```

Второй сеанс немедленно отобразит сообщение об ошибке:

```
ops$tkyte session(6,479)> update id_table
  2      set id_value = id_value+1
  3      where id_name = 'MY_KEY';
update id_table
*
ERROR at line 1:
ORA-08177: can't serialize access for this transaction
ОШИБКА В СТРОКЕ 1:
ORA-08177: невозможно обеспечить последовательный доступ для выполнения
этой транзакции
```

Продемонстрированная ошибка возникает независимо от порядка следования оператора фиксации. Все, что для этого необходимо сделать внутри транзакции — попытаться модифицировать любую запись, которая была изменена каким-то другим сеансом после того, как транзакция началась.

Таким образом, приведенный фрагмент кода, независимого от базы данных, в действительности таковым не является. В зависимости от применяемого уровня изоляции он может ненадежно выполняться даже в какой-то одной базе данных! Иногда возникает блокировка с ожиданием, а иногда и ошибка. Не стоит и говорить, что в любом случае (длительное ожидание или длительное ожидание с получением в итоге сообщения об ошибке) конечный пользователь будет недоволен.

Проблема усугубляется тем, что в реальной ситуации транзакция значительно сложнее. Операторы UPDATE и SELECT, приведенные в примере — лишь два из множества тех, которые могут входить в состав транзакции. Мы должны еще вставить в таблицу строку с только что сгенерированным ключом и выполнить любую другую работу, требуемую для завершения этой транзакции. Такое последовательное выполнение операций будет фактором, который всерьез ограничит масштабирование. Подумайте о последствиях использования этой технологии на веб-сайтах, обрабатывающих заказы, и ее применения для генерации номеров заказов. Параллельная работа множества пользователей оказалась бы невозможной, так что пришлось бы выполнять все действия последовательно.

Правильный подход к решению этой проблемы заключается в использовании кода, который наилучшим образом подходит для каждой базы данных. В Oracle 12c это будет следующий код (предполагая, что генерируемый первичный ключ необходим таблице T):

```
EODA@ORA12CR1> create sequence s;  
Sequence created.  
Последовательность создана.  
  
EODA@ORA12CR1> create table t  
2  ( x          number  
3          default s.nextval  
4          constraint t_pk primary key,  
5  other_data varchar2(20)  
6  )  
7  /  
Table created.  
Таблица создана.
```

В качестве альтернативы можно было бы применить атрибут IDENTITY и пропустить генерацию последовательности:

```
EODA@ORA12CR1> create table t  
2  ( x          number  
3          generated as identity  
4          constraint t_pk primary key,  
5  other_data varchar2(20)  
6  )  
7  /  
Table created.  
Таблица создана.
```

На заметку! Столбец IDENTITY должен иметь числовой тип данных.

Но обратите внимание, что генерация последовательности на самом деле не пропускается — последовательность генерируется автоматически базой данных. В ранних выпусках Oracle (11g и предшествующих версиях) для получения автоматически сгенерированного суррогатного первичного ключа обычно использовался следующий код:

```

EODA@ORA12CR1> create table t
2  ( pk number primary key,
3    other_data varchar2(20)
4  )
5  /
Table created.
Таблица создана.

EODA@ORA12CR1> create sequence t_seq;
Sequence created.
Последовательность создана.

EODA@ORA12CR1> create trigger t before insert on t
2  for each row
3  begin
4      :new.pk := t_seq.nextval;
5  end;
6  /
Trigger created.
Триггер создан.

```

На заметку! В выпусках, предшествующих Oracle 11g, вместо присваивания необходимо было применять оператор `SELECT T_SEQ.NEXTVAL INTO :NEW.PK FROM DUAL;`. Прямое присваивание последовательности в PL/SQL является новой возможностью версии 11g.

В результате уникальный ключ будет автоматически — и прозрачно — присваиваться каждой вставленной строке. Ориентированный на достижение более высокой производительности подход в Oracle 11g и предшествующих версиях был бы простым:

```
insert into t ( pk, ... ) values ( t_seq.NEXTVAL, ... );
```

То есть мы полностью избавляемся от накладных расходов, связанных с триггером (именно этому подходу я отдаю предпочтение). Того же самого эффекта можно достичь и в других базах данных, используя их типы. Синтаксис `CREATE TABLE` будет отличаться, но общий результат останется таким же. Здесь мы воспользовались возможностью каждой базы данных по генерации *неблокирующего* и обладающего высокой степенью параллелизма уникального ключа, причем не внесли никаких существенных изменений в код приложения — в этом случае вся логика содержится в DDL-коде.

Уровневое программирование

Теперь, когда вы понимаете, что каждая база данных *будет реализовывать функциональные средства по-разному*, в качестве еще одного примера безопасного программирования, направленного на обеспечение переносимости, можно привести уровневую организацию доступа к базе данных, когда это необходимо.

Давайте предположим, что вы программируете с применением JDBC. Если нужно выполнять только традиционные операции SELECT, INSERT, UPDATE и DELETE, то, скорее всего, уровень абстракции не понадобится. В этом случае SQL-код вполне можно встраивать непосредственно в приложение, по крайней мере, если используемые конструкции ограничены лишь теми, которые поддерживаются каждой целевой базой данных — и, конечно, если вся обработка выполняется одинаково (вспомните обсуждение `NULL=NULL`). Однако это означает, что вы будете располагать неэффективным SQL-кодом. К тому же вам определенно придется обладать большим объемом знаний о многих базах данных, чем большинство их тех, кого я знаю (в конце концов, это единственный способ узнать, есть ли у какого-нибудь кода шанс работать одинаково во всех базах данных!). Еще один подход, обеспечивающий как более высокую переносимость, так и более высокую производительность, предусматривает применение хранимых процедур для возврата результирующих наборов. Вы обнаружите, что база данных каждого поставщика может возвращать результирующие наборы из хранимых процедур, но способы их возврата отличаются. Действительный исходный код хранимых процедур, который вам придется писать, будет выглядеть по-другому для разных баз данных.

У вас есть два варианта — отказаться от использования хранимых процедур для возвращения результирующих наборов, либо реализовать отличающийся код для разных баз данных. Я предпочитаю следовать подходу с написанием разного кода и интенсивно применять хранимые процедуры. На первый взгляд кажется, что при таком подходе реализация приложения в другой базе данных потребует большего времени. Однако вы убедитесь, что этот подход в действительности облегчает реализацию для нескольких СУБД. Вместо того чтобы искать идеальный SQL-код, который работает во *всех* базах данных (возможно, в одних лучше, в других хуже), мы реализуем SQL-код, который наиболее эффективно функционирует в конкретной СУБД. Это можно делать вне самого приложения, что обеспечит дополнительную гибкость в настройке приложения. Вы можете исправить плохо выполняющийся запрос в базе данных и развернуть исправление немедленно, без необходимости во внесении корректировок в приложение. Кроме того, этот метод позволяет свободно применять расширения SQL, предоставляемые поставщиком. Например, Oracle поддерживает широкое разнообразие расширений SQL, среди которых аналитические функции, SQL-конструкция MODEL, сопоставление строк по шаблонам и многие другие. В Oracle можно свободно использовать эти расширения SQL, поскольку они действуют “вне” приложения (т.е. скрыты в базе данных). Для достижения аналогичных целей в других базах данных, скорее всего, придется применять какие-то другие средства, предлагаемые этими СУБД. Вы заплатили за эти возможности, так что вполне можете ими пользоваться.

Еще один аргумент в пользу подхода с разработкой специализированного кода для базы данных, в которой он будет развертываться, связан с тем, что практически невозможно найти одного разработчика (о команде не стоит даже упоминать), достаточно профессионального для того, чтобы понимать все нюансы отличий между Oracle, SQL Server и DB2 (давайте ограничимся всего тремя СУБД). В течение последних 20 лет мне приходилось работать, в основном, с базой данных Oracle (правда, не только с ней). И *ежедневно* я узнавал о ней что-то новое. Весьма сомнительно, что я мог бы выступать в качестве эксперта одновременно по трем упомянутым СУБД и понимать как отличия между всеми ними, так и влияние этих отличий на

уровень “обобщенного кода”, который пришлось бы разрабатывать. Я не уверен, что смог бы качественно или эффективно справиться с такой задачей. Кроме того, примите во внимание, что здесь речь идет об отдельных представителях; сколько разработчиков действительно полностью понимают либо задействуют ту базу данных, которая имеется в их распоряжении, не говоря уже о трех СУБД? Поиск гения, который способен разрабатывать безопасные, масштабируемые, не зависящие от баз данных процедуры, сродни поиску Святого Грааля. Формирование команды таких разработчиков — задача невыполнимая. Найти эксперта по Oracle, эксперта по DB2 и эксперта по SQL Server и сказать им: “Нам нужна транзакция, чтобы сделать X, Y и Z” относительно легко. Далее им следует сообщить: “Здесь находятся входные данные, здесь — данные, которые нужно получить на выходе, а вот — то, к чему приведет этот бизнес-процесс”, после чего они смогут создать транзакционные API-интерфейсы (хранимые процедуры), которые отвечают всем требованиям. Каждый API-интерфейс будет реализован наилучшим образом для конкретной базы данных согласно ее уникальному набору возможностей. Эти разработчики вольны задействовать полную мощь (или в зависимости от случая отказаться от нее) лежащей в основе платформы базы данных.

Те же самые приемы используют разработчики, реализующие многоплатформенный код. Например, в Oracle Corporation эта технология применяется при разработке собственной базы данных. Существует огромный объем кода (хотя он составляет малую долю всего кода базы данных), именуемый кодом *OSD* (Operation System Dependent — зависимый от операционной системы), который реализуется специфическим образом для каждой платформы. За счет применения этого уровня абстракции в Oracle появляется возможность задействовать многие низкоуровневые средств операционной системы для обеспечения высокой производительности и интеграции, что не требует переписывания большей части самой базы данных. Это подтверждается тем фактом, что Oracle может функционировать как многопоточное приложение в Windows и многопроцессное — в UNIX/Linux. Механизмы взаимодействия между процессами абстрагированы до такого уровня, что они могут быть реализованы повторно на основе операционной системы, позволяя радикально отличающимся реализациям вести себя так, как приложение, написанное напрямую и специально для данной платформы.

В дополнение к синтаксическим различиям в SQL, реализации и производительности одного и того же запроса в упомянутых ранее разных базах данных существуют проблемы управления параллельным доступом, уровнями изоляции, согласованности запросов и т.д. Некоторые подробности будут раскрыты в главе 7 настоящей книги, и вы увидите, как эти различия могут коснуться вас. Стандарт SQL92/SQL99 пытается предоставить прямолинейное определение того, каким образом должна работать транзакция и какие уровни изоляции должны быть реализованы, но, в конечном счете, вы все равно получите отличающиеся результаты от разных баз данных. Причина связана с реализацией. В одной базе данных внутри приложения возникает взаимоблокировка, и его работа полностью парализуется. В другой базе данных то же самое приложение работает гладко. В одной базе данных факт установки блокировки (физически последовательной обработки) обеспечивает преимущество, но при развертывании в другой базе данных и отсутствии блокирования получается некорректный ответ. Перенос приложения в другую базу данных требует массы затрат и усилий, даже если вы на 100% следовали стандарту.

Средства и функции

Естественным продолжением довода о том, что вовсе не обязательно стремиться к независимости от базы данных, является идея того, что вам необходимо досконально понимать возможности, которые должна предлагать конкретная база данных, и обеспечить полноценное их использование. Это не раздел, посвященный всем функциональным средствам, которые может предоставить Oracle 12c — данная тема потребовала бы отдельной книги большого объема. В наборе документации Oracle описание новых возможностей Oracle9i, Oracle 10g, Oracle 11g и Oracle 12c занимает целую книгу. При наличии более 10 000 страниц документации, предлагаемой Oracle, раскрытие каждого средства и функции было бы довольно сложным предприятием. Вместо этого цель настоящего раздела заключается в выяснении преимуществ, которые дает даже поверхностное знание доступных возможностей.

Как упоминалось ранее, я отвечаю на вопросы по Oracle в Интернете. Могу сказать, что около 80% этих ответов — просто URL-ссылки на документацию (на каждый опубликованный ответ на вопрос, обычно представляющий собой ссылку на документацию, находятся минимум два неопубликованных вопроса с ответами в форме “прочтите это”). Люди спрашивают, как можно реализовать те или иные сложные функции в базе данных (или за ее пределами), а я лишь указываю на место в документации, где объясняется, каким образом необходимая функция реализована в Oracle, и описан способ работы с ней. Часто встречаются вопросы по репликации. Вот типичный пример вопроса, который мне задают.

Существует ли представление, которое показало бы выполняемый буквальный SQL-оператор? Я имею в виду, что когда делаю выборку из V\$SQL, значение в столбце SQL_TEXT выглядит так: INSERT INTO TABLE1 (COL1,COL2) VALUES (:1,:2). А мне нужно видеть реальные данные, переданные в базу, например: INSERT INTO TABLE1 (COL1,COL2) VALUES ('FirstVal',12). Дело в том, что я пытаюсь получить список операторов вставки, обновления или удаления, выполненных в одной схеме, и запустить те же самые операторы SQL в другой схеме в таком же порядке. Мне бы хотелось иметь код следующего вида:

```
Select SQL_FULLTEXT from V$SQL where FIRST_LOAD_TIME > SYSDATE-(1/24) AND  
➡ (SQL_TEXT like 'INSERT%'...) order by FIRST_LOAD_TIME
```

Такой набор записей можно было бы отправить через веб-службу в другую схему, которая обработала бы эти же операторы. Возможно ли подобное?

Здесь кто-то пытается заново изобрести репликацию! Он не может получить буквальный SQL (и это хорошо), но даже если бы смог, то такой подход никогда бы не заработал. Вы не можете просто взять параллельно выполняемый набор операторов SQL (что произойдет на многопроцессорной машине, где два оператора SQL выполняются в точности одновременно?) и выполнить их последовательно (могут получиться разные ответы!). Вам пришлось бы воспроизвести их с применением той же степени параллелизма, которая имела место в исходной системе.

Например, если два пользователя в одно время выполняют оператор INSERT INTO A_TABLE SELECT * FROM A_TABLE;, то в таблице A_TABLE получится утроенное количество строк по сравнению с ее исходным состоянием. Скажем, если первоначально в таблице A_TABLE было 100 строк, и первый пользователь произведет такую вставку, в ней окажется 200 строк. Если сразу за ним вставку выполнит второй

пользователь (до того, как первый пользователь зафиксирует операцию), то он не увидит 200 строк первого пользователя, а вставит в A_TABLE еще 100 строк, приводя в результате к появлению в таблице 300 строк. Если теперь сделать так, чтобы вместо первого пользователя вставку осуществляла веб-служба (увеличивая A_TABLE со ста до двухсот строк) и затем вставку выполнял второй пользователь (увеличивая A_TABLE с двухсот до четырехсот строк), то возникнет проблема. Репликация — вещь непростая, и на самом деле реализовать ее довольно трудно. В Oracle (и других базах данных) репликация делалась на протяжении более двух десятилетий; ее реализация и сопровождение требовали немалых усилий.

Действительно, вы можете создать собственный механизм репликации, и это может даже оказаться интересным, но по большому счету поступать так не слишком разумно. База данных выполняет массу разнообразных функций. В целом она способна решать задачи лучше, чем это удастся нам самим. Например, репликация встроена в ядро, написанное на языке С. Это быстрый, довольно простой в использовании и надежный механизм. Он работает в разных версиях и на различных платформах. Кроме того, для него обеспечивается поддержка, поэтому при возникновении проблем команда поддержки Oracle всегда окажет необходимую помощь. После модернизации репликация поддерживается также и в новой версии, возможно, предлагая какие-то новые свойства. А теперь представьте, что вы разработали этот механизм самостоятельно. В таком случае придется обеспечить поддержку для всех новых версий, с которыми вы собираетесь иметь дело. А как насчет достижения совместимости между новыми и старыми версиями? Это тоже ваша работа. В случае отказа вы не сможете позвонить в службу поддержки — по крайней мере, до тех пор, пока вы не получите достаточно компактный тестовый сценарий, демонстрирующий основную проблему. При появлении нового выпуска базы данных Oracle вам придется самостоятельно переносить в него собственный код репликации.

Понимание происходящего

Отсутствие полного понимания того, что доступно, в конце концов, может обернуться крупными неприятностями. Как-то мне пришлось сотрудничать с несколькими разработчиками, обладающими многолетним опытом создания приложений баз данных, но в других СУБД. Они построили программное обеспечение анализа (оценки тенденций, создания отчетов и визуализации результатов). Программа работала с клиническими данными, имеющими отношение к здравоохранению. Разработчики ничего не знали о таких синтаксических конструкциях SQL, как встроенные представления, аналитические функции и скалярные подзапросы. Основная проблема заключалась в том, что им требовалось анализировать данные из одной родительской таблицы и двух дочерних таблиц. Диаграмма “сущность-отношение” (Entity Relation Diagram — ERD) могла бы выглядеть похожей на показанную на рис. 1.1.

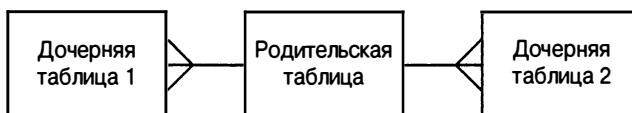


Рис. 1.1. Простая диаграмма “сущность-отношение”

Разработчикам была нужна возможность создания отчета на основе родительской записи с агрегированием данных из каждой дочерней таблицы. СУБД, с которыми разработчикам доводилось работать в прошлом, не поддерживали ни разложение подзапросов на составляющие (конструкция WITH), ни встроенные представления — возможность “запрашивать запросы” вместо запрашивания таблиц. Не зная о существовании этих функциональных средств, разработчики создали на промежуточном уровне своего рода собственную СУБД. Они собирались запрашивать родительскую таблицу и для каждой возвращенной строки выполнять запрос агрегирования с каждой дочерней таблицей. В результате получались тысячи мелких запросов для каждого запроса, выдаваемого конечным пользователем. Или они должны были загружать целые агрегированные дочерние таблицы промежуточного уровня в хеш-таблицы в памяти и реализовать хешированное соединение.

Короче говоря, разработчики заново изобретали СУБД, выполняя функциональные эквиваленты соединения посредством вложенных циклов или хешированного соединения, и не располагая преимуществами временных табличных пространств, оптимизации сложных запросов и т.п. Они потратили свое время на разработку, проектирование, настройку и совершенствование программы, которая пыталась делать то же, что и СУБД, которую они уже приобрели! Между тем конечные пользователи нуждались в новых функциональных средствах, но не получали их, поскольку основную часть времени разработчики посвящали построению этого “механизма” составления отчетов, который действительно был замаскированным механизмом базы данных.

Я показал им, что для сравнения данных, которые были сохранены на разных уровнях детализации, можно выполнять такие действия, как соединение двух агрегаций. При этом возможны различные подходы, как продемонстрировано в листингах 1.1–1.3.

Листинг 1.1. Встроенные представления для запрашивания из запроса

```
select p.id, c1_sum1, c2_sum2
  from p,
       (select id, sum(q1) c1_sum1
        from c1
        group by id) c1,
       (select id, sum(q2) c2_sum2
        from c2
        group by id) c2
 where p.id = c1.id
    and p.id = c2.id
/
```

Листинг 1.2. Скалярные подзапросы, которые выполняют отдельный запрос для каждой строки

```
select p.id,
       (select sum(q1) from c1 where c1.id = p.id) c1_sum1,
       (select sum(q2) from c2 where c2.id = p.id) c2_sum2
  from p
 where p.name = '1234'
/
```

Листинг 1.3. Разложение подзапросов с помощью конструкции WITH

```

with c1_vw as
(select id, sum(q1) c1_sum1
  from c1
 group by id),
c2_vw as
(select id, sum(q2) c2_sum2
  from c2
 group by id),
c1_c2 as
(select c1.id, c1.c1_sum1, c2.c2_sum2
  from c1_vw c1, c2_vw c2
 where c1.id = c2.id)
select p.id, c1_sum1, c2_sum2
  from p, c1_c2
 where p.id = c1_c2.id
/

```

В дополнение к тому, что вы видите в этих листингах, можно также делать замечательные вещи с помощью аналитических функций, таких как LAG, LEAD, ROW_NUMBER, функций ранжирования и многих других. Короче говоря, вместо того, чтобы потратить остаток дня на выяснение приемлемого способа настройки их механизма базы данных промежуточного уровня, мы провели остаток дня за изучением выведенного на экран *справочного руководства по SQL* (вместе с SQL*Plus для создания разовых демонстраций работы тех или иных функций). Конечной целью больше не была настройка промежуточного уровня; теперь они стремились как можно быстрее от него избавиться.

А вот еще один пример. Мне приходилось наблюдать, как люди настраивали в базе данных Oracle процессы-демоны, которые читали сообщения из каналов (механизм взаимодействия между процессами базы данных, реализованный через DBMS_PIPE). Эти процессы-демоны выполняли SQL-код, содержащийся внутри сообщения канала, и фиксировали работу. Они делали это так, что можно было проводить аудит и регистрацию ошибок в транзакции, которая не должна была отменяться при выполнении транзакции более высокого уровня. Обычно если для аудита доступа к каким-то данным используется триггер или что-то подобное, но позже оператор завершается неудачей, то произойдет откат всей работы. Таким образом, за счет отправки сообщения другому процессу они могли поручить выполнение работы и фиксацию отдельной транзакции. Запись аудита сохранялась бы даже в случае отката родительской транзакции. В версиях Oracle, предшествовавших 8i, это был подходящий (и по большому счету единственный) способ реализации такой функциональности. Когда я рассказал им о средстве базы данных под названием *автономные транзакции*, разработчики были обескуражены. Автономные транзакции, реализуемые с помощью единственной строки кода, делали именно то, что разработчики выполняли самостоятельно. Это означало, что можно было отказаться от множества строк кода и, следовательно, от его сопровождения. Вдобавок система в целом начала функционировать быстрее и стала проще для понимания. Но вместе с тем разработчикам было жаль времени, напрасно потраченного на повторное изобретение колеса. В частности, разработчик, который создал эти процессы-демоны, был весьма огорчен тем, что плоды его труда оказались не нужны.

Мне часто приходится сталкиваться с подобными ситуациями — крупные и сложные решения задач, которые уже решены самой СУБД. *Я и сам грешил этим.* Хорошо помню тот день, когда мой консультант по Oracle (тогда я выступал в роли клиента) застал меня окруженным ворохом документации по Oracle. Я взглянул на него и спросил: “Неужели все это правда?” Несколько последующих дней я потратил только на изучение документации. Я самонадеянно полагал, что знаю о базах данных абсолютно все, поскольку ранее работал с SQL/DS, DB2, Ingress, Sybase, Informix, SQLBase, Oracle и другими СУБД. Вместо того чтобы найти время для выяснения возможностей, предлагаемых каждой из них, я лишь применял знания о других базах данных к той, с которой имел дело в текущий момент. (Переход на Sybase/SQL Server был для меня огромным шоком — эта СУБД работала совершенно не так, как остальные.) По мере выяснения того, что могла делать СУБД Oracle (и, откровенно говоря, другие базы данных), я начал пользоваться предоставляемыми ею преимуществами и смог работать быстрее, реализуя меньший объем кода. Это происходило в 1993 году. Только представьте себе, что можно делать с помощью этого программного обеспечения в наши дни, спустя более двадцати лет.

Выделите время на изучение доступного арсенала средств. Пренебрегая этим, вы очень многое теряете. Практически ежедневно я узнаю об Oracle что-то новое. Это должно войти в привычку; я по-прежнему читаю документацию.

Простое решение проблем

Любую проблему всегда можно решить двумя способами: простым и сложным. То и дело мне приходится встречать людей, выбирающих сложный путь. Это не всегда происходит осознанно. Чаше такой выбор обусловлен недостатком знаний. Люди просто не подозревают, что база данных может делать то или иное. Я же, напротив, убежден, что база данных способна решать что угодно, и выбираю сложный путь (самостоятельно пишу код), только когда обнаруживаю, что она что-то не может сделать.

Например, меня часто спрашивают: “Как можно обеспечить наличие у конечного пользователя только одного сеанса в базе данных?” (Я бы мог здесь привести сотни других аналогичных примеров.) Такое требование должно присутствовать во многих приложениях, но мне никогда не приходилось заниматься решением этой задачи — я просто не вижу веской причины для подобного ограничения возможностей, доступных пользователям. Тем не менее, люди хотят делать это и, как правило, строят решение сложным способом. Например, речь может идти о пакетном задании, запускаемом операционной системой, которое будет просматривать таблицу V\$SESSION и произвольно удалять сеансы пользователей, открывших более одного сеанса. В качестве альтернативы они будут создавать собственные таблицы, в которые приложение будет вставлять строку при входе пользователя в систему и удалять ее при его выходе. Такая реализация неминуемо приведет к множеству обращений в службу поддержки, поскольку при аварийном отказе приложения строка останется на месте. Мне встречалось немало “креативных” путей решения этой задачи, но ни один из них не может сравниться по простоте с приведенным ниже:

```
EODA@ORA12CR1> create profile one_session limit sessions_per_user 1;
Profile created.
Профиль создан.
```

```

EODA@ORA12CR1> alter user scott profile one_session;
User altered.

```

Пользователь изменен.

```

EODA@ORA12CR1> alter system set resource_limit=true;
System altered.

```

Система изменена.

Теперь мы попытаемся подключить пользователя SCOTT два раза; вторая попытка должна потерпеть неудачу:

```

EODA@ORA12CR1> connect scott/tiger

```

Connected.

Подключено.

```

SCOTT@ORA12CR1> host sqlplus scott/tiger

```

SQL*Plus: Release 12.1.0.1.0 Production on Fri Mar 14 11:12:04 2014

Copyright (c) 1982, 2013, Oracle. All rights reserved.

ERROR:

ORA-02391: exceeded simultaneous SESSIONS_PER_USER limit

ОШИБКА:

ORA-02391: превышен лимит одновременных сеансов для пользователя

Вот и все — теперь любой пользователь с профилем ONE_SESSION может входить в систему только один раз. Когда я показываю это решение, собеседник обычно хлопает себя по лбу и восклицает: “Я даже не подозревал, что это можно было сделать так”. Ознакомление с возможностями инструментов, с которыми приходится работать, может сэкономить массу времени и усилий в процессе разработки.

Тот же самый довод “сохраняйте простоту” применим и на более широком архитектурном уровне. Я всегда советую хорошо подумать, прежде чем решаться на очень сложные реализации. Чем больше в системе движущихся частей, тем больше элементов, которые могут отказать, а выявление точного места ошибки в излишне сложной архитектуре — задача не простая. Может быть, действительно “круто” реализовать систему с использованием бесчисленных связующих звеньев, но это не является правильным вариантом, если простая хранимая процедура позволяет решить задачу лучше, быстрее и с меньшим числом ресурсов.

Мне доводилось встречать проекты, в которых разработка приложений тянулась долгие месяцы, а конца все не было видно. Разработчики применяли новейшие и наилучшие технологии и языки, но разработка все равно шла не особенно быстро. Приложение было не таким уж большим — вероятно, в этом и заключалась проблема. Когда вы строите конуру (небольшая деревообделочная работа), то не станете пригонять тяжелую технику. Вы будете использовать несколько небольших электроинструментов, а не что-то крупное и сложное. С другой стороны, при построении жилого комплекса вам придется нанять сотни специалистов для работы над проектом и располагать крупной техникой — т.е. применять для решения этой задачи совершенно другие инструменты. То же самое относится и к разработке приложений. Не существует единственной “идеальной архитектуры”. Не существует единственного “идеального языка”. И, конечно, не существует единственного “идеального подхода”.

Например, для построения своего веб-сайта я использовал среду APEX (Application Express). Веб-сайт был крошечным приложением, разработкой которого

занимался один разработчик (или два). Он содержал около 20 экранов. Применение PL/SQL и APEX было правильным выбором для этой реализации — я не нуждался в найме десятков программистов, кодирующих на Java, создающих EJB, использующих Hibernate и т.д. Задача была простой и решалась просто. Существует несколько сложных, крупномасштабных, огромных приложений (большинство из них мы приобретаем в наши дни: системы учета кадров, системы планирования и управления ресурсами предприятия и т.п.), но есть и тысячи небольших приложений. Для выполнения каждой работы необходимо применять соответствующий подход и инструменты.

Я всегда стремлюсь использовать самую простую архитектуру, которая позволяет полностью решить нужную задачу. Окупаемость может оказаться гигантской. Каждая технология занимает свое место. Не все задачи сводятся к забиванию гвоздей, так что в нашем арсенале должен быть не только молоток.

Открытость

Я часто встречаю людей, которые решают задачи сложным путем по другой причине, и это снова связано с идеей о том, что мы должны стремиться к открытости и независимости от базы данных любой ценой. Разработчики хотят избегать применения закрытых, патентованных функциональных возможностей базы данных (даже когда они настолько просты, как хранимые процедуры или последовательности), поскольку их применение приведет к замыканию на конкретную СУБД. Что ж, позвольте мне выдвинуть идею о том, что даже при разработке приложения, читающего и записывающего данные, вы уже в какой-то мере замыкаетесь на СУБД. Вы обнаружите тонкие (а иногда и не настолько тонкие) отличия между базами данных, как только приступите к выполнению запросов и внесению изменений. Например, в одной СУБД может оказаться, что запрос `SELECT COUNT (*) FROM T` попадает в состояние взаимоблокировки при простом обновлении двух строк. Вы обнаружите, что в Oracle запрос `SELECT COUNT (*)` никогда не блокируется для процесса, выполняющего запись. Вы столкнетесь с ситуацией, когда из-за побочных эффектов модели блокировки бизнес-правило выполняется в одной базе данных, но не выполняется в другой. Выяснится, что при одном и том же наборе транзакций запросы дают отличающиеся ответы в разных базах данных — и все это обусловлено фундаментальными различиями в реализации. Вы обнаружите, что лишь очень немногие приложения могут быть перенесены из одной базы данных в другую без каких-либо изменений. Всегда будут присутствовать отличия в интерпретации и обработке SQL-кода (например, вспомните обсуждение `NULL=NULL`).

В одном из проектов разработчики строили веб-приложение, используя Visual Basic, элементы управления ActiveX, сервер IIS и базу данных Oracle. Мне сказали, что разработчики обеспокоены тем, что поскольку бизнес-логика была реализована на PL/SQL, приложение стало зависеть от базы данных, и спросили: “Как мы можем это исправить?”

Я был немного озадачен этим вопросом. Глядя на список выбранных технологий, я не мог понять, почему зависимость от базы данных вызвала обеспокоенность.

- Разработчики выбрали язык, который привязал их к единственной операционной системе, предлагаемой единственным поставщиком (они вполне могли бы выбрать Java).

- Разработчики выбрали технологию построения компонентов, которая привязала их к единственной операционной системе и поставщику (можно было выбрать J2EE).
- Разработчики выбрали веб-сервер, который привязал их к единственному поставщику и одной платформе (почему бы не отдать предпочтение чему-то более открытому?).

Каждая из выбранных разработчиками технологий привязывала их к очень специфичной конфигурации — фактически, *единственной* технологией, которая предоставляла широкий выбор в отношении операционных систем, была база данных.

Невзирая на это (для такого выбора технологий должны были существовать веские причины), целая группа разработчиков сознательно решила не пользоваться функциональностью одного из основных компонентов архитектуры во имя открытости. Я искренне убежден, что следует тщательно выбирать технологии, а затем максимально их эксплуатировать. Вы заплатили за них огромные деньги — так разве не в ваших интересах выжать из них максимум? Полагаю, что разработчики собирались максимально задействовать потенциал других технологий, так почему же база данных должна была стать исключением? Ответить на этот вопрос было еще труднее в свете того, что наиболее полное использование возможностей базы данных имело решающее значение для успеха всего проекта.

С точки зрения открытости на довод можно взглянуть под несколько иным углом. Все данные помещаются в базу. СУБД — очень открытый инструмент. Она поддерживает доступ к данным с помощью широкого многообразия открытых системных протоколов и механизмов доступа. Пока все звучит прекрасно: СУБД — наиболее открытый инструмент в мире.

После этого вы размещаете всю логику приложения и, что еще более важно — *средства безопасности* — за пределами базы данных. Возможно, это будут Java-бины или JSP-страницы, которые осуществляют доступ к данным. Может быть, это будет код на Visual Basic или код, сгенерированный Hibernate. В любом случае вы закрываете базу данных — делаете ее “не открытой”. Пользователи больше не могут выбирать существующие технологии, чтобы работать с этими данными; они *обязаны* применять ваши методы доступа (или полностью обходить средства безопасности). Все это звучит хорошо в текущий момент, но всегда следует помнить, что новейшее достижение на сегодняшний день завтра станет устаревшей технологией. Единственное что сохраняется в реляционном мире на протяжении более 30 лет (и, вероятно, в большинстве объектных реализаций) — это сама база данных. Интерфейсы доступа к данным меняются почти ежегодно, в результате чего приложения, все средства безопасности которых встроены внутрь них, а не вынесены в базу данных, устаревают и становятся препятствием на пути дальнейшего прогресса.

СУБД Oracle предоставляет средство под названием *детальный контроль доступа* (fine-grained access control — FGAC). Если кратко, то эта технология позволяет разработчикам встраивать в базу данных процедуры, которые могут изменять запросы во время их передачи базе данных. Такое изменение запросов направлено на ограничение строк, которые клиент будет получать или модифицировать. Процедура может следить за тем, кто выполняет запрос, когда он запускается, какое приложение запрашивает данные, из какого терминала поступает запрос и т.п., и соответствующим образом ограничивать доступ к данным.

Средство FGAC дает возможность реализовать правила безопасности, такие как описанные ниже.

- Любой запрос, выполненный вне обычных рабочих часов определенным классом пользователей, будет возвращать ноль записей.
- Любые данные могут быть выведены на терминал в защищенном помещении, но только несекретная информация может быть возвращена на терминал удаленного клиента.

По существу FGAC позволяет размещать элемент управления доступом в базе данных, *непосредственно рядом с данными*. Больше не играет роли, как пользователь обращается к данным — из Java-бина, JSP-страницы или приложения Visual Basic через ODBC или SQL*Plus; в любом случае будут использоваться одни и те же протоколы безопасности. Это позволит успешно применить технологию, которая появится следующей.

А теперь я спрашиваю: какая же реализация более “открыта”? Та, которая делает возможным доступ к данным только посредством кода Visual Basic и элементов управления ActiveX (при желании можете заменить Visual Basic на Java, а ActiveX на EJB — здесь речь идет не о выборе конкретной технологии, а о реализации), или же решение, которое обеспечивает доступ из любого приложения, способного взаимодействовать с базой данных, через такие протоколы, как SSL, HTTP и Oracle Net (и другие), либо с использованием API-интерфейсов, подобных ODBC, JDBC, OCI и т.д.? Мне еще не приходилось встречать специализированного инструмента генерации отчетов, который мог бы “запрашивать” код на Visual Basic. Однако мне известны десятки приложений, которые запрашивают SQL-код.

Конечно, люди вольны принимать решение о соблюдении полной открытости и независимости от баз данных. Многие пытаются это делать, но я уверен, что такое решение ошибочно. Независимо от применяемой СУБД, ее следует эксплуатировать максимально полно, задействовав даже самые мелкие функциональные возможности продукта. В любом случае вы обнаружите, что делаете это на этапе настройки (которую, как правило, приходится проводить сразу после развертывания). Просто удивительно, насколько быстро может быть забыто требование независимости от базы данных, когда появляется возможность создания приложения, которое работает в пять раз быстрее всего лишь за счет использования всех возможностей программного обеспечения базы данных.

Как заставить приложение выполняться быстрее?

Вопрос, вынесенный в заголовок этого раздела, мне задают постоянно. Все ищут переключатель вроде “ускорить”, полагая при этом, что “настройка базы данных” означает выполнение настройки самой базы данных. На самом деле, согласно моему опыту, более 80% (а часто все 100%) увеличения производительности должно обеспечиваться на уровне проектирования и реализации приложения — не на уровне базы данных. Вам не удастся провести настройку базы данных до тех пор, пока вы не настроите приложения, работающие с СУБД.

По мере развития технологии появляются переключатели, которые можно устанавливать на уровне базы данных, и которые способствуют уменьшению влияния недостатков программирования. Например, в Oracle 8.1.6 был добавлен новый па-

параметр `CURSOR_SHARING=FORCE`. Этот параметр позволяет при желании реализовать *автоматическую привязку*. Он будет молча преобразовывать запросы вида `SELECT * FROM EMP WHERE EMPNO = 1234` в `SELECT * FROM EMP WHERE EMPNO = :x`. Применение этого параметра *может* радикально уменьшить количество полных разборов и количество ожиданий освобождения библиотечных защепок, которые рассматривались в разделе “Архитектура Oracle” ранее в главе, *но* (всегда есть “но”) могут возникать побочные эффекты. Распространенный побочный эффект с разделением курсора выглядит примерно так:

```
EODA@ORA12CR1> select /* TAG */ substr( username, 1, 1 )
2   from all_users aul
3   where rownum = 1;
```

```
S
-
S
```

```
EODA@ORA12CR1> alter session set cursor_sharing=force;
Session altered.
Сеанс изменен.
```

```
EODA@ORA12CR1> select /* TAG */ substr( username, 1, 1 )
2   from all_users au2
3   where rownum = 1;
```

```
SUBSTR(USERNAME,1,1)
```

```
-----
S
```

Что здесь случилось? Почему столбец, о котором сообщил SQL*Plus, неожиданно оказался настолько большим во втором запросе, который на первый взгляд ничем не отличается от первого? Это (и кое-что еще) станет очевидным, если взглянуть на то, к чему приводит установка разделения курсора:

```
EODA@ORA12CR1> select sql_text from v$sql where sql_text like 'select /* TAG */ %';
SQL_TEXT
```

```
-----
select /* TAG */ substr( username, 1, 1 )
from all_users aul
where rownum = 1

select /* TAG */ substr( username, : "SYS_B_0", : "SYS_B_1" )
from all_users au2
where rownum = : "SYS_B_2"
```

Разделение курсора удаляет информацию из запроса. Запрос нашел *все* литералы, в том числе константы встроенной функции `substr`, которые мы использовали. Они были удалены из запроса и заменены переменными привязки. Теперь механизм SQL не знает, что столбец является подстрокой с длиной 1 — его длина не определена. Кроме того, как видите, конструкция `where rownum = 1` также стала привязанной. На первый взгляд это кажется хорошей идеей; тем не менее, оптимизатор только что лишился некоторой важной информации. Он больше не знает, что “этот запрос будет извлекать одиночную строку”; он уверен, что “этот запрос будет возвращать первые *N* строк, и *N* может быть любым числом”. Такой результат может отрицательно повлиять на сгенерированные планы выполнения запросов.

Кроме того, я убедился, что хотя режим `CURSOR_SHARING=FORCE` обеспечивает значительно более быстрое выполнение, чем синтаксический разбор и оптимизация множества уникальных запросов (см. раздел “Используйте переменные привязки”), он работает все же медленнее, чем в случае применения запросов, где разработчик делал привязку самостоятельно. Это происходит не из-за какой-то неэффективности в коде разделения курсора, а скорее из-за неэффективности внутри самой программы. Во многих случаях приложение, не использующее переменные привязки, не позволяет ни эффективно выполнить синтаксический разбор, ни повторно использовать курсоры. Поскольку приложение уверено, что каждый запрос уникален (запросы построены в форме уникальных операторов), оно никогда не будет применять курсор более одного раза. Реальность заключается в том, что если программист с самого начала использовал переменные привязки, то он может выполнить разбор запроса один раз, а затем применять его многократно. Именно накладные расходы по синтаксическому разбору снижают общую производительность.

На заметку! Важно также отметить, что установка `CURSOR_SHARING=FORCE` не устраняет ошибку внедрения SQL. Привязка выполняется после того, как запрос был переписан конечным пользователем; внедрение SQL уже произошло. Посредством `CURSOR_SHARING=FORCE` нельзя получить более высокую степень защиты, чем она была. Сопrotивляемость внедрению SQL можно обеспечить только за счет использования самих переменных привязки.

В основном важно иметь в виду, что сама по себе настройка `CURSOR_SHARING=FORCE` не обязательно решит проблемы. Она вполне может породить новые сложности. В некоторых случаях параметр `CURSOR_SHARING` — очень полезный инструмент, но он отнюдь не является панацеей от всех бед. В тщательно разработанном приложении он никогда не понадобится. С точки зрения долговременной перспективы правильным подходом является применение переменных привязки там, где это уместно, и констант — когда это необходимо.

На заметку! Не существует никаких магических средств — *вообще никаких*. Если бы они были, то превратились бы в стандартное поведение, и вы никогда бы о них не услышали.

Хотя есть переключатели, которые можно использовать на уровне базы данных (они действительно немногочисленные и редко встречающиеся), то проблемы, связанные с параллельной обработкой и неэффективным выполнением запросов (из-за плохо написанных запросов или плохо структурированных данных), решить с их помощью не удастся. Подобные ситуации требуют переписывания кода (и нередко изменения архитектуры). Перемещение файлов данных, настройка параметров и других переключателей уровня базы данных часто оказывает лишь незначительное влияние на общую производительность приложения. Во всяком случае, эти действия определенно не обеспечат двух-, трех- или *n*-кратного повышения производительности, которое требуется для того, чтобы приложение стало приемлемым. Сколько раз вы получали приложение, которое работало на 10% медленнее желаемого? Никто не жалуется на 10-процентное замедление. А вот *пятикратное* замедление непременно вызовет недовольство пользователей. Повторю еще раз: пятикратного увеличения производительности не удастся достичь перемещением файлов данных.

Этого можно добиться только путем внесения в приложение корректировок, возможно, направленных на значительное уменьшение количества выполняемых операций ввода-вывода.

На заметку! Просто чтобы отметить, как вещи меняются со временем. Я часто пишу о том, что вы не получите пятикратного роста производительности, перемещая файлы данных с места на место. С появлением аппаратных решений, таких как Oracle Exadata (устройство сети хранения данных, спроектированное в качестве расширения базы данных), вы на самом деле можете получить 5-, 10-, 50-кратное и более сокращение времени ответа, просто переместив туда файлы данных. Но это скорее история о том, как “мы полностью сменили нашу аппаратную архитектуру”, чем о том, как “мы реорганизовали наше хранилище”. Кроме того, увеличение скорости работы приложения лишь в 5 или 10 раз на Exadata лично меня бы разочаровало — хотелось бы получить 50-кратный рост и более — и потребовало переосмысления реализации самого приложения.

Производительность — это характеристика, которую нужно учитывать при проектировании, построении и тестировании на протяжении всего периода разработки. Ее учет никогда не следует переносить на момент, когда приложение готово. Меня удивляет, насколько часто люди дожидаются поставки приложения клиенту, его установки и начала эксплуатации, прежде чем приступают к его тонкой настройке. Мне приходилось сталкиваться с ситуациями, когда приложения поставлялись только с первичными ключами — вообще без каких-либо других индексов. Запросы никогда не настраивались и не подвергались нагрузочному тестированию. Приложения никогда не испытывались при наличии достаточно большого количества пользователей. Настройка считалась частью процесса установки программного продукта. На мой взгляд, этот подход совершенно неприемлем. Конечные пользователи должны с первого дня получить работоспособную и полностью настроенную систему. “Проблем с продуктом” будет более чем достаточно и без того, чтобы пользователи с самого начала сталкивались с низкой производительностью. Пользователи готовы к тому, что новое приложение будет содержать несколько программных ошибок, но не заставляйте их мучительно долго дожидаться появления результатов на экране.

Отношения между администратором базы данных и разработчиком

Не вызывает сомнений, что большинство успешных информационных систем основано на отношениях сотрудничества между администратором базы данных и разработчиком приложения. В этом разделе я хочу привести точку зрения разработчика на разделение труда между разработчиком и администратором базы данных (исходя из предположения, что команда разработки каждого серьезного проекта включает в себя группу администраторов базы данных).

Разработчику не обязательно знать, как устанавливать и конфигурировать программное обеспечение. Это задача администратора базы данных и, возможно, системного администратора. Настройка Oracle Net, запуск прослушивающего процесса, конфигурирование сервера совместного использования, включение пула соединений, установка СУБД, создание базы данных и т.п. — все эти функции я передаю в руки администратора базы данных/системного администратора.

В общем случае разработчику не обязательно знать, каким образом настраивается операционная система. Обычно я оставляю эту работу системным администраторам.

торам. Разработчик приложений баз данных должен быть компетентен в вопросах использования выбранной операционной системы, но от него не следует ожидать умения ее настраивать.

Единственная главнейшая ответственность администратора базы данных — восстановление базы данных. Обратите внимание, речь идет не о “резервном копировании”, а о “восстановлении”, и я готов присягнуть, что это — единственная ответственность администратора базы данных. Понимание того, как работает откат и повторение — да, это то, что разработчик должен знать. Знание того, как выполнять восстановление табличной области на конкретный момент времени — нечто такое, что без чего разработчик вполне может обойтись. Знать о том, что подобное возможно, будет нелишним, но действительно делать его — нет никакой необходимости.

Настройка на уровне экземпляра базы данных и выяснение оптимального значения для параметра `PGA_AGGREGATE_TARGET`, как правило, является задачей администраторов базы данных (и СУБД располагает достаточным набором средств, облегчающих определение корректной конфигурации). В ряде исключительных случаев разработчику может понадобиться изменить некоторые настройки сеанса, но на уровне базы данных за это отвечает администратор. Типичная СУБД поддерживает приложения, реализованные более чем одним разработчиком. Только администратор базы данных, который поддерживает все приложения, может принять правильное решение.

Выделение дискового пространства и управление файлами — также задача администратора базы данных. Разработчики могут вносить свои предложения относительно дискового пространства (указывать, сколько, по их мнению, им потребуется места), но обо всем остальном позаботится администратор базы данных/системный администратор.

В общем случае от разработчиков не требуется знание способов эксплуатации СУБД. Они должны знать, как запускать программы в СУБД. Разработчик и администратор базы данных совместно работают над различными частями одной и той же головоломки. Администратор базы данных будет обращаться к разработчику, если запросы потребляют слишком много ресурсов. Разработчик будет обращаться к администратору, если ему не удастся понять, каким образом еще больше ускорить работу системы (это именно тот случай, когда можно прибегнуть к настройке экземпляра после того, как приложение полностью настроено).

Все это варьируется в зависимости от среды, но я — приверженец разделения труда. Хороший разработчик обычно является плохим администратором базы данных и наоборот. По моему убеждению, эти занятия требуют двух разных профессиональных навыков, двух отличающихся образов мышления и двух разных личностных качеств.

Резюме

В этой главе мы несколько обобщенно рассмотрели причины, которые обуславливают необходимость знания СУБД. Приведенные мною примеры отнюдь не выдуманы — они встречаются ежедневно. С подобными проблемами мне приходится сталкиваться постоянно.

Давайте кратко вспомним ключевые моменты. Если разработка выполняется с применением Oracle, важны перечисленные ниже аспекты.

- Необходимо понимать архитектуру Oracle. Не обязательно знать ее до такой степени, чтобы быть в состоянии заново написать код сервера, но в ней необходимо разбираться в достаточной мере, чтобы иметь понятие о последствиях использования конкретного функционального средства.
- Необходимо знать особенности блокировки и управления параллельной обработкой и то, что в каждой базе данных они реализованы по-другому. В противном случае база данных будет выдавать “неправильные” ответы и у вас возникнут большие проблемы соперничества, ведущие к низкой производительности.
- Не следует трактовать базу данных как черный ящик — т.е. то, что не требует понимания. База данных является наиболее важным компонентом большинства приложений. Попытка ее игнорирования приведет к фатальным последствиям.
- Не изобретайте колесо. Я видел не одну команду разработчиков, которые сталкивались с проблемами — не только на техническом, но и на личностном уровне — из-за неосведомленности о том, что свободно предлагается Oracle. Это происходит, когда кто-то обнаруживает, что средство, на реализацию которого потрачено пару месяцев, уже есть в ядре базы данных. Читайте документацию, которая поставляется вместе с программным обеспечением!
- Решайте задачи наиболее простым способом, максимально используя встроенные функциональные возможности Oracle. За них были уплачены немалые деньги.
- Программные проекты приходят и уходят, равно как языки программирования и платформы. Разработчики могут заниматься обеспечением работоспособности систем несколько недель, возможно, месяцев, а затем переходить к решению следующей задачи. Если снова и снова заниматься изобретением колеса, то никогда не удастся добиться значительных успехов в разработке. Подобно тому, как вы никогда не станете строить собственный класс хеш-таблицы в Java, поскольку в Java такой класс уже есть, вы должны применять имеющиеся в вашем распоряжении функциональные возможности базы данных. Естественно, первый шаг на пути к этому — понимание того, чем вы располагаете. Читайте документацию.

Продолжая последнюю мысль — проекты программного обеспечения и языки программирования приходят и уходят, но *данные* существуют всегда. Мы строим приложения, которые используют данные, а со временем эти данные будут применяться многими приложениями. Таким образом, речь идет не о приложениях, а о данных. Задействуйте технологии и реализации, которые допускают использование и многократное использование данных. Если вы применяете базу данных в качестве хранилища, но открываете доступ к любым данным только через приложение, конечная цель не будет достигнута. Вы не сможете произвольно запрашивать свое приложение и не сможете построить новое приложение поверх старого. Но если вы используете базу данных, то обнаружите, что с течением времени добавление новых приложений, отчетов или чего-нибудь другого становится значительно проще.

ГЛАВА 2

Обзор архитектуры

Система управления базами данных Oracle была спроектирована как в высшей степени переносимая — она доступна на всех значимых платформах, от Windows и UNIX/Linux до мэйнфреймов. Однако физическая архитектура Oracle в разных операционных системах отличается. Например, в UNIX/Linux вы увидите, что СУБД Oracle реализована как совокупность процессов операционной системы, фактически по одному процессу на каждую крупную функцию. Для UNIX/Linux такая реализация вполне естественна, поскольку работа происходит в среде с множеством процессов. Тем не менее, в Windows эта архитектура была бы неподходящей и работала бы не особенно эффективно (оказалась бы медленной и не поддающейся масштабированию). Поэтому на платформе Windows СУБД Oracle реализована в виде одного процесса с множеством потоков. В прошлом на мэйнфреймах, работающих под управлением OS/390 и z/OS, архитектура Oracle, специфичная для этих операционных систем, задействует набор адресных пространств OS/390, которые все функционируют как единственный экземпляр Oracle. Для одного экземпляра базы данных можно сконфигурировать до 255 адресных пространств. Кроме того, СУБД Oracle работает совместно с диспетчером рабочей загрузки (Workload Manager — WLM) при установлении приоритета выполнения конкретных рабочих нагрузок Oracle по отношению друг к другу и ко всем другим задачам, выполняемым в системе OS/390. Несмотря на то что физические механизмы, применяемые в реализациях Oracle, варьируются от платформы к платформе, архитектура в достаточной степени обобщена, чтобы можно было получить хорошее представление о ее работе на всех платформах.

В настоящей главе я представлю обширное описание этой архитектуры. Мы взглянем на сервер Oracle и определим ряд терминов, таких как *база данных*, *подключаемая база данных*, *контейнерная база данных* и *экземпляр* (термины, которые практически всегда вызывают путаницу). Мы посмотрим, что происходит при “подключении” к Oracle, и ознакомимся с тем, как сервер управляет памятью на высоком уровне. В трех последующих главах мы рассмотрим три основных компонента архитектуры Oracle.

- В главе 3 освещаются файлы. В ней мы взглянем на пять основных категорий файлов, которые образуют базу данных: файлы параметров, файлы данных, временные файлы, управляющие файлы и журнальные файлы для восстановления. Мы также рассмотрим другие типы файлов, включая файлы трассировки, файлы предупреждений, файлы дампа (DMP), файлы помпы данных и простые плоские файлы. Вы узнаете о файловой области (в Oracle 10g и пос-

ледующих версиях), которая называется Fast Recovery Area (Область быстрого восстановления). Мы также обсудим влияние автоматического управления памятью (Automatic Storage Management — ASM) на файловое хранилище.

- В главе 4 раскрываются структуры памяти Oracle, называемые SGA (System Global Area — системная глобальная область), PGA (Process Global Area — глобальная область процесса) и UGA (User Global Area — глобальная область пользователя). Мы исследуем взаимосвязь между этими структурами, а также обсудим разделяемый пул, большой пул, пул Java и другие компоненты структуры SGA.
- В главе 5 рассматриваются физические процессы или потоки Oracle. Мы взглянем на три отличающихся типа процессов, которые будут выполняться в базе данных: серверные процессы, фоновые процессы и подчиненные процессы.

Мне было нелегко решить, какие из этих компонентов раскрывать первыми. Структура SGA используется процессами, поэтому рассмотрение SGA до ознакомления с процессами может быть лишено смысла. С другой стороны, при обсуждении процессов и выполняемых ими действий неизбежны ссылки на SGA. Эти два компонента тесно связаны: процессы воздействуют на файлы, и рассмотрение файлов без предварительного ознакомления с тем, что делают процессы, было бы лишено смысла.

Поэтому я определяю некоторые термины и предоставляю общий обзор о том, на что похожа база данных Oracle (так, как если бы изучение начиналось с нуля). Обсудить предстоит две архитектуры. Первая из них — это архитектура, которую база данных Oracle применяла всецело, начиная с версии 6 и заканчивая версией 11g (здесь она называется архитектурой с *единственным владельцем* (single tenant)), и новой архитектурой с *множеством владельцев* (multi-tenant), появившейся в версии Oracle 12c. После этого можно будет приступить к ознакомлению с некоторыми подробностями.

Определение базы данных и экземпляра

Существуют два термина, которые при использовании в контексте Oracle, похоже, вызывают наибольшую путаницу: *база данных* и *экземпляр*. В терминологии Oracle они определяются следующим образом.

- **База данных.** Коллекция физических файлов или дисков операционной системы. В случае применения разделов ASM или RAW база данных может не выглядеть как отдельные файлы в среде операционной системы, однако это определение остается в силе. В версии Oracle Database 12c существуют три типа баз данных.
- **База данных с единственным владельцем (single tenant).** Это обособленный набор файлов данных, управляющих файлов, журнальных файлов для восстановления, файлов параметров и т.д., который вдобавок ко всем метаданным, данным и коду приложения содержит все метаданные Oracle (например, определение ALL_OBJECTS), данные Oracle и код Oracle (такой как код для DBMS_OUTPUT). В версиях, предшествующих 12c, это был единственный тип базы данных.

- **Контейнерная (container) или корневая (root) база данных.** Это обособленный набор файлов данных, управляющих файлов, журнальных файлов для восстановления, файлов параметров и т.д., который включает только метаданные Oracle, данные Oracle и код Oracle. В таких файлах данных отсутствуют какие-либо объекты или код приложения — только метаданные и объекты кода, предоставленные Oracle. База данных является обособленной в том смысле, что она может быть смонтирована и открыта без дополнительных поддерживающих физических структур.
- **Подключаемая (pluggable) база данных.** Это набор одних лишь файлов данных. Такая база данных не является обособленной. Чтобы появилась возможность открытия и доступа, подключаемая база данных должна быть “вставлена” в контейнерную базу данных. Файлы данных содержат только метаданные для объектов приложений, данные приложений и код для приложений. Метаданные или любой код Oracle в этих файлах данных отсутствует. Нет каких-либо журнальных файлов для восстановления, управляющих файлов, файлов параметров и т.п. — только файлы данных, ассоциированные с подключаемой базой данных. Подключаемая база данных наследует другие типы файлов от контейнерной базы данных, к которой она подключена в текущий момент.
- **Экземпляр.** Набор фоновых процессов или потоков и области разделяемой памяти Oracle, представляющей собой память, которая совместно используется этими потоками или процессами, выполняющимися на одном компьютере. Это место для изменяющихся, непостоянных данных, часть которых сбрасывается на диск. Экземпляр базы данных может существовать вообще без какого-либо дискового хранилища. Возможно, это и не самая полезная вещь в мире, но такое представление о нем позволяет проводить линию между экземпляром и базой данных.

Иногда эти два термина, *экземпляр* и *база данных*, применяют взаимозаменяемо, но они отражают очень разные концепции, особенно теперь, в архитектуре с множеством владельцев. Отношение между ними заключается в том, что база данных с единственным владельцем или контейнерная база данных (в главе под просто *базой данных* подразумевается либо база данных с единственным владельцем, либо контейнерная база данных; при обсуждении подключаемых баз данных характеристика *подключаемая* указывается явно) может быть смонтирована и открыта несколькими экземплярами. Экземпляр может монтировать и открывать только одну базу данных в каждый момент времени. Фактически, правильно сказать, что экземпляр будет монтировать и открывать самое большее одну базу данных на протяжении всего времени своего существования! Вскоре мы рассмотрим пример.

Подключаемая база данных будет ассоциирована с одной контейнерной базой данных за раз и только непрямо связана с экземпляром; она будет совместно использовать экземпляр, созданный для монтирования и открытия контейнерной базы данных. Таким образом, подобно контейнерной базе данных, в каждый момент времени подключаемая база данных может быть ассоциирована с одним или несколькими экземплярами. Однако в отличие от базы данных с единственным владельцем, экземпляр может предоставлять доступ ко многим (приблизительно до 250) подключаемым базам данных одновременно. То есть единственный экземпляр может

предлагать службы для множества подключаемых баз данных, но только для одной контейнерной базы данных или базы данных с единственным владельцем.

Еще больше запутались? Дальнейшие пояснения должны облегчить понимание этих концепций.

Экземпляр представляет собой просто набор процессов операционной системы или один процесс с множеством потоков и определенной областью памяти. Эти процессы могут действовать в единственной базе данных. База данных является всего лишь коллекцией файлов (файлов данных, временных файлов, журнальных файлов для восстановления и управляющих файлов). В любое время экземпляр будет связан только с одним набором файлов (одной контейнерной базой данных или базой данных с единственным владельцем). Можно открывать и получать доступ к нескольким подключаемым базам данных, которые подчиняются контейнерной базе данных, но все они будут совместно использовать единственный экземпляр, созданный для открытия этой контейнерной базы данных.

В большинстве случаев обратное утверждение также справедливо: контейнерная база данных или база данных с единственным владельцем будет иметь только один работающий с ней экземпляр. Однако в особом случае применения технологии Oracle RAC (Real Application Clusters — кластеры для реальных приложений), которая позволяет Oracle функционировать на множестве компьютеров в кластеризованной среде, мы можем иметь несколько экземпляров, одновременно монтирующих и открывающих одну базу данных, которая размещена на наборе разделяемых физических дисков. Это позволяет получать доступ к единственной базе данных из множества разных компьютеров в одно и то же время. Технология Oracle RAC предназначена для систем с исключительно высокой готовностью и обладает потенциалом для проектирования в высшей степени масштабируемых решений.

Давайте начнем с рассмотрения простого примера. Предположим, что мы только что установили Oracle 12c версии 12.1.0.1 на компьютер, функционирующий под управлением UNIX/Linux. Было установлено только программное обеспечение без каких-либо начальных баз данных и других компонентов — т.е. ничего кроме программного обеспечения.

Команда `pwd` отображает текущий рабочий каталог, `db`s (в Windows это был бы каталог `database`), и команда `ls -l` показывает, что он пуст. Он не содержит ни файла `init.ora`, ни файлов `SPFILE` (файлы хранимых параметров, которые подробно обсуждаются в главе 3):

```
[oral2crl@dellpe dbs]$ pwd
/home/oral2crl/app/oral2crl/product/12.1.0/dbhome_1/dbs
[oral2crl@dellpe dbs]$ ls -l
total 0
```

С помощью команды `ps` (`process status` — состояние процесса) можно вывести информацию обо всех процессах, выполняемых пользователем `oral2crl` (владелец программного обеспечения Oracle в этом случае). В данный момент процессы базы данных Oracle вообще отсутствуют:

```
[oral2crl@dellpe dbs]$ ps -aef | grep oral2crl
root      18392 15416  0 14:31 pts/1    00:00:00 su - oral2crl
oral2crl  18401 18392  0 14:31 pts/1    00:00:00 -bash
oral2crl  18461 18401  0 14:34 pts/1    00:00:00 ps -aef
oral2crl  18462 18401  0 14:34 pts/1    00:00:00 grep oral2crl
```

Теперь можно ввести `ipcs` — команду UNIX/Linux для отображения сведений об устройствах взаимодействия между процессами, таких как разделяемая память, семафоры и т.п. В настоящее время в системе такие устройства не используются:

```
[oral2crl@dellpe dbs]$ ipcs -a
----- Shared Memory Segments -----
----- Сегменты разделяемой памяти ----
key      shmid      owner      perms      bytes      nattch      status
----- Semaphore Arrays -----
----- Массивы семафоров -----
key      semid      owner      perms      nsems
----- Message Queues -----
----- Очереди сообщений -----
key      msqid      owner      perms      used-bytes  messages
```

Теперь запустим SQL*Plus (интерфейс командной строки Oracle) и подключимся к базе данных как `sysdba` (учетная запись, которая позволяет делать в базе данных практически что угодно). Изначально, предполагая, что переменная среды `ORACLE_SID` еще не установлена, вы увидите следующий вывод:

```
[oral2crl@dellpe dbs]$ sqlplus / as sysdba
SQL*Plus: Release 12.1.0.1.0 Production on Mon Sep 2 14:35:52 2013
Copyright (c) 1982, 2013, Oracle. All rights reserved.
ERROR:
ORA-12162: TNS:net service name is incorrectly specified
ОШИБКА:
ORA-12162: TNS: имя сетевой службы указано некорректно
```

Эта ошибка происходит из-за того, что программному обеспечению базы данных ничего не известно о том, куда следует подключаться. Когда производится подключение, программное обеспечение Oracle ищет строку соединения TNS (сетевое соединение). Если строка не предоставлена, как в этом примере, Oracle ищет переменную среды `ORACLE_SID` (в Windows также проверяется параметр реестра `ORACLE_SID`). Переменная `ORACLE_SID` — это “идентификатор сайта” Oracle, т.е. своего рода ключ для получения доступа к экземпляру. Если установить `ORACLE_SID`, как показано ниже:

```
[oral2crl@dellpe dbs]$ export ORACLE_SID=oral2c
```

то подключение будет успешным и SQL*Plus сообщит о соединении с простаивающим экземпляром:

```
[oral2crl@dellpe dbs]$ sqlplus / as sysdba
SQL*Plus: Release 12.1.0.1.0 Production on Mon Sep 2 14:36:54 2013
Copyright (c) 1982, 2013, Oracle. All rights reserved.
Connected to an idle instance.
Подключено к простаивающему экземпляру.
SQL>
```

Прямо сейчас наш “экземпляр” состоит только из процесса сервера Oracle, который в следующем выводе выделен полужирным. Пока еще не существует ни выделенной разделяемой памяти, ни других процессов.

```
SQL> !ps -aef | grep ora12crl
root      18392 15416  0 14:31 pts/1    00:00:00 su - ora12crl
ora12crl 18401 18392  0 14:31 pts/1    00:00:00 -bash
ora12crl 18474 18473  0 14:36 pts/0    00:00:00 ../dbhome_1/bin/sqlplus as
                                     sysdba
ora12crl 18475 18474  0 14:36 ?          00:00:00 oracleora12c
(DESCRIPTION=(LOCAL=YES) (ADDRESS=(PROTOCOL=beq) ))
ora12crl 18482 18474  0 14:38 pts/0    00:00:00 /bin/bash -c ps -aef |
                                     grep ora12crl
ora12crl 18483 18482  0 14:38 pts/0    00:00:00 ps -aef
ora12crl 18484 18482  0 14:38 pts/0    00:00:00 grep ora12crl

SQL> !ipcs -a
----- Shared Memory Segments -----
----- Сегменты разделяемой памяти ---
key      shmid      owner      perms      bytes      nattch     status
----- Semaphore Arrays -----
----- Массивы семафоров -----
key      semid      owner      perms      nsems
----- Message Queues -----
----- Очереди сообщений -----
key      msqid      owner      perms      used-bytes  messages
```

На заметку! В среде Windows программное обеспечение Oracle выполняется как единственный процесс с потоками; вы не увидите отдельных процессов, как в среде UNIX/Linux. Более того, потоки Windows не будут иметь таких же имен, как у только что показанных процессов. Я специально использую здесь UNIX/Linux, чтобы мы могли различать индивидуальные процессы и четко “видеть” их.

В этом выводе команды `ps` интересно отметить процесс по имени `oracleora12c`. Как бы тщательно вы ни просматривали систему, вы не найдете исполняемого файла с таким именем. На самом деле этот процесс запускается двоичным файлом `$ORACLE_HOME/bin/oracle`.

На заметку! Предполагается, что в переменной среды (в UNIX/Linux) или в параметре реестра (в Windows) по имени `ORACLE_HOME` был указан полностью определенный путь к месту, где установлено программное обеспечение Oracle.

Разработчики Oracle просто переименовывают этот процесс при его загрузке в память. Именем единственного процесса Oracle, который выполняется прямо сейчас (наш *выделенный серверный процесс*, о котором речь пойдет позже), является `oracle$ORACLE_SID`. Такое соглашение об именовании позволяет очень легко увидеть, каким образом процессы ассоциируются с экземплярами. Итак, попробуем запустить экземпляр:

```
SQL> startup
ORA-01078: failure in processing system parameters
LRM-00109: could not open parameter file '/home/oral2crl/app/oral2crl/
product/12.1.0/dbhome_1/dbs/initoral2c.ora'
ORA-01078: отказ при обработке системных параметров
LRM-00109: не удалось открыть файл параметров '/home/oral2crl/app/
oral2crl/product/12.1.0/dbhome_1/dbs/initoral2c.ora'
```

Обратите внимание на ошибку, возникшую из-за отсутствия файла `initoral2c.ora`. Этот файл, в просторечии называемый *файлом* `init.ora` или более точно *файлом параметров*, является единственным файлом, который должен существовать, чтобы можно было запустить экземпляр — необходимо иметь либо файл параметров (простой плоский файл, который будет более подробно описан ниже), либо файл хранимых параметров.

Давайте создадим файл параметров и поместим в него минимальную информацию, требуемую для запуска экземпляра базы данных. (Обычно приходится задавать значительно больше параметров, в числе которых размер блока базы данных, местоположения управляющих файлов и т.п.) По умолчанию этот файл находится в каталоге `$ORACLE_HOME/dbs` и имеет имя `init${ORACLE_SID}.ora`:

```
[oral2crl@dellpe dbs]$ cd $ORACLE_HOME/dbs
[oral2crl@dellpe dbs]$ echo db_name=oral2c > initoral2c.ora
[oral2crl@dellpe dbs]$ cat initoral2c.ora
db_name=oral2c
```

Теперь возвратимся в SQL*Plus:

```
[oral2crl@dellpe dbs]$ sqlplus / as sysdba
SQL*Plus: Release 12.1.0.1.0 Production on Mon Sep 2 14:42:27 2013
Copyright (c) 1982, 2013, Oracle. All rights reserved.
Connected to an idle instance.
Подключено к простаивающему экземпляру.
```

```
SQL> startup nomount
ORACLE instance started.
Экземпляр ORACLE запущен.
```

```
Total System Global Area    329895936 bytes
Вся системная глобальная область
Fixed Size                    2287960 bytes
Фиксированный размер
Variable Size                 272631464 bytes
Переменный размер
Database Buffers              50331648 bytes
Буферы базы данных
Redo Buffers                  4644864 bytes
Журнальные буферы
```

Мы используем опцию `nomount` команды `startup`, т.к. пока еще нет базы данных, подлежащей монтированию (все опции команд запуска и завершения можно найти в документации по SQL*Plus).

На заметку! В среде Windows перед запуском команды startup понадобится выполнить оператор создания службы с применением утилиты oradim.exe.

Теперь мы получили то, что можно было бы назвать экземпляром. Все фоновые процессы, необходимые для действительного запуска базы данных, такие как монитор процессов (process monitor — pmon), процесс записи в журнал (lgwr) и т.д. (все они подробно описаны в главе 5), запущены. Давайте взглянем на них:

```
SQL> !ps -aef | grep ora12cr1
root      18392 15416  0 14:31 pts/1    00:00:00 su - ora12cr1
ora12cr1 18401 18392  0 14:31 pts/1    00:00:00 -bash
ora12cr1 18499 18401  0 14:42 pts/1    00:00:00 rlwrap /home/ora12cr1/app/
ora12cr1 18500 18499  0 14:42 pts/0    00:00:00 /home/ora12cr1/app/
ora12cr1 18508      1 0 14:43 ?        00:00:00 ora_pmon_ora12c
ora12cr1 18510      1 0 14:43 ?        00:00:00 ora_psp0_ora12c
ora12cr1 18512      1 0 14:43 ?        00:00:00 ora_vktm_ora12c
ora12cr1 18516      1 0 14:43 ?        00:00:00 ora_gen0_ora12c
ora12cr1 18518      1 0 14:43 ?        00:00:00 ora_mman_ora12c
ora12cr1 18522      1 0 14:43 ?        00:00:00 ora_diag_ora12c
ora12cr1 18524      1 0 14:43 ?        00:00:00 ora_dbrm_ora12c
ora12cr1 18526      1 0 14:43 ?        00:00:00 ora_dia0_ora12c
ora12cr1 18528      1 0 14:43 ?        00:00:00 ora_dbw0_ora12c
ora12cr1 18530      1 0 14:43 ?        00:00:00 ora_lgwr_ora12c
ora12cr1 18532      1 0 14:43 ?        00:00:00 ora_ckpt_ora12c
ora12cr1 18534      1 0 14:43 ?        00:00:00 ora_lg00_ora12c
ora12cr1 18536      1 0 14:43 ?        00:00:00 ora_lg01_ora12c
ora12cr1 18538      1 0 14:43 ?        00:00:00 ora_smon_ora12c
ora12cr1 18540      1 0 14:43 ?        00:00:00 ora_reco_ora12c
ora12cr1 18542      1 0 14:43 ?        00:00:00 ora_lreg_ora12c
ora12cr1 18544      1 0 14:43 ?        00:00:00 ora_mmon_ora12c
ora12cr1 18546      1 0 14:43 ?        00:00:00 ora_mml_ora12c
ora12cr1 18547 18500  0 14:43 ?        00:00:00 oracleora12c
                                (DESCRIPTION=(LOCAL=YES)
                                (ADDRESS=(PROTOCOL=beq)))
ora12cr1 18566 18500  0 14:45 pts/0    00:00:00 /bin/bash -c ps -aef |
                                grep ora12cr1
ora12cr1 18567 18566  0 14:45 pts/0    00:00:00 ps -aef
ora12cr1 18568 18566  0 14:45 pts/0    00:00:00 grep ora12cr1
```

Кроме того, команда `ipcs` впервые сообщает об использовании разделяемой памяти и семафоров — двух важных устройств взаимодействия между процессами в системе UNIX/Linux:

```
SQL> !ipcs -a
----- Shared Memory Segments -----
----- Сегменты разделяемой памяти ---
key          shmid      owner      perms      bytes      nattch     status
0x10d1c894 13074435  ora12cr1   660        10485760   38
0x00000000 13107204  ora12cr1   660        322961408 19
```

----- Semaphore Arrays -----

----- Массивы семафоров -----

| key | semid | owner | perms | nsems |
|------------|--------|----------|-------|-------|
| 0xfc46e83c | 425986 | ora12cr1 | 660 | 171 |
| 0xfc46e83d | 458755 | ora12cr1 | 660 | 171 |
| 0xfc46e83e | 491524 | ora12cr1 | 660 | 171 |
| 0xfc46e83f | 524293 | ora12cr1 | 660 | 171 |
| 0xfc46e840 | 557062 | ora12cr1 | 660 | 171 |

----- Message Queues -----

----- Очереди сообщений -----

| key | msqid | owner | perms | used-bytes | messages |
|-----|-------|-------|-------|------------|----------|
|-----|-------|-------|-------|------------|----------|

Обратите внимание, что пока еще нет какой-либо “базы данных”. Мы имеем имя базы данных (в созданном файле параметров), но не саму базу данных. Если попробовать “смонтировать” эту базу данных, то попытка закончится неудачей, поскольку база данных попросту не существует. Давайте создадим ее. Как уже упоминалось, создание базы данных Oracle требует выполнения совсем небольшого количества действий, но судите сами:

```
SQL> create database;
```

```
Database created.
```

```
База данных создана.
```

Это действительно все, что нужно для создания базы данных. Однако в реальности, скорее всего, придется применять несколько более сложную форму команды CREATE DATABASE, т.к. программному обеспечению Oracle потребуются указать, куда поместить оперативные журнальные файлы, файлы данных, управляющие файлы и т.п. Но мы теперь имеем полноценно функционирующую базу данных. Для построения остальных словарей данных, используемых при повседневной работе (например, представления наподобие ALL_OBJECTS пока что не присутствуют в этой базе данных), понадобится еще запустить сценарий \$ORACLE_HOME/rdbms/admin/catalog.sql и другие сценарии, связанные с каталогами, однако сама база данных уже существует. Мы можем применять простой запрос к некоторым представлениям V\$ в Oracle, в частности V\$DATAFILE, V\$LOGFILE и V\$CONTROLFILE, для получения списка файлов, образующих эту базу данных:

```
SQL> select name from v$datafile;
```

```
NAME
```

```
-----
/home/ora12cr1/app/ora12cr1/product/12.1.0/dbhome_1/dbs/dbs1ora12c.dbf
/home/ora12cr1/app/ora12cr1/product/12.1.0/dbhome_1/dbs/dbx1ora12c.dbf
/home/ora12cr1/app/ora12cr1/product/12.1.0/dbhome_1/dbs/dbu1ora12c.dbf
```

```
SQL> select member from v$logfile;
```

```
MEMBER
```

```
-----
/home/ora12cr1/app/ora12cr1/product/12.1.0/dbhome_1/dbs/log1ora12c.dbf
/home/ora12cr1/app/ora12cr1/product/12.1.0/dbhome_1/dbs/log2ora12c.dbf
```

```
SQL> select name from v$controlfile;
```

```
NAME
```

```
-----
/home/ora12cr1/app/ora12cr1/product/12.1.0/dbhome_1/dbs/cntrlora12c.dbf
```

СУБД Oracle использовала стандартные значения параметров для сбора всей сведений воедино и создания базы данных в виде набора постоянных файлов. Если закрыть эту базу данных и попробовать открыть ее снова, выяснится, что это невозможно:

```
SQL> alter database close;
Database altered.
База данных изменена.
SQL> alter database open;
alter database open
*
ERROR at line 1:
ORA-16196: database has been previously opened and closed
ОШИБКА в строке 1:
ORA-16196: база данных была ранее открыта и закрыта
```

На протяжении периода своего существования экземпляр может монтировать и открывать максимум одну базу данных — с единственным владельцем или контейнерную. Вспомните, что экземпляр состоит просто из процессов и разделяемой памяти. Они запущены и работают. Все, что мы сделали — закрыли базу данных, т.е. физические файлы. Чтобы открыть текущую или любую другую базу данных, потребуется отбросить этот экземпляр (завершить) и создать новый.

Ниже перечислены моменты, которые необходимо запомнить.

- Экземпляр — это набор фоновых процессов и разделяемой памяти.
- База данных (с единственным владельцем или контейнерная) — это обособленная коллекция данных, хранящихся на диске.
- Экземпляр всегда может монтировать и открывать только одну базу данных. Как будет показано позже, он может предоставлять доступ к нескольким подключаемым базам данных, а также “открывать” и “закрывать” их много раз; тем не менее, экземпляр будет открывать только одну обособленную базу данных.
- База данных может быть смонтирована и открыта одним или большим числом экземпляров (с применением RAC), а количество экземпляров, монтирующих одну базу данных, со временем может изменяться.

Как отмечалось ранее, в большинстве случаев экземпляр и база данных связаны отношением “один к одному”. Вероятно, именно этим обстоятельством объясняется путаница в использовании этих терминов. Согласно личному опыту большинства людей, база данных — это экземпляр, а экземпляр — база данных.

Однако во многих тестовых средах это не так. Скажем, на моем диске может присутствовать пять отдельных баз данных. На тестовом компьютере в каждый конкретный момент времени функционирует только один экземпляр Oracle, но база данных, к которой он обращается, может отличаться изо дня в день или от часа к часу в зависимости от текущих потребностей. Просто имея множество разных файлов конфигурации, я могу монтировать и открывать любую одну из этих баз данных. Таким образом, существует один экземпляр, но много баз данных, только одна из которых доступна в каждый конкретный момент.

Итак, если теперь кто-то заговорит об экземпляре, то вы будете знать, что речь идет о процессах и памяти Oracle. Когда же упоминают базу данных, то имеют в виду физические файлы, которые хранят данные. База данных может быть доступна из множества экземпляров, но в каждый момент времени экземпляр будет предоставлять доступ только к одной базе данных (с единственным владельцем или контейнерной).

Системная глобальная область и фоновые процессы

Теперь вы, пожалуй, готовы к восприятию абстрактной картины того, на что похожи экземпляр и база данных Oracle (рис. 2.1).

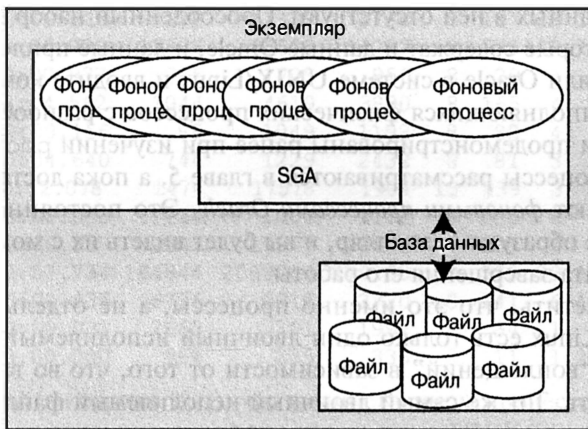


Рис. 2.1. Экземпляр и база данных Oracle

На рис. 2.1 показаны экземпляр и база данных Oracle в своей простейшей форме. В Oracle имеется большая область памяти под названием SGA (system global area — системная глобальная область), которая применяется, например, для выполнения следующих действий.

- Поддержка многих внутренних структур данных, к которым должны обращаться все процессы.
- Кеширование данных с диска; буферизация журнальных данных для восстановления перед их записью на диск.
- Хранение планов SQL-операторов, прошедших разбор.
- И т.д.

В Oracle существует набор процессов, которые “присоединены” к этой области SGA, и механизмы, с помощью которых производится присоединение, зависят от операционной системы. В среде UNIX/Linux процессы будут физически присоединяться к большому сегменту разделяемой памяти — области памяти, выделенной в операционной системе, которая одновременно может быть доступна множеству процессов (в основном, с использованием функций `shmget()` и `shmat()`).

В среде Windows для выделения памяти эти процессы просто вызывают С-функцию `malloc()`, поскольку в действительности они являются потоками внутри одного большого процесса и, следовательно, разделяют одно и то же пространство виртуальной памяти.

СУБД Oracle также имеет набор файлов, в которых процессы или потоки базы данных выполняют чтение и запись (чтение или запись в эти файлы имеют право осуществлять только процессы Oracle). В архитектуре с единственным владельцем эти файлы содержат все табличные данные, индексы, временные области, журналы для восстановления, код PL/SQL и т.д. В архитектуре с множеством владельцев контейнерная база данных будет хранить все относящиеся к Oracle метаданные, данные и код; данные вашего приложения будут содержаться отдельно в подключаемой базе данных, которую нам еще предстоит создать. База данных, созданная нами ранее, по умолчанию является базой данных с единственным владельцем; доступные подключаемые базы данных в ней отсутствуют. Обособленный набор файлов включает файлы данных, которые содержат и данные Oracle, и данные приложения.

Если вы запустили Oracle в системе UNIX/Linux и вводите команду `ps`, то увидите множество выполняющихся физических процессов с разнообразными именами. Примеры были продемонстрированы ранее при изучении `pmon`, `smon` и других процессов. Эти процессы рассматриваются в главе 5, а пока достаточно знать, что обычно их называют *фоновыми процессами Oracle*. Это постоянно существующие процессы, которые образуют экземпляр, и вы будете видеть их с момента запуска экземпляра до момента завершения его работы.

Интересно отметить, что это именно процессы, а не отдельные программы. В системе UNIX/Linux есть только один двоичный исполняемый файл Oracle; он имеет множество “воплощений” в зависимости от того, что во время запуска ему было указано делать. Тот же самый двоичный исполняемый файл, который запускался для создания процесса `ora_pmon_oral2c`, также применялся для создания процесса `ora_ckpt_oral2c`. Существует только одна двоичная исполняемая программа, которая называется просто `oracle`. Она всего лишь выполняется много раз с разными именами.

В среде Windows с помощью инструмента `pstat` (часть набора Windows XP Resource Kit; если вы еще не располагаете этим набором, поищите в Интернете “`pstat download`”) мы найдем только один процесс — `oracle.exe`. Опять-таки, в среде Windows существует только один двоичный исполняемый файл (`oracle.exe`). Внутри этого процесса мы обнаружим множество потоков, представляющих фоновые процессы Oracle.

С помощью `pslist` (или любого другого аналогичного инструмента, такого как `tasklist`, который поступает со многими версиями Windows) эти потоки можно просмотреть:

```
C:\WINDOWS> pstat
```

```
Pstat version 0.3: memory: 523760 kb uptime: 0 1:37:54.375
```

```
PageFile: \??\C:\pagefile.sys
```

```
Current Size: 678912 kb Total Used: 228316 kb Peak Used 605488 kb
```

```
Memory: 523760K Avail: 224492K TotalWs: 276932K InRam Kernel: 872K P:20540K  
Commit: 418468K/ 372204K Limit:1169048K Peak:1187396K Pool N:10620K P:24588K
```

| User Time | Kernel Time | Ws | Faults | Commit | Pri | HndThd | Pid | Name |
|-------------|-------------|--------|---------|--------|-----|--------|-----|----------------------|
| | | 56860 | 2348193 | | | | | File Cache |
| 0:00:00.000 | 1:02:23.109 | 28 | 0 | 0 | 0 | 0 | 1 | 0Idle Process |
| 0:00:00.000 | 0:01:50.812 | 32 | 4385 | 28 | 8 | 694 | 52 | 4System |
| 0:00:00.015 | 0:00:00.109 | 60 | 224 | 172 | 11 | 19 | 3 | 332smss.exe |
| 0:00:33.234 | 0:00:32.046 | 2144 | 33467 | 1980 | 13 | 396 | 14 | 556csrss.exe |
| 0:00:00.343 | 0:00:01.750 | 3684 | 6811 | 7792 | 13 | 578 | 20 | 580winlogon.exe |
| 0:00:00.078 | 0:00:01.734 | 1948 | 3022 | 1680 | 9 | 275 | 16 | 624services.exe |
| 0:00:00.218 | 0:00:03.515 | 1896 | 5958 | 3932 | 9 | 363 | 25 | 636lsass.exe |
| 0:00:00.015 | 0:00:00.078 | 80 | 804 | 592 | 8 | 25 | 1 | 812vmacthlp.exe |
| 0:00:00.093 | 0:00:00.359 | 1416 | 2765 | 3016 | 8 | 195 | 17 | 828svchost.exe |
| 0:00:00.062 | 0:00:00.453 | 1340 | 3566 | 1764 | 8 | 244 | 10 | 896svchost.exe |
| 0:00:00.828 | 0:01:16.593 | 9632 | 36387 | 11708 | 8 | 1206 | 59 | 1024svchost.exe |
| 0:00:00.046 | 0:00:00.640 | 1020 | 2315 | 1300 | 8 | 81 | 6 | 1100svchost.exe |
| 0:00:00.015 | 0:00:00.234 | 736 | 2330 | 1492 | 8 | 165 | 11 | 1272svchost.exe |
| 0:00:00.015 | 0:00:00.218 | 128 | 1959 | 3788 | 8 | 117 | 10 | 1440spoolsv.exe |
| 0:00:01.312 | 0:00:19.828 | 13636 | 35525 | 14732 | 8 | 575 | 19 | 1952explorer.exe |
| 0:00:00.250 | 0:00:00.937 | 956 | 1705 | 856 | 8 | 29 | 1 | 228VMwareTray.exe |
| 0:00:00.812 | 0:00:04.562 | 1044 | 4619 | 3800 | 8 | 165 | 4 | 240VMwareUser.exe |
| 0:00:00.015 | 0:00:00.156 | 88 | 1049 | 1192 | 8 | 88 | 4 | 396svchost.exe |
| 0:00:00.109 | 0:00:04.640 | 744 | 1229 | 2432 | 8 | 81 | 3 | 460cvpnd.exe |
| 0:00:02.015 | 0:00:12.078 | 1476 | 17578 | 1904 | 13 | 139 | 3 | 600VMwareService.exe |
| 0:00:00.031 | 0:00:00.093 | 124 | 1004 | 1172 | 8 | 105 | 6 | 192alg.exe |
| 0:00:00.062 | 0:00:00.937 | 2648 | 13977 | 22656 | 8 | 101 | 3 | 720TNSLSNR.EXE |
| 0:04:00.359 | 0:02:57.734 | 164844 | 2009785 | 279168 | 8 | 550 | 29 | 1928oracle.exe |
| 0:00:00.093 | 0:00:00.437 | 6736 | 2316 | 2720 | 8 | 141 | 6 | 1224msiexec.exe |
| 0:00:00.015 | 0:00:00.031 | 2668 | 701 | 1992 | 8 | 34 | 1 | 804cmd.exe |
| 0:00:00.015 | 0:00:00.000 | 964 | 235 | 336 | 8 | 11 | 1 | 2856pstat.exe |

Здесь мы видим (в колонке Thd), что в единственном процессе Oracle содержатся 29 потоков. Эти потоки представляют то, что в среде UNIX/Linux является процессами — pmon, arch, lgwr и т.д. Каждый из них представляет отдельную часть процесса Oracle. Продвигаясь далее по отчету pstat, мы можем просмотреть дополнительные сведения о каждом потоке:

```
pid:788 pri: 8 Hnd: 550 Pf:2009785 Ws: 164844K oracle.exe
tid pri Ctx Swtch StrtAddr User Time Kernel Time State
498 9 651 7C810705 0:00:00.000 0:00:00.203 Wait:Executive
164 8 91 7C8106F9 0:00:00.000 0:00:00.000 Wait:UserRequest
...
a68 8 42 7C8106F9 0:00:00.000 0:00:00.031 Wait:UserRequest
```

В этом выводе отсутствуют “имена” потоков, подобные отображаемым в UNIX/Linux (ora_pmon_oral2c и т.п.), но есть идентификаторы потоков (Tid), приоритеты (Pri) и другая учетная информация о потоках операционной системы.

Подключение к Oracle

В этом разделе мы рассмотрим механизмы, лежащие в основе двух наиболее распространенных методов обслуживания запросов сервером Oracle: подключения к *выделенному серверу* и подключения к *разделяемому серверу*. Мы выясним, что происходит на компьютерах клиента и сервера для установки подключений, чтобы можно было выполнить вход и делать действительную работу в базе данных. Наконец, мы кратко

рассмотрим установку соединений TCP/IP (TCP/IP — это основной сетевой протокол, используемый для подключения к базе данных Oracle по сети) и особенности работы прослушивающего процесса, который отвечает за установление физического подключения к серверу, в случае подключений к выделенному и разделяемому серверу.

Выделенный сервер

Диаграмма на рис. 2.1 и следующий вывод команды `ps` представляет картину того, как база данных Oracle по имени `oral2cr1` может выглядеть непосредственно после запуска.

```
[tkyte@dellpe]$ ps -aef | grep _$ORACLE_SID
oral2cr1 19607      1  0 15:16 ?        00:00:00 ora_pmon_oral2cr1
oral2cr1 19609      1  0 15:16 ?        00:00:00 ora_psp0_oral2cr1
oral2cr1 19611      1  0 15:16 ?        00:00:00 ora_vktm_oral2cr1
oral2cr1 19615      1  0 15:16 ?        00:00:00 ora_gen0_oral2cr1
oral2cr1 19617      1  1 15:16 ?        00:00:00 ora_mman_oral2cr1
oral2cr1 19621      1  0 15:16 ?        00:00:00 ora_diag_oral2cr1
oral2cr1 19623      1  0 15:16 ?        00:00:00 ora_dbrm_oral2cr1
oral2cr1 19625      1  0 15:16 ?        00:00:00 ora_dia0_oral2cr1
oral2cr1 19627      1  0 15:16 ?        00:00:00 ora_dbw0_oral2cr1
oral2cr1 19629      1  0 15:16 ?        00:00:00 ora_lgwr_oral2cr1
oral2cr1 19631      1  0 15:16 ?        00:00:00 ora_ckpt_oral2cr1
oral2cr1 19633      1  0 15:16 ?        00:00:00 ora_lg00_oral2cr1
oral2cr1 19635      1  0 15:16 ?        00:00:00 ora_lg01_oral2cr1
oral2cr1 19637      1  0 15:16 ?        00:00:00 ora_smon_oral2cr1
oral2cr1 19639      1  0 15:16 ?        00:00:00 ora_reco_oral2cr1
oral2cr1 19641      1  0 15:16 ?        00:00:00 ora_lreg_oral2cr1
oral2cr1 19643      1  1 15:16 ?        00:00:00 ora_mmon_oral2cr1
oral2cr1 19645      1  0 15:16 ?        00:00:00 ora_mnnl_oral2cr1
oral2cr1 19647      1  0 15:16 ?        00:00:00 ora_d000_oral2cr1
oral2cr1 19649      1  0 15:16 ?        00:00:00 ora_s000_oral2cr1
oral2cr1 19661      1  0 15:16 ?        00:00:00 ora_tmon_oral2cr1
oral2cr1 19663      1  0 15:16 ?        00:00:00 ora_tt00_oral2cr1
oral2cr1 19665      1  0 15:16 ?        00:00:00 ora_smco_oral2cr1
oral2cr1 19667      1  0 15:16 ?        00:00:00 ora_fbda_oral2cr1
oral2cr1 19671      1  0 15:16 ?        00:00:00 ora_agpc_oral2cr1
oral2cr1 19675      1  0 15:16 ?        00:00:00 ora_cjq0_oral2cr1
oral2cr1 19705      1  0 15:16 ?        00:00:00 ora_w000_oral2cr1
oral2cr1 19708      1  0 15:16 ?        00:00:00 ora_qm02_oral2cr1
oral2cr1 19710      1  0 15:16 ?        00:00:00 ora_qm03_oral2cr1
oral2cr1 19712      1  0 15:16 ?        00:00:00 ora_q002_oral2cr1
oral2cr1 19714      1  0 15:16 ?        00:00:00 ora_q003_oral2cr1
tkyte    19720 15416  0 15:16 pts/1    00:00:00 grep _oral2cr1
```

Если вы вошли в эту базу данных с применением выделенного сервера, то увидите новый процесс (или поток в другой операционной системе), который был создан как раз для вашего обслуживания:

```
OPS$TKYTE@ORA12CR1> !ps -aef | grep $ORACLE_SID
oral2cr1 19607      1  0 15:16 ?        00:00:00 ora_pmon_oral2cr1
oral2cr1 19609      1  0 15:16 ?        00:00:00 ora_psp0_oral2cr1
oral2cr1 19611      1  0 15:16 ?        00:00:00 ora_vktm_oral2cr1
oral2cr1 19615      1  0 15:16 ?        00:00:00 ora_gen0_oral2cr1
```

```

ora12cr1 19617      1 0 15:16 ?      00:00:00 ora_mman_ora12cr1
ora12cr1 19621      1 0 15:16 ?      00:00:00 ora_diag_ora12cr1
ora12cr1 19623      1 0 15:16 ?      00:00:00 ora_dbrm_ora12cr1
ora12cr1 19625      1 0 15:16 ?      00:00:00 ora_dia0_ora12cr1
ora12cr1 19627      1 0 15:16 ?      00:00:00 ora_dbw0_ora12cr1
ora12cr1 19629      1 0 15:16 ?      00:00:00 ora_lgwr_ora12cr1
ora12cr1 19631      1 0 15:16 ?      00:00:00 ora_ckpt_ora12cr1
ora12cr1 19633      1 0 15:16 ?      00:00:00 ora_lg00_ora12cr1
ora12cr1 19635      1 0 15:16 ?      00:00:00 ora_lg01_ora12cr1
ora12cr1 19637      1 0 15:16 ?      00:00:00 ora_smon_ora12cr1
ora12cr1 19639      1 0 15:16 ?      00:00:00 ora_reco_ora12cr1
ora12cr1 19641      1 0 15:16 ?      00:00:00 ora_lreg_ora12cr1
ora12cr1 19643      1 0 15:16 ?      00:00:00 ora_mmon_ora12cr1
ora12cr1 19645      1 0 15:16 ?      00:00:00 ora_mnnl_ora12cr1
ora12cr1 19647      1 0 15:16 ?      00:00:00 ora_d000_ora12cr1
ora12cr1 19649      1 0 15:16 ?      00:00:00 ora_s000_ora12cr1
ora12cr1 19661      1 0 15:16 ?      00:00:00 ora_tmon_ora12cr1
ora12cr1 19663      1 0 15:16 ?      00:00:00 ora_tt00_ora12cr1
ora12cr1 19665      1 0 15:16 ?      00:00:00 ora_smco_ora12cr1
ora12cr1 19667      1 0 15:16 ?      00:00:00 ora_fbda_ora12cr1
ora12cr1 19671      1 0 15:16 ?      00:00:00 ora_aqpc_ora12cr1
ora12cr1 19675      1 0 15:16 ?      00:00:00 ora_cjq0_ora12cr1
ora12cr1 19705      1 0 15:16 ?      00:00:00 ora_w000_ora12cr1
ora12cr1 19708      1 0 15:16 ?      00:00:00 ora_qm02_ora12cr1
ora12cr1 19712      1 0 15:16 ?      00:00:00 ora_q002_ora12cr1
ora12cr1 19714      1 0 15:16 ?      00:00:00 ora_q003_ora12cr1
tkyte      19732 19731 0 15:17 pts/0 00:00:00 /home/ora12cr1/app/ora12cr1/
product/12.1.0/dbhome_1/bin/sqlplus
ora12cr1 19733 19732 0 15:17 ?      00:00:00 oracleora12cr1
(DESCRIPTION=(LOCAL=YES) (ADDRESS=(PROTOCOL=beq) ) )
tkyte      19744 19742 0 15:18 pts/0 00:00:00 grep ora12cr1

```

Теперь вы можете отметить наличие нового процесса `oracleora12cr1`, созданного как процесс выделенного сервера. После выхода из базы данных дополнительный поток/процесс исчезнет.

Все сказанное подводит нас к следующей итерации диаграммы. Подключение к базе данных Oracle в ее наиболее распространенной конфигурации в результате дает то, что показано на рис. 2.2.

Как уже отмечалось, обычно Oracle будет создавать новый процесс после входа пользователя в базу данных. В общем случае это называют конфигурацией с выделенным сервером, поскольку процесс сервера будет предназначен специально для обслуживания пользователя на протяжении всего его сеанса. Для каждого сеанса будет создаваться новый выделенный сервер. По определению такой процесс выделенного сервера не является частью экземпляра. Клиентский процесс (любая программа, пытающаяся подключиться к базе данных) будет напрямую взаимодействовать с этим выделенным сервером через определенный сетевой канал, подобный сокету TCP/IP. Именно этот процесс сервера будет принимать от пользователя SQL-команды и выполнять их. При необходимости он будет читать файлы данных и искать нужные данные в кеше базы данных. Он будет выполнять операторы обновления и код PL/SQL. Его единственное назначение — реагировать на отправляемые пользователем SQL-операторы.

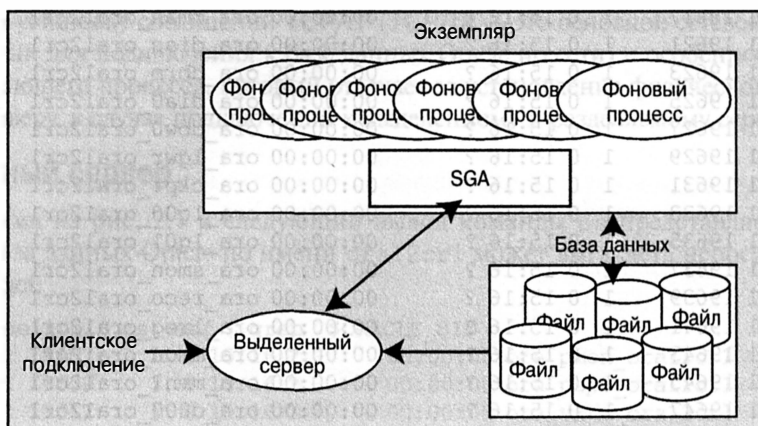


Рис. 2.2. Типичная конфигурация с выделенным сервером

Разделяемый сервер

СУБД Oracle может также принимать подключения в стиле так называемого *разделяемого сервера*, при котором для каждого нового подключения пользователя дополнительный поток или новый процесс UNIX/Linux появляться не будет.

На заметку! В версиях Oracle 7.x и Oracle 8.x разделяемый сервер назывался MTS (multithreaded server — многопоточный сервер). Это устаревшее название больше не используется.

В случае разделяемого сервера Oracle применяет пул разделяемых процессов для большой группы пользователей. Разделяемые серверы являются просто механизмом организации пула подключений. Вместо использования для 10 000 сеансов базы данных 10 000 выделенных серверов (что предполагает наличие многочисленных процессов или потоков) разделяемый сервер позволяет иметь небольшой процент от этого множества процессов или потоков, которые совместно применяются (разделяются) всеми сеансами. В итоге у Oracle появляется возможность подключения значительно большего числа пользователей, чем это было бы в противном случае. Компьютер может не справиться с нагрузкой, связанной с управлением 10 000 процессами, но управление 100 или 1000 процессами — вполне реальная задача. В режиме разделяемого сервера разделяемые процессы в основном запускаются базой данных и появляются в списке `ps`.

Существенное различие между подключениями посредством разделяемого и выделенного сервера состоит в том, что клиентский процесс, подключенный к базе данных, никогда не обращается напрямую к разделяемому серверу, как это было бы в случае выделенного сервера. Он не может взаимодействовать с разделяемым сервером, потому что фактически этот процесс является разделяемым. Для совместного использования таких процессов необходим еще один механизм, через который производится взаимодействие. В Oracle для этой цели применяется процесс (или набор процессов), называемый *диспетчером*. Клиентский процесс будет обмениваться данными с процессом диспетчера по сети. Процесс диспетчера поместит запрос клиента в очередь запросов внутри области SGA (в этом заключается одна из многих

функций SGA). Первый незанятый разделяемый сервер примет этот запрос и обработает его (например, запросом мог бы служить `UPDATE T SET X=X+5 WHERE Y=2`). По завершении этой команды разделяемый сервер поместит ответ в очередь ответов вызывающего диспетчера. Процесс диспетчера будет отслеживать эту очередь и, обнаружив результат, переправит его клиенту. В общих чертах прохождение запроса к разделяемому серверу выглядит так, как показано на рис. 2.3.

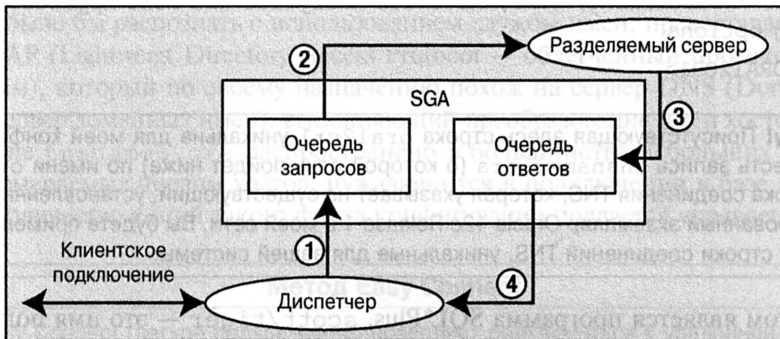


Рис. 2.3. Шаги выполнения запроса к разделяемому серверу

На рис. 2.3 видно, что клиентское подключение отправляет запрос диспетчеру. Диспетчер сначала поместит этот запрос в очередь запросов внутри SGA (1). Первый свободный разделяемый сервер извлечет этот запрос из очереди (2) и обработает его. Когда разделяемый сервер завершит обработку, ответ (коды возврата, данные и т.д.) будет помещен в очередь ответов (3), после чего извлечен оттуда диспетчером (4) и передан обратно клиенту.

С точки зрения разработчика подключение посредством выделенного сервера концептуально ничем не отличается от подключения с использованием разделяемого сервера. Архитектурно они существенно разнятся, но приложению это не видно.

Теперь, когда вы понимаете, что собой представляют подключения через выделенный и через разделяемый сервер, могут возникнуть перечисленные ниже вопросы.

- Как осуществить собственно подключение?
- Как запустить выделенный сервер?
- Как связаться с диспетчером?

Ответы зависят от применяемой платформы, но общее в последующих разделах приведено общее описание процесса.

Механизмы подключения через TCP/IP

Рассмотрим наиболее распространенный случай обмена данными по сети — запрос сетевого подключения через TCP/IP. В этом случае клиент находится на одном компьютере, а сервер — на другом, и оба компьютера подключены к сети TCP/IP. Все действия инициируются клиентом. Клиент делает запрос на подключение к базе данных, используя клиентское программное обеспечение Oracle (набор готовых программных интерфейсов приложения, или API-интерфейсов). Например, клиент выдает следующий запрос:

```
[tkyte@dellpe ~]$ sqlplus scott/tiger@oral2cr1
SQL*Plus: Release 12.1.0.1.0 Production on Mon Sep 2 15:25:06 2013
Copyright (c) 1982, 2013, Oracle. All rights reserved.
Last Successful login time: Mon Sep 02 2013 13:44:49 -04:00
Connected to:
Oracle Database 12c Enterprise Edition Release 12.1.0.1.0 - 64bit Production
With the Partitioning, OLAP, Advanced Analytics and Real Application
Testing options
SCOTT@ORA12CR1>
```

На заметку! Присутствующая здесь строка `oral2cr1` уникальна для моей конфигурации. У меня есть запись `tnsnames.ora` (о которой речь пойдет ниже) по имени `oral2cr1`. Это строка соединения TNS, которая указывает на существующий, установленный и сконфигурированный экземпляр Oracle 12c Release 1 в моей сети. Вы будете применять собственные строки соединений TNS, уникальные для вашей системы.

Клиентом является программа SQL*Plus, `scott/tiger` — это имя пользователя и пароль, а `oral2cr1` — имя службы TNS. Аббревиатура TNS расшифровывается как Transparent Network Substrate (прозрачная сетевая среда) и представляет собой “базисное” программное обеспечение, встроенное в клиент Oracle, которое обрабатывает удаленные подключения, обеспечивая возможность одноранговых коммуникаций. Строка подключения TNS сообщает программному обеспечению Oracle, каким образом следует подключаться к удаленной базе данных. Как правило, клиентская программа, выполняющаяся на локальном компьютере, будет читать файл по имени `tnsnames.ora`. Это простой текстовый конфигурационный файл, который обычно расположен в каталоге `$ORACLE_HOME/network/admin` (`$ORACLE_HOME` представляет полный путь к каталогу с установленной копией Oracle). Он будет содержать записи, подобные показанным ниже:

```
SCOTT@ORA12CR1> !cat $ORACLE_HOME/network/admin/tnsnames.ora
# tnsnames.ora Network Configuration File:
# /home/oral2cr1/app/oral2cr1/product/12.1.0/dbhome_1/network/admin/tnsnames.ora
# Generated by Oracle configuration tools.
# tnsnames.ora Файл конфигурации сети:
# Сгенерировано инструментами конфигурирования Oracle.

ORA12CR1 =
  (DESCRIPTION =
    (ADDRESS = (PROTOCOL = TCP) (HOST = somehost.somewhere.com) (PORT = 1521))
    (CONNECT_DATA =
      (SERVER = DEDICATED)
      (SERVICE_NAME = oral2cr1)
    )
  )
```

Эта конфигурационная информация позволяет клиентской программе Oracle сопоставить указанную строку соединения TNS (`oral2cr1`) с более полезной информацией: именем хоста, номером порта на этом хосте, через который прослушивающий процесс будет принимать подключения, именем службы базы данных на этом хосте, к которой нужно подключиться, и т.д. Имя службы представляет группы

приложений с общими атрибутами, порогами уровня обслуживания и приоритетами. Круг экземпляров, предлагающих службу, прозрачен для приложения, и каждый экземпляр базы данных может быть зарегистрирован прослушивающим процессом как готовый предоставлять множество служб. Таким образом, службы отображаются на экземпляры физической базы данных и позволяют администратору базы данных связать с ними определенные пороговые значения и приоритеты.

Эта строка, `ora12cr1`, может быть распознана и другими способами. Например, ее можно было бы распознать с использованием службы имен, предоставляемой сервером LDAP (Lightweight Directory Access Protocol — облегченный протокол доступа к каталогам), который по своему назначению похож на сервер DNS (Domain Name Server — сервер доменных имен), выполняющий преобразование имен хостов. Однако применение `tnsnames.ora` распространено в большинстве небольших и средних систем (измеряемых числом хостов, нуждающихся в подключении к базе данных), в которых количество копий такого файла конфигурации поддается управлению.

Метод Easy Connect

Метод Easy Connect (легкое подключение) позволяет подключаться к удаленной базе данных без необходимости в наличии файла `tnsnames.ora` (или других методов распознавания местоположения базы данных). Если известны имя хоста, сервер, порт и имя службы, то их можно ввести непосредственно в командной строке. Синтаксис выглядит следующим образом:

```
sqlplus имя_пользователя[/[/]хост[:порт] [/имя_службы] [:сервер] [/имя_экземпляра]
```

Например, предположим, что именем хоста является `hesta`, номером порта — 1521, а именем службы — `ora12cr1`. Тогда подключиться можно так:

```
$ sqlplus user/pass@hesta:1521/ora12cr1
```

Метод Easy Connect удобно использовать в ситуациях, когда проводится поиск и устранение проблем с подключаемостью или нет доступного файла `tnsnames.ora` (либо других путей распознавания удаленного подключения).

Теперь, когда клиентской программе известно, куда подключаться, она откроет сокетное подключение TCP/IP к серверу с именем хоста `somehost.somewhere.com` через порт 1521 (стандартный номер порта для прослушивания; можно сконфигурировать запуск и на другом порте). Если администратор базы данных сервера установил и сконфигурировал сеть Oracle Net, а прослушивающий процесс прослушивает порт 1521 на предмет поступления запросов о подключении, то это подключение может быть принято. В сетевой среде на сервере будет действовать процесс, называемый *прослушивающим процессом TNS*. Именно этот прослушивающий процесс будет обеспечивать физическое подключение к базе данных. Получив входящий запрос на подключение, прослушивающий процесс первым делом анализирует запрос. Затем, основываясь на собственных файлах конфигурации, он либо отклоняет запрос (например, из-за отсутствия такой службы или по причине запрета IP-адресу подключаться к этому хосту), либо принимает и производит подключение.

При подключении посредством выделенного сервера прослушивающий процесс создаст выделенный сервер. В системе UNIX/Linux это осуществляется с помощью системных вызовов `fork()` и `exec()` (единственным способом создания но-

вого процесса после инициализации в UNIX/Linux является применение `fork()`. Новый процесс выделенного сервера наследует подключение, установленное прослушивающим процессом, в результате появляется физическое подключение к базе данных. В среде Windows прослушивающий процесс запрашивает у процесса базы данных создание нового потока для подключения. После того как этот поток создан, клиент “перенаправляется” на него, устанавливая в итоге физическое подключение. Описанные выше действия для системы UNIX/Linux можно проиллюстрировать диаграммой, представленной на рис. 2.4.

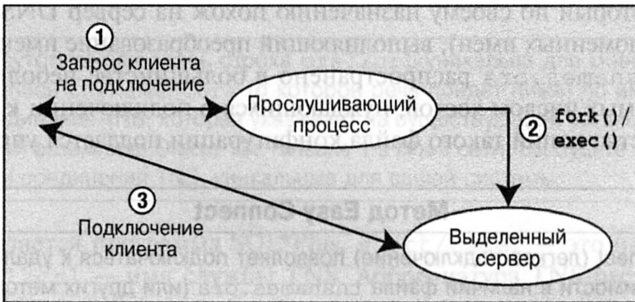


Рис. 2.4. Прослушивающий процесс и подключения посредством выделенного сервера

Тем не менее, прослушивающий процесс будет вести себя по-другому, если выдается запрос на подключение с использованием разделяемого сервера. Прослушивающему процессу известен диспетчер (или диспетчеры), функционирующий в экземпляре. При получении запросов на подключение прослушивающий процесс будет выбирать процесс диспетчера из пула доступных диспетчеров. Он либо возвратит клиенту информацию о подключении, описывающую то, как клиент может подключиться к процессу диспетчера, либо, если это возможно, передаст подключение процессу диспетчера (поведение зависит от версий операционной системы и СУБД, но конечный результат будет таким же). Когда прослушивающий процесс отправляет информацию о подключении, это делается потому, что он выполняется на хорошо известном хосте и порте этого хоста, но диспетчеры также принимают подключения от произвольно назначаемых портов на данном сервере. Диспетчер осведомлен об этих произвольных назначениях портов диспетчером и самостоятельно выберет диспетчера. Затем клиент отключается от прослушивающего процесса и подключается непосредственно к диспетчеру. В итоге устанавливается физическое подключение к базе данных. Описанный процесс показан на рис. 2.5.

Подключаемые базы данных

Начиная с версии Oracle Database 12c, появился новый подход к трактовке концепции “база данных”.

До сих пор в этой главе мы были сосредоточены на базах данных с единственным владельцем или контейнерных базах данных и их связью с экземплярами (база данных может иметь один и более экземпляров; экземпляр будет монтировать и открывать максимум одну базу данных).

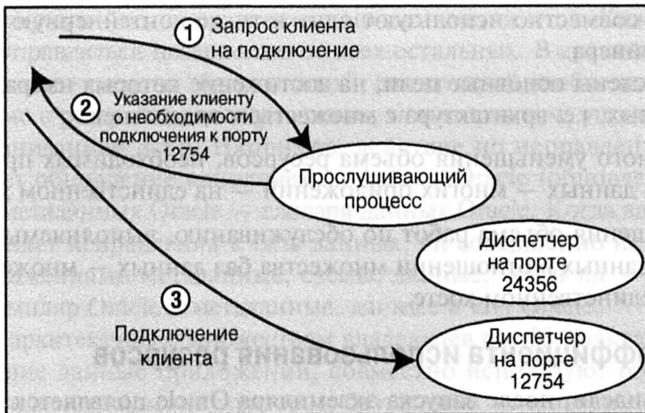


Рис. 2.5. Прослушивающий процесс и подключения посредством разделяемого сервера

Теперь мы готовы взглянуть на концепцию подключаемой базы данных — что она собой представляет и как работает. Подключаемые базы данных в архитектуре с множеством владельцев являются необособленными наборами файлов данных, которые состоят исключительно из данных и метаданных приложений. В них нет каких-либо специфичных для Oracle данных/метаданных; эта информация находится в контейнерных базах данных, которые в текущий момент ассоциированы с ними. Для того чтобы подключаемую базу данных можно было использовать и запрашивать, она должна быть связана с контейнерной базой данных. Эта контейнерная база данных будет содержать в себе только данные и метаданные Oracle — просто информацию, которая нужна Oracle для “запуска”. Подключаемые базы данных хранят “остальные” метаданные и данные базы данных.

Таким образом, например, контейнерная база данных могла бы иметь определение пользователя SYS (метаданные для пользователя SYS), а также скомпилированный и исходный код для объектов, подобных DBMS_OUTPUT и UTL_FILE. С другой стороны, подключаемая база данных содержала бы определение схемы приложения, такой как SCOTT, все метаданные, описывающие таблицы в схеме SCOTT, весь исходный код PL/SQL для схемы SCOTT, все привилегии, выданные схеме SCOTT, и т.д. Одним словом, подключаемая база данных имеет все, что описывает набор схем приложения — метаданные для учетных записей, метаданные для таблиц в этих учетных записях и действительные данные для этих таблиц. Подключаемая база данных является обособленной в отношении содержащихся в ней учетных записей приложения, но чтобы иметь возможность “открываться” и запрашиваться, ей необходима контейнерная база данных. Следовательно, можно сказать, что подключаемая база данных “необособленна”; для открытия и работы с ней требуется кое-что еще.

Подключаемая база данных не открывается экземпляром напрямую — взамен должен быть запущен экземпляр Oracle, который смонтирует и откроет контейнерную базу данных. После того, как экземпляр контейнера запущен, а контейнерная база данных открыта, она может открывать до 250 отдельных подключаемых баз данных. Каждая такая подключаемая база данных действует так, как если бы она была “автономной” базой данных. То есть все они выглядят так, как будто бы являются обособленными, автономными базами данных “с единственным владельцем”.

Однако все они совместно используют одну и ту же контейнерную базу данных и экземпляр контейнера.

Ниже перечислены основные цели, на достижение которых направлена подключаемая база данных, т.е. архитектура с множеством владельцев.

- Для заметного уменьшения объема ресурсов, необходимых при размещении многих баз данных — многих приложений — на единственном хосте.
- Для сокращения объема работ по обслуживанию, выполняемых администратором баз данных в отношении множества баз данных — множества приложений — на единственном хосте.

Снижение коэффициента использования ресурсов

Как вы уже видели, после запуска экземпляра Oracle появляется много ассоциированных с ним процессов. В главе 5 вы ознакомитесь с каждым из них и узнаете, что они делают; но уже сейчас несложно заметить, что каждый экземпляр поддерживается с помощью 20–40 процессов. Если вы попытаетесь настроить 50 баз данных с единственным владельцем (где каждая база данных имеет связанный с ней экземпляр или собственный экземпляр), то получите свыше 1000 процессов, которые нужны просто для запуска этих баз данных! Это порождает исключительно высокую нагрузку на операционную систему, как в плане создания множества процессов, так и в плане управления ими.

Кроме того, каждый экземпляр может иметь собственную область SGA. В главе 4 будет показано, что находится в SGA, но достаточно сказать, что там присутствует множество дублированных данных. Каждая область SGA будет иметь кешированную копию DBMS_OUTPUT в своем разделяемом пуле, а также журнальный буфер и много других дублирующих структур данных.

Благодаря подключаемым базам данных, вы получаете возможность разделить относящиеся к приложению метаданные, пользователей, данные, код и тому подобное, но при этом избежать наличия избыточных экземпляров. Другими словами, вы можете иметь единственный экземпляр с одной контейнерной базой данных (метаданные, код и данные Oracle), предоставляющей доступ к 250 подключаемым базам данных, в каждой из которых размещается отдельное приложение. Вместо 1000 процессов для поддержки 50 разных баз данных приложений у вас будет возможность располагать 51 базой данных (одной контейнерной и 50 базами данных приложений), которые разделяют те же самые 20–40 процессов. Это крупное снижение коэффициента утилизации ресурсов. Вдобавок все они разделяют общую область SGA, позволяя совместно использовать повторяющиеся части 50 отдельных областей SGA, которые существовали бы в противном случае. В целом размер одной области SGA, которую понадобилось бы выделить для этих 50 баз данных приложений, будет меньше суммарного размера 50 отдельных областей SGA, которые пришлось бы выделять иначе.

Сокращение объема работ по обслуживанию

Если администратору баз данных было поручено управление 50 отдельными базами данных в рамках архитектуры с единственным владельцем, у него окажется 50 баз данных, для которых придется выполнять действия по резервному копированию,

мониторингу, управлению, исправлению, модернизации и т.д. Каждая база данных должна будет управляться независимо от всех остальных. В архитектуре с множеством владельцев имеется одна “база данных”, для которой необходимо проводить работы, связанные с резервным копированием, мониторингом, управлением, исправлением, модернизацией и т.п. Например, действие по исправлению базы данных Oracle включает обновление исполняемых файлов Oracle (обновление экземпляра) и обновление метаданных Oracle — словаря данных Oracle. Когда администратор баз данных применяет исправления к базе данных, он совершенно не затрагивает связанные с приложениями метаданные, схемы, данные, код и т.д. — он воздействует только на экземпляр Oracle и метаданные, данные и код Oracle.

В условиях архитектуры с множеством владельцев все 50 подключаемых баз данных, содержащие данные приложений, совместно используют общий экземпляр; поэтому при внесении администратором баз данных исправлений в экземпляр ко всем 50 подключаемым базам данных применяются исправления со стороны экземпляра. Аналогично, когда администратор баз данных вносит исправления в контейнерную базу данных, т.е. обособленный набор файлов, которые хранят метаданные, код и данные Oracle, эти исправления наследуются всеми 50 подключаемыми базами данных. Они разделяют метаданные, код и данные Oracle. Таким образом, действие по исправлению единственной контейнерной базы данных фактически приводит к применению исправлений ко всем связанным подключаемым базам данных.

Если эта ситуация нежелательна (исправление всех 50 баз данных за раз), можно воспользоваться альтернативным подходом. Вместо внесения исправлений в существующую контейнерную базу данных администратор может просто создать новую контейнерную базу данных, которая уже исправлена. Теперь есть два экземпляра Oracle с двумя контейнерными базами данных: одна версии X, а другая версии Y. Для применения исправлений к подключаемой базе данных администратор “отключает” ее от контейнерной базы данных, к которой она в текущий момент присоединена. Отключение приводит к созданию XML-файла манифеста с описаниями файлов, принадлежащих к подключаемой базе данных. После этого администратор “подключает” эти подключаемые данные к новой контейнерной базе данных. Действие по отключению подключаемой базы данных является легковесным и очень быстрым — необходимо лишь создать XML-файл манифеста. Действие по подключению тоже отличается легковесностью и высокой скоростью — XML-файл манифеста читается, файлы, которые ассоциированы с подключаемой базой данных, регистрируются в контейнерной базе данных и подключаемая база данных может использоваться снова. Подключаемая база данных исправляется или модернизируется просто путем ее подключения к контейнерной базе данных, к которой уже применены исправления. В итоге администратору приходится управлять, исправлять, модернизировать и выполнять другие действия по отношению к одной контейнерной базе данных. Подключаемые базы данных всего лишь наследуют эту работу.

Отличия подключаемой базы данных

С точки зрения разработчика подключаемая база данных ничем не отличается от базы данных с единственным владельцем. Приложение подключается к такой базе данных тем же самым способом, как это делалось для базы данных с единственным владельцем в ранних выпусках. Предшествующие примеры создания подключений с использованием

разделяемых и выделенных серверов по-прежнему применимы. Отличия касаются лежащей в основе архитектуры — наличие единственного экземпляра для многих подключаемых баз данных и, как результат, снижение коэффициента использования ресурсов на сервере, а также простота управления для администратора баз данных.

С точки зрения администратора баз данных в методах администрирования базы данных появилось много изменений — позитивных изменений. Например, если администратор сконфигурировал контейнерную базу данных для RAC, то любой подключаемой базе данных, относящейся к этому контейнеру, будут доступны возможности RAC. Аналогичная ситуация и со средством Data Guard, резервными копиями RMAN и т.д. У администратора баз данных есть один экземпляр, предназначенный для конфигурирования и работы с ним, а не по одному экземпляру на приложение, как было в прошлом.

Резюме

На этом краткое знакомство с архитектурой Oracle завершено. В этой главе были определены термины *экземпляр*, *база данных* и *подключаемая база данных*, а также показано, каким образом выполняется подключение к базе данных посредством выделенного или разделяемого сервера. На рис. 2.6 обобщены материалы, раскрытые в главе, и показано взаимодействие с базой данных клиента, применяющего подключение через разделяемый сервер, и клиента, который использует подключение посредством выделенного сервера. На нем также видно, что экземпляр Oracle может применять оба типа подключений одновременно. (В действительности база данных Oracle всегда поддерживает подключения с использованием выделенного сервера — даже если она сконфигурирована для разделяемого сервера.)

Теперь можно приступить к более подробному анализу файлов, образующих базу данных, и процессов, действующих на сервере — выполняемых ими функций и взаимодействия друг с другом. Вы также готовы к ознакомлению с областью SGA, ее содержимым и назначением. Следующая глава начнется с описания типов файлов, которые Oracle применяет для управления данными, и ролью, которую играет каждый тип.

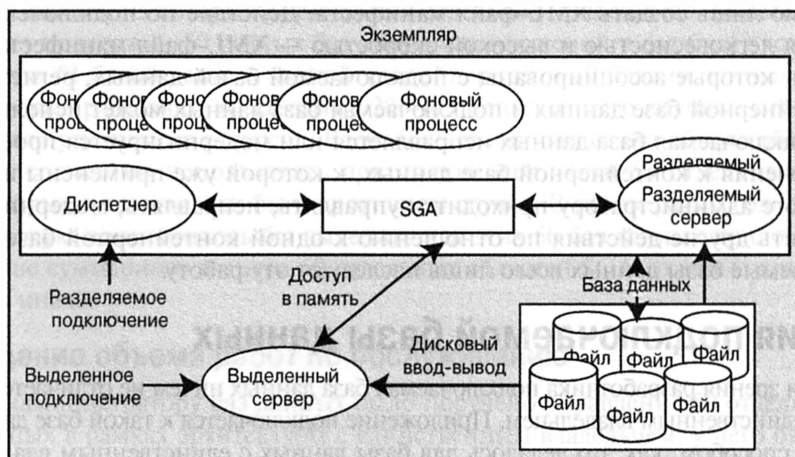


Рис. 2.6. Обзор подключений

ГЛАВА 3

Файлы

В этой главе мы рассмотрим восемь главных типов файлов, которые образуют базу данных и экземпляры. Перечисленные ниже файлы связаны с экземпляром.

- **Файлы параметров (parameter file).** Эти файлы указывают экземпляру базы данных Oracle, где искать управляющие файлы, а также задают набор параметров инициализации, которые определяют размеры ряда структур памяти и т.д. Мы рассмотрим два доступных варианта хранения файлов параметров.
- **Трассировочные файлы (trace file).** Эти диагностические файлы создаются процессом сервера по большей части в ответ на возникновение некоторых необычных ошибок.
- **Сигнальный файл (alert file).** Этот файл подобен трассировочным файлам, но содержит информацию об “ожидаемых” событиях, а также служит единым централизованным файлом для хранения оповещений администратора базы данных о многих событиях, связанных с базой данных.

Следующие файлы образуют базу данных.

- **Файлы данных (data file).** Эти файлы предназначены для базы данных; они хранят таблицы, индексы и все остальные сегменты.
- **Временные файлы (temp file).** Эти файлы используются для выполнения дисковых сортировок и в качестве временного хранилища.
- **Управляющие файлы (control file).** Эти файлы указывают местоположение файлов данных, временных файлов, журнальных файлов повторения действий (файлов с измененными данными), а также содержат другие важные метаданные об их состоянии. Они также хранят информацию резервных копий, поддерживаемую диспетчером восстановления (Recovery Manager — RMAN).
- **Журнальные файлы повторения действий (redo log file).** Эти файлы представляют собой журналы транзакций.
- **Файлы паролей (password file).** Эти файлы применяются для аутентификации пользователей, выполняющих административные действия через сеть. Мы не будем подробно обсуждать файлы такого типа, поскольку они не являются обязательным компонентом любой базы данных Oracle.

В Oracle 10g и последующих версиях существует пара необязательных типов файлов, которые используются Oracle для содействия ускоренным операциям резервного копия и восстановления.

- **Файл отслеживания изменений (change-tracking file).** Этот файл обеспечивает по-настоящему инкрементное резервное копирование данных Oracle. Он не обязательно должен находиться в области для быстрого восстановления (Fast Recovery Area), но из-за того, что он связан исключительно с резервным копированием и восстановлением данных, мы рассмотрим его в контексте этой области.
- **Файлы ретроспективного журнала (flashback log file).** Эти файлы хранят “предшествующие образы” блоков базы данных для содействия команде FLASHBACK DATABASE.

Мы также взглянем на другие типы файлов, обычно связанных с базой данных.

- **Файлы дампа (dump (DMP) file).** Эти файлы генерируются утилитой экспорта базы данных (Export) и используются утилитой импорта базы данных (Import). Следует отметить, что в текущих выпусках Oracle утилита Export устарела — полностью поддерживается только утилита Import. Она поддерживается для того, чтобы обеспечить перемещение данных из старых выпусков Oracle (где утилита Export поддерживалась полностью) в более новые выпуски базы данных.
- **Файлы помпы данных (Data Pump file).** Такие файлы генерируются в Oracle процессом Data Pump Export (EXPDP) и потребляются процессом Data Pump Import (IMPDP). Этот файловый формат может также создаваться и использоваться внешними таблицами.
- **Плоские файлы (flat file).** Это простые старые файлы, которые можно просматривать в текстовом редакторе. Обычно их применяют для загрузки данных в базу.

Из всех перечисленных файлов наиболее важными являются файлы данных и журнальные файлы повторения действий, т.к. они содержат данные, к сбору которых вы приложили большие усилия. Утеря любых или даже всех остальных файлов не приведет к безвозвратной утрате данных. Утеря журнальных файлов может привести к потере некоторых данных. Утеря файлов данных и их резервных копий неизбежно приведет к безвозвратной потере данных.

Далее мы рассмотрим каждый из упомянутых типов файлов: где они обычно располагаются, как именуются и что в них можно ожидать найти.

Файлы параметров

С базой данных Oracle связано много разных файлов параметров, например, файл `tnsname.ora` на рабочей станции клиента (используемый для “нахождения” сервера в сети), файл `listener.ora` на сервере (предназначенный для начального запуска прослушивающего процесса), файлы `sqlnet.ora`, `scan.ora` и `ldap.ora` — и это далеко не полный перечень. Однако наиболее важным из них является файл параметров базы данных. Как было продемонстрировано в главе 2, без него невозможно даже запустить экземпляр. Остальные файлы также важны;

все они имеют отношение к работе в сети и подключению к базе данных. Тем не менее, они выходят за рамки нашего обсуждения. Информацию по их установке и конфигурированию можно найти в руководстве администратора сетевых служб (*Net Services Administrator's Guide*). Как правило, разработчику не приходится настраивать эти файлы — все уже сделано.

Файл параметров базы данных обычно называют `init`-файлом или файлом `init.ora`. Это объясняется его исторически сложившимся стандартным именем, которое выглядит как `init<ORACLE_SID>.ora`. Я называю это имя “исторически сложившимся”, потому что в версии Oracle9i Release 1 был введен чрезвычайно усовершенствованный метод сохранения параметров настройки базы данных — файл параметров сервера, или просто `SPFILE`. По умолчанию этот файл имеет имя `spfile<ORACLE_SID>.ora`. Мы рассмотрим оба вида файлов параметров.

На заметку! Для тех, кто не знаком с термином `SID` или `ORACLE_SID`, здесь приводится его полное определение. Аббревиатура `SID` означает *site identifier* (идентификатор сайта). В среде UNIX/Linux идентификатор `SID` и `ORACLE_HOME` (место, где установлено программное обеспечение Oracle) совместно хешируются с целью образования уникального имени ключа для создания или присоединения к области памяти `SGA` (System Global Area — системная глобальная область). Если значения `ORACLE_SID` или `ORACLE_HOME` установлены некорректно, и вы используете локальное (не сетевое) подключение (сведения о локальных/удаленных подключениях ищите в главе 2), то это приведет к возникновению ошибки “ORACLE NOT AVAILABLE” (база данных Oracle недоступна), т.к. невозможно присоединиться к сегменту разделяемой памяти, который идентифицирован этим уникальным ключом. В среде Windows разделяемая память эксплуатируется не так, как в UNIX/Linux, но параметр `SID` по-прежнему важен. В одном и том же месте, указанном в `ORACLE_HOME`, может существовать более одной базы данных, поэтому необходим способ уникальной идентификации каждой из них наряду с их файлами конфигурации.

Без файла параметров запуск базы данных Oracle невозможен. Это обстоятельство делает файл параметров очень важным. Начиная с Oracle9i Release 2 (версии 9.2 и последующих), инструмент резервного копирования и восстановления `RMAN` (Recovery Manager — диспетчер восстановления) учитывает важность этого файла и позволяет включать в набор резервного копирования файл параметров сервера (но не унаследованный файл параметров `init.ora`). Тем не менее, поскольку `init.ora` является простым текстовым файлом, который можно создавать с помощью любого текстового редактора, его не стоит защищать любой ценой. Он всегда может быть воссоздан при условии, что его содержимое известно (например, эту информацию можно извлечь из сигнального файла базы данных, если к нему есть доступ, и реконструировать полный файл `init.ora`).

Ниже мы по очереди рассмотрим оба типа файла параметров (`init.ora` и `SPFILE`), но сначала давайте выясним, на что похож файл параметров базы данных.

Что собой представляют параметры

Выражаясь простыми терминами, параметр базы данных можно представлять как пару “ключ–значение”. В предыдущей главе вы сталкивались с важным параметром `db_name`. Этот параметр был сохранен как `db_name = ora12c`. Ключом здесь является `db_name`, а значением — `ora12c`. Это и есть наша пара “ключ–значение”.

Для выяснения текущего значения параметра экземпляра можно запросить представление V\$ по имени V\$PARAMETER. В качестве альтернативы в SQL*Plus можно ввести команду SHOW PARAMETER, например:

```
EODA@ORA12CR1> select value
2   from v$parameter
3  where name = 'db_block_size'
4  /
```

VALUE

8192

```
EODA@ORA12CR1> show parameter db_block_s
```

| NAME | TYPE | VALUE |
|---------------|---------|-------|
| db_block_size | integer | 8192 |

Оба вывода отображают в основном ту же самую информацию, хотя представление V\$PARAMETER предоставляет значительно больше сведений (для выбора доступно значительно больше столбцов, чем показано в этом примере). Однако я отдаю предпочтение команде SHOW PARAMETER из-за простоты ее применения и того факта, что она автоматически использует “групповые символы”. Обратите внимание, что было введено только db_block_size; команда SHOW PARAMETER добавляет символ % в начало и конец.

На заметку! Все представления V\$ и все словарные представления полностью документированы в руководстве *Oracle Database Reference* (Справочник по СУБД Oracle). Рассматривайте это руководство как наиболее полный источник сведений о том, что доступно в заданном представлении.

Если бы вы выполнили предыдущий пример от имени менее привилегированного пользователя (для работы с примерами этой книги пользователю EODA была выдана роль DBA), то увидели бы взамен следующий вывод:

```
EODA@ORA12CR1> connect scott/tiger
```

Connected.

Подключено.

```
SCOTT@ORA12CR1> select value
2   from v$parameter
3  where name = 'db_block_size'
4  /
```

from v\$parameter

*

ERROR at line 2:

ORA-00942: table or view does not exist

ОШИБКА в строке 2:

ORA-00942: таблица или представление не существует

```
SCOTT@ORA12CR1> show parameter db_block_s
```

ORA-00942: table or view does not exist

ORA-00942: таблица или представление не существует

По умолчанию “обычным” учетным записям не предоставлен доступ к представлениям производительности V\$. Однако пусть вас это не сбивает с толку. Существует документированный API-интерфейс, как правило, доступный всем пользователям, который позволяет видеть содержимое V\$PARAMETER; следующая небольшая вспомогательная функция поможет увидеть, что установлено в качестве параметра, например:

```
SCOTT@ORA12CR1> create or replace
 2 function get_param( p_name in varchar2 )
 3 return varchar2
 4 as
 5     l_param_type      number;
 6     l_intval          binary_integer;
 7     l_strval          varchar2(256);
 8     invalid_parameter exception;
 9     pragma exception_init( invalid_parameter, -20000 );
10 begin
11     begin
12         l_param_type :=
13             dbms_utility.get_parameter_value
14             ( parnam => p_name,
15               intval => l_intval,
16               strval => l_strval );
17     exception
18         when invalid_parameter
19         then
20             return '*access denied*';
21     end;
22     if ( l_param_type = 0 )
23     then
24         l_strval := to_char(l_intval);
25     end if;
26     return l_strval;
27 end get_param;
28 /
```

Function created.

Функция создана.

На заметку! Если вы применили последнее исправление безопасности для версии 11g R2 или 12c, то вам может понадобиться выдать пользователю, выполняющему эту функцию, привилегию select на v_\$parameter или select_catalog_role.

Если вы запустите эту функцию в SQL*Plus, то увидите такой вывод:

```
SCOTT@ORA12CR1> exec dbms_output.put_line( get_param( 'db_block_size' ) );
8192
```

PL/SQL procedure successfully completed.

Процедура PL/SQL успешно выполнена.

Через API-вызов dbms_utility.get_parameter_value доступен не каждый параметр. В частности, параметры, связанные с памятью, такие как sga_max_size, db_cache_size, pga_aggregate_target и тому подобные, видны не будут.

Мы обрабатываем эту ситуацию в строках кода с 17 по 21, возвращая '*access denied*' (доступ запрещен) при обращении к параметру, видеть который не разрешено. Если вас интересует полный список параметров с ограниченным доступом, можете (как и любая учетная запись, которой выдано право EXECUTE на эту функцию) запустить следующий запрос:

```

EODA@ORA12CR1> select name, scott.get_param( name ) val
2   from v$parameter
3   where scott.get_param( name ) = '*access denied*';

```

| NAME | VAL |
|-------------------------|-----------------|
| ----- | ----- |
| sga_max_size | *access denied* |
| shared_pool_size | *access denied* |
| large_pool_size | *access denied* |
| java_pool_size | *access denied* |
| streams_pool_size | *access denied* |
| ... | |
| client_result_cache_lag | *access denied* |
| olap_page_pool_size | *access denied* |
| 25 rows selected. | |
| 25 строк выбрано. | |

На заметку! В разных версиях будут получаться разные результаты. С течением времени вам следует ожидать изменений количества и значений недоступных параметров в большую или меньшую сторону по мере того, как изменяется общее число параметров.

Если подсчитать количество документированных параметров, которые можно устанавливать в каждой версии базы данных — 9i Release 2, 10g Release 2, 11g Release 1, 11g Release 2 и 12c Release 1, — то получится, соответственно, 258, 259, 294, 342 и 368 разных параметров (я уверен, что могут быть и дополнительные параметры, доступные в конкретной операционной системе). Другими словами, количество параметров (и их имена) варьируются в зависимости от выпуска. Большая часть параметров вроде `db_block_size` существует очень долго (они не исчезают от выпуска к выпуску), но многие другие параметры со временем устаревают из-за изменений, вносимых в реализации.

Например, в Oracle 9.0.1 и предшествующих версиях — вплоть до Oracle 6 — присутствовал параметр `distributed_transactions`, который можно было устанавливать в какое-то положительное значение и управлять количеством параллельных распределенных транзакций, выполняемых базой данных. Он был доступен в более ранних выпусках, но отсутствует во всех выпусках, вышедших после 9.0.1. Попытка его использования в этих выпусках вызовет ошибку, например:

```

EODA@ORA12CR1> alter system set distributed_transactions = 10;
alter system set distributed_transactions = 10
*
ERROR at line 1:
ORA-25138: DISTRIBUTED_TRANSACTIONS initialization parameter has been made
obsolete
ОШИБКА в строке 1:
ORA-25138: Параметр инициализации DISTRIBUTED_TRANSACTIONS устарел

```

Если вы хотите просмотреть эти параметры и получить представление о том, какие из них доступны и что они делают, обращайтесь в руководство *Oracle Database Reference*. В первой главе этого руководства подробно описаны все документированные параметры. В большинстве случаев значения, присваиваемые этим параметрам по умолчанию (либо производные значения, получаемые на основе других параметров), вполне подходят для большинства систем. В то же время, значения параметров, таких как `control_files` (определяет местоположение управляющих файлов системы), `db_block_size` и различных параметров, связанных с памятью, требуют индивидуальной настройки для каждой базы данных.

Обратите внимание, что в предыдущем абзаце я использовал термин “документированные”. Существует также множество *недокументированных* параметров. Их можно определить по именам, начинающимся с символа подчеркивания (`_`). Относительно них ведется множество споров. Поскольку они не документированы, некоторые верят в их “магические возможности”, и многие считают, что они хорошо известны и применяются инсайдерами Oracle. На самом деле я убежден в обратном. Эти параметры не очень хорошо исследованы, и с ними трудно работать. Большинство из недокументированных параметров довольно скучны, т.к. представляют устаревшую функциональность и флаги обратной совместимости. Другие помогают в восстановлении данных, а не самой базы данных; например, некоторые из этих параметров позволяют выполнить запуск базы данных в определенных экстремальных условиях, но только на время, достаточное для вывода данных. После этого база данных требует повторного построения.

Если только вы не получили прямого указания от службы поддержки Oracle, то нет никакой необходимости включать недокументированные параметры в реальную конфигурацию. Многие из них вызывают побочные эффекты, которые могут приводить к разрушительным последствиям. В своей производственной базе данных я не использую недокументированные параметры.

Внимание! Применяйте недокументированные параметры только по запросу из службы поддержки Oracle (Oracle Support). Их использование может привести к повреждению базы данных, а их реализация может — и будет — меняться от версии к версии.

Разнообразные значения параметров можно устанавливать одним из двух способов: либо только для текущего экземпляра, либо постоянно. Обеспечение наличия нужных значений в файле параметров — как раз ваша задача. В случае применения унаследованных файлов параметров `init.ora` это делается вручную. Для изменения значения параметра на постоянной основе, чтобы оно сохранялось между перезапусками сервера, вы должны вручную редактировать файл параметров `init.ora`. Используя файлы параметров сервера, вы увидите, что эта работа более-менее автоматизирована в единственной команде.

Унаследованные файлы параметров `init.ora`

Структура унаследованного файла `init.ora` очень проста. Он представляет собой последовательность пар “ключ–значение”.

Ниже приведен пример содержимого этого файла:

```
control_files='/u01/dbfile/ORA12CR1/control01.ctl','/u02/dbfile/
ORA12CR1/control02.ctl'
db_block_size=8192
db_name='ORA12CR1'
```

По существу этот пример близко отражает наиболее элементарное содержимое файла `init.ora`, с которым можно было бы столкнуться в реальности. Если на конкретной платформе можно применять стандартный размер блока (который варьируется в зависимости от платформы), то строку с параметром `db_block_size` можно удалить. Файл параметров используется, по меньшей мере, для получения имени базы данных и местоположения управляющих файлов. Управляющие файлы сообщают Oracle местоположение всех остальных файлов, поэтому они очень важны для процесса “начальной загрузки”, который запускает экземпляр.

Теперь, когда вы знаете, что собой представляют унаследованные файлы параметров базы данных, и где можно найти дополнительную информацию о допустимых параметрах, которые можно устанавливать, осталось также выяснить, где эти файлы искать на диске. По умолчанию для этих файлов приняты следующие соглашения, касающиеся их имен:

- `init$ORACLE_SID.ora` (переменная среды UNIX/Linux)
- `init%ORACLE_SID%.ora` (переменная среды Windows)

По умолчанию они будут находиться в следующих местах:

- `$ORACLE_HOME/dbs` (UNIX/Linux)
- `%ORACLE_HOME%\DATABASE` (Windows)

Интересно отметить, что во многих случаях все содержимое этого файла параметров будет выглядеть приблизительно так:

```
IFILE='/некоторый/путь/к/каталогу/init.ora'
```

Директива `IFILE` работает подобно директиве `#include` в языке C. Она включает содержимое указанного файла в текущий файл. Здесь директива включает файл `init.ora` из нестандартного местоположения.

Следует отметить, что файл параметров не обязательно должен размещаться в каком-то определенном местоположении. При запуске экземпляра можно использовать опцию `ptile=имя_файла` команды запуска. Это особенно удобно, если требуется оценить влияние различных параметров `init.ora` на базу данных.

Унаследованные файлы параметров можно обслуживать с помощью любого текстового редактора. Например, в среде UNIX/Linux я использую `vi`; во многих версиях Windows я применяю Notepad, а в среде мейнфрейма, вероятно, воспользовался бы Xedit. Важно понимать, что именно вы полностью отвечаете за редактирование и обслуживание этого файла. Внутри самой базы данных нет никаких команд, которые можно было бы применить для изменения значений в файле `init.ora`. Например, когда используется файл параметров `init.ora`, то выполнение команды `ALTER SYSTEM` для изменения размера компонента SGA не приведет к постоянному изменению в этом файле. Если нужно, чтобы изменение стало постоянным (другими словами, чтобы оно стало стандартным при последующих перезапусках базы

данных), вы должны обеспечить ручное обновление содержимого всех файлов параметров `init.ora`, которые могут применяться для запуска этой базы данных.

И последняя интересная особенность унаследованного файла параметров связана с тем, что он не обязательно должен располагаться на сервере базы данных. Одной из причин ввода файла параметров (которые мы вскоре обсудим) было исправление этой ситуации. Унаследованный файл параметров должен присутствовать на компьютере клиента, пытающегося запустить базу данных. Это значит, что в случае использования сервера UNIX/Linux, но его администрировании через сеть с помощью интерфейса SQL*Plus, установленного на настольном компьютере Windows, файл параметров базы данных должен находиться на настольном компьютере.

До сих пор помню ситуацию, когда я сделал болезненное открытие, что файлы параметров не хранятся на сервере. Это произошло много лет назад, как только появился совершенно новый (теперь изъятый из обращения) инструмент под названием *SQL*DBA*. Он позволял выполнять удаленные операции, в частности, удаленное администрирование. Со своего сервера (функционирующего тогда под управлением SunOS) я мог удаленно подключаться к серверу базы данных мейнфрейма. Также была возможность выдачи команды завершения. Тем не менее, именно тогда я и понял, что не все так просто — при попытке запуска экземпляра инструмент SQL*DBA сообщал о невозможности нахождения файла параметров. Я знал, что эти файлы параметров — простые текстовые файлы `init.ora` — размещались на компьютере клиента; они должны были существовать на стороне *клиента*, а не сервера. Инструмент SQL*DBA искал файл параметров в моей локальной системе, чтобы запустить базу данных мейнфрейма. Мало того, что на моем компьютере этот файл отсутствовал, я еще и не имел ни малейшего понятия о том, что в него нужно было поместить, чтобы запустить систему! Я не знал ни значения параметра `db_name`, ни местоположений управляющих файлов (даже получение информации о соглашениях по именованию файлов на мейнфрейме оказалось довольно сложной задачей), и не имел доступа для входа в саму систему мейнфрейма. С тех пор я больше не повторял эту ошибку; урок был достаточно болезненным, чтобы хорошо его усвоить.

Когда администраторы баз данных осознали, что файл параметров `init.ora` должен находиться на компьютере клиента, запускающего базу данных, это привело к быстрому размножению таких файлов. Каждый администратор базы данных желал запускать инструменты администрирования со своего рабочего стола, поэтому ему требовалась копия файла параметров на его настольном компьютере. Такие средства, как OEM (Oracle Enterprise Manager — диспетчер предприятия Oracle), добавляли в общую смесь еще один файл параметров. Инструменты подобного рода пытались централизовать задачи администрирования всех баз данных предприятия на одном компьютере, который иногда называли *сервером управления*. На этом единственном компьютере должно было функционировать программное обеспечение, которое все администраторы базы данных применяли бы для запуска, завершения, резервного копирования и выполнения других действий по администрированию баз данных. На первый взгляд такое решение — сосредоточение всех файлов параметров в одном месте и использование инструментов с графическим пользовательским интерфейсом для выполнения всех операций — выглядит идеальным. Однако реальность состоит в том, что иногда при проведении некоторых административных задач значительно удобнее выдавать команду администрирования из среды SQL*Plus

самого сервера баз данных. В результате снова возникла ситуация существования нескольких файлов параметров: одного на сервере управления и одного на сервере баз данных. Кроме того, эти файлы параметров требовали бы синхронизации друг с другом, и пользователям пришлось бы удивляться, почему параметр, который они изменили в прошлом месяце, вдруг “исчез”, а затем по непонятной причине снова “появился”.

В итоге это привело к введению файла параметров сервера (SPFILE), который теперь может служить единственным “критерием истины” для базы данных.

Файлы параметров сервера (SPFILE)

Файлы SPFILE представляют фундаментальное изменение в способе доступа Oracle к параметрам настройки для экземпляра и их сопровождения. Файл SPFILE устраняет две серьезные проблемы, связанные с унаследованными файлами параметров.

- Он предотвращает размножение файлов параметров. Файл SPFILE всегда хранится на сервере базы данных; он должен находиться на компьютере самого сервера и не может быть размещен на компьютере клиента. Это обстоятельство обуславливает целесообразность наличия единственного источника “истинных” значений параметров.
- Он устраняет необходимость (фактически лишает возможности) поддерживать файлы параметров вручную с помощью текстовых редакторов вне базы данных. Команда ALTER SYSTEM позволяет записывать значения непосредственно в файл SPFILE. Администраторам больше не нужно находить и обслуживать все файлы параметров вручную.

По умолчанию для этого файла действует следующее соглашение по именованию:

- \$ORACLE_HOME/dbs/spfile\$ORACLE_SID.ora (переменная среды Unix/Linux)
- %ORACLE_HOME/database/spfile%ORACLE_SID%.ora (переменная среды Windows)

Я настоятельно рекомендую применять стандартное местоположение; в противном случае простота использования файла SPFILE будет сведена на нет. Когда этот файл находится в стандартном каталоге, все действия выполняются более или менее автоматически. Перемещение файла SPFILE в другой каталог требует указания Oracle этого местоположения, снова приводя к возникновению проблем, которые присущи унаследованным файлам параметров!

Преобразование в файлы SPFILE

Предположим, что у вас есть база данных, в которой применяется унаследованный файл параметров. Переход на файл SPFILE довольно прост — вы используете команду CREATE SPFILE.

На заметку! Можно также применять “обратную” команду для создания файла параметров (PFILE) из файла SPFILE. Вскоре я объясню, чем может быть вызвана такая необходимость.

Итак, при наличии файла параметров `init.ora` в стандартном местоположении на сервере достаточно выдать команду `CREATE SPFILE` и перезапустить экземпляр сервера:

```
EODA@ORA12CR1> show parameter spfile;
```

| NAME | TYPE | VALUE |
|--------|--------|-------|
| spfile | string | |

```
EODA@ORA12CR1> create spfile from pfile;
```

```
create spfile from pfile
```

```
*
```

```
ERROR at line 1:
```

```
ORA-01031: insufficient privileges
```

```
ОШИБКА в строке 1:
```

```
ORA-01031: недостаточно привилегий
```

Что же, команда `SHOW PARAMETER` показывает, что мы не создали `SPFILE`: значение в колонке `VALUE` пустое. Мы не располагаем достаточными привилегиями, чтобы создать `SPFILE`, хотя вход был совершен от имени администратора базы данных. Создание `SPFILE` считается очень высокопривилегированной операцией, и ее можно выполнять, только в случае подключения с использованием учетных данных, которые позволяют запускать и останавливать базу данных. Давайте сделаем это:

```
EODA@ORA12CR1> connect / as sysoper;
```

```
Connected.
```

```
Подключено.
```

```
PUBLIC@ORA12CR1> create spfile from pfile;
```

```
File created.
```

```
Файл создан.
```

```
PUBLIC@ORA12CR1> startup force;
```

```
ORACLE instance started.
```

```
Database mounted.
```

```
Database opened.
```

```
Экземпляр ORACLE запущен.
```

```
База данных смонтирована.
```

```
База данных открыта.
```

Я воспользовался наименее привилегированной учетной записью, от имени которой можно выполнить эту операцию — записью (принадлежащей мне) с административной привилегией `SYSOPER`. Привилегия `SYSOPER` позволяет управлять файлами параметров, запускать и останавливать базу данных, но ничего сверх этого; вот почему вывод команды `startup` выглядит иначе — отсутствует поддержка `SG`, и фактически не видны параметры памяти:

```
PUBLIC@ORA12CR1> show parameter spfile;
```

```
ORA-00942: table or view does not exist
```

```
ORA-00942: таблица или представление не существует
```

Хотя привилегия `SYSOPER` дает право запускать и останавливать базу данных, она не имеет доступа к представлениям `V$` и т.д. Ее возможности весьма ограничены. Мы можем проверить, что файл `SPFILE` применяется, подключившись от имени учетной записи с достаточными привилегиями:


```
EODA@ORA12CR1> show parameter spfile;
```

| NAME | TYPE | VALUE |
|--------|--------|--|
| spfile | string | /home/ora12cr1/app/ora12cr1/pr oduct/12.1.0/dbhome_1/dbs/spfi leora12cr1.ora |

Мы использовали здесь команду SHOW PARAMETER, чтобы показать, что изначально файл SPFILE не применялся, но после его создания и перезапуска экземпляра он используется и имеет стандартное имя.

На заметку! В кластеризованной среде, в которой применяется Oracle RAC, все экземпляры совместно используют тот же самый файл SPFILE, поэтому процесс преобразования файла PFILE в файл SPFILE должен производиться в управляемой манере. Единственный файл SPFILE может содержать все значения параметров, даже специфичные для экземпляра, но все необходимые файлы параметров придется объединить в единый файл PFILE с описанным ниже форматом.

В кластеризованной среде для перехода от отдельных файлов SPFILE к общему файлу SPFILE индивидуальные файлы SPFILE понадобится объединить в единый файл, подобный показанному ниже:

```
*.cluster_database_instances=2
*.cluster_database=TRUE
*.cluster_interconnects='10.10.10.0'
*.compatible='12.1.0.0.0'
*.control_files='/u1/d1/O12C/control01.ctl','/u1/d2/O12C/control02.ctl'
*.db_name='O12C'
...
*.processes=150
*.undo_management='AUTO'
O12C1.instance_number=1
O12C2.instance_number=2
O12C1.local_listener='LISTENER_O12C1'
O12C2.local_listener='LISTENER_O12C2'
O12C1.remote_listener='LISTENER_O12C2'
O12C2.remote_listener='LISTENER_O12C1'
O12C1.thread=1
O12C2.thread=2
O12C1.undo_tablespace='UNDOTBS1'
O12C2.undo_tablespace='UNDOTBS2'
```

То есть имена параметров, общих для всех экземпляров в кластере, будут начинаться с символов *. Имена параметров, характерных для отдельного экземпляра, вроде INSTANCE_NUMBER и THREAD, должны быть предварены именем экземпляра (идентификатором Oracle SID). В предыдущем примере:

- файл PFILE будет предназначен для двухузлового кластера, который содержит экземпляры с именами O12C1 и O12C2;
- присваивание *.db_name = 'O12C' указывает, что все экземпляры, использующие этот файл SPFILE, будут монтировать базу данных по имени O12C;
- присваивание O12C1.undo_tablespace='UNDOTBS1' отражает то, что экземпляр O12C1 будет применять эту конкретную табличное пространство, и т.д.

Установка значений в файлах SPFILE

После того как база данных запущена с участием SPFILE, возникает вопрос о том, как устанавливать и изменять значения, содержащиеся в этом файле. Вспомните, что файлы SPFILE являются двоичными, и их нельзя редактировать в простом текстовом редакторе. Ответом будет команда ALTER SYSTEM с показанным ниже синтаксисом (элементы, заключенные в скобки <>, являются необязательными, а символ канала | указывает на необходимость выбора одного элемента из списка):

```
Alter system set параметр=значение <comment='текст'> <deferred>
                <scope=memory|spfile|both> <sid='системный_идентификатор|* '>
                <container=current|all>
```

По умолчанию команда ALTER SYSTEM SET будет обновлять выполняющийся в текущий момент экземпляр, а также вносить изменения в файл SPFILE или в случае подключаемой базы данных — в словарь данных этой подключаемой базы данных (дополнительные сведения будут приведены в разделе “Подключаемые базы данных” далее в главе). Это значительно облегчает администрирование и устраняет проблемы, возникающие во время установки или модификации параметров посредством команды ALTER SYSTEM, если при этом администратор забыл обновить файл параметров init.ora или этот файл вообще отсутствует.

Давайте взглянем на каждый элемент этой команды.

- Присваивание параметр=значение предоставляет имя параметра и его новое значение. Например, элемент pga_aggregate_target=1024m установил бы значение параметра PGA_AGGREGATE_TARGET в 1024 Мбайт (1 Гбайт).
- Присваивание comment='текст' — необязательный комментарий, который можно связать с данным значением параметра. Этот комментарий будет отображаться в поле UPDATE_COMMENT представления V\$PARAMETER. При использовании варианта сохранения изменения в SPFILE комментарий будет записан в этот файл, и также сохраняться при последующих перезапусках сервера, так что впоследствии он будет видимым для базы данных.
- Параметр deferred указывает на то, что системные изменения должны вступить в силу только в последующих сеансах (т.е. не в текущих сеансах, включая тот, в котором производится изменение). По умолчанию команда ALTER SYSTEM будет вступать в силу сразу же, но некоторые параметры не могут быть изменены немедленно — это допускается только для вновь созданных сеансов. Выяснить, какие параметры требуют применения deferred, можно с помощью следующего запроса:

```
EODA@ORA12CR1> select name
  2   from v$parameter
  3   where issys_modifiable='DEFERRED'
  4   /
```

NAME

```
-----
backup_tape_io_slaves
recyclebin
audit_file_dest
object_cache_optimal_size
```

```
object_cache_max_size_percent
sort_area_size
sort_area_retained_size
olap_page_pool_size

8 rows selected.
8 строк выбрано.
```

На заметку! Результаты, полученные вами, могут отличаться; список параметров, которые могут быть установлены в оперативном режиме, но их установка должна быть отложена, от версии к версии может — и будет — отличаться.

В выводе видно, что параметр `sort_area_size` можно изменять на системном уровне, но только отложенным образом. Показанный ниже код демонстрирует, к чему приводят попытки изменения значения `sort_area_size` с и без `deferred`:

```
EODA@ORA12CR1> alter system set sort_area_size = 65536;
alter system set sort_area_size = 65536
*
ERROR at line 1:
ORA-02096: specified initialization parameter is not modifiable with this
option
ОШИБКА в строке 1:
ORA-02096: указанный параметр инициализации не может быть изменен
подобным образом

EODA@ORA12CR1> alter system set sort_area_size = 65536 deferred;
System altered.
Система изменена.
```

- `scope=memory|spfile|both` указывает “область действия” этого параметра. Рассмотрим варианты установки значения.
 - `scope=memory` изменяет значение только в экземпляре (экземплярах); после перезапуска базы данных это изменение будет утеряно. При следующем запуске значение будет таким, каким оно зафиксировано в файле `SPFILE`.
 - `scope=spfile` изменяет значение только в файле `SPFILE`. Изменение не вступит в силу, пока база данных не будет перезапущена и файл `SPFILE` не обработается повторно. Некоторые параметры могут быть изменены только с применением этого варианта — например, параметр `processes` должен использовать `scope=spfile`, т.к. изменить значение активного экземпляра невозможно.
 - `scope=both` означает, что изменение параметра происходит как в памяти, так и в файле `SPFILE`. Изменение отразится в текущем экземпляре, а также будет действовать при следующем запуске базы данных. При использовании файла `SPFILE` эта область действия считается стандартной. В случае применения файла параметров `init.ora` стандартным и единственным допустимым значением является `scope=memory`.
- Присваивание `sid='системный_идентификатор|*'` главным образом удобно в кластеризованной среде. Вариант `sid='*'` используется по умолчанию. Это позволяет указывать значение параметра для любого конкретного экземпляра в кластере. Указывать `sid=` нужно только в случае применения технологии Oracle RAC.

- Присваивание `container=current|all` используется в базе данных с множеством владельцев для определения области действия изменения. Если команда `ALTER SYSTEM` выполняется в корневой контейнерной базе данных, то за счет применения параметра `all` значение может распространяться в каждую подключаемую базу данных. В противном случае по умолчанию изменение воздействует только на текущую контейнерную или подключаемую базу данных. Обратите внимание, что настройки, специфичные для подключаемой базы данных, не записываются в файл `SPFILE`, а сохраняются в словаре данных этой подключаемой базы данных, так что в случае перемещения в другой контейнер ее настройки переместятся вместе с ней.

Команда обычно используется следующим образом:

```
EODA@ORA12CR1> alter system set pga_aggregate_target=512m;
System altered.
Система изменена.
```

На заметку! Предыдущая команда — и фактически многие команды `ALTER SYSTEM` в этой книге — может дать сбой в вашей системе. Если вы применяете другие параметры, которые несовместимы с приведенным примером (скажем, другие параметры памяти), может возникнуть ошибка. Это вовсе не означает, что команда не работает, а скорее то, что параметры, которые вы пытаетесь использовать, несовместимы с общей установкой.

Вероятно, даже еще лучше применять присваивание `COMMENT=`, чтобы документировать, когда и почему было внесено конкретное изменение:

```
EODA@ORA12CR1> alter system set pga_aggregate_target=512m
  2 comment = 'Changed 14-aug-2013, AWR recommendation';
System altered.
Система изменена.
```

```
EODA@ORA12CR1> select value, update_comment
  2 from v$parameter
  3 where name = 'pga_aggregate_target'
  4 /
```

| VALUE | UPDATE_COMMENT |
|-----------|---|
| 536870912 | Changed 14-aug-2013, AWR recommendation |

Отмена установки значений в файлах `SPFILE`

Теперь возникает следующий вопрос: как отменить установленное значение? Другими словами, нам больше не нужен такой параметр в файле `SPFILE`. Каким образом решить эту задачу, если файл нельзя редактировать в текстовом редакторе? Для этого также используется команда `ALTER SYSTEM`, но на этот раз с конструкцией `RESET`:

```
Alter system reset parameter <scope=memory|spfile|both> sid='sid|*'
```

Итак, например, если необходимо удалить параметр `sort_area_size`, чтобы дать ему возможность принимать стандартное значение, введите следующую команду:

```
EODA@ORA12CR1> alter system reset sort_area_size scope=spfile;
System altered.
Система изменена.
```

На заметку! В предшествующих выпусках, в частности, в Oracle 10g Release 2 и более ранних, конструкция `SID=` не была необязательной, как сейчас. Тогда нужно было помещать конструкцию `SID='*'` в конец команды `ALTER SYSTEM`, чтобы сбросить параметр для всех экземпляров в `SPFILE`. Или же требовалось указать `SID='некоторый_sid'`, чтобы сбросить его для одного экземпляра.

В результате параметр `sort_area_size` удаляется из `SPFILE`, в чем можно удостовериться так, как показано ниже:

```
EODA@ORA12CR1> connect / as sysoper;
Connected.
Подключено.

PUBLIC@ORA12CR1> create pfile='/tmp/pfile.tst' from spfile;
File created.
Файл создан.
```

После этого можете проверить содержимое файла `/tmp/pfile.tst`, который будет сгенерирован на сервере базы данных. Вы увидите, что `sort_area_size` в файле параметров больше не существует.

Создание файлов *PFILE* из файлов *SPFILE*

Показанная только что команда `CREATE PFILE...FROM SPFILE` является противоположностью команды `CREATE SPFILE`. Она принимает двоичный файл `SPFILE` и создает из него простой текстовый файл — файл, который можно редактировать в любом текстовом редакторе, а затем применять для запуска базы данных. Эту команду можно использовать для выполнения, по меньшей мере, двух описанных ниже действий на регулярной основе.

- **Создание одноразового файла параметров с рядом специальных настроек для запуска базы данных в целях обслуживания.** Сначала вы должны запустить команду `CREATE PFILE...FROM SPFILE` и отредактировать результирующий текстовый файл `PFILE`, изменив нужные параметры. Затем потребуется запустить базу данных с применением `PFILE=<ИМЯ_ФАЙЛА>`, чтобы указать файл `PFILE` вместо `SPFILE`. По завершении достаточно будет запустить базу данных обычным образом, и она будет использовать файл `SPFILE`.
- **Поддержание хронологии комментированных изменений.** В прошлом многие администраторы баз данных интенсивно снабжали свои файлы параметров комментариями, отражающими хронологию изменений. Например, если они изменяли размер кеша буфера 20 раз, то в файле `init.ora` параметру `db_cache_size` могло предшествовать 20 комментариев с указанием даты и причины внесения изменений. Файл `SPFILE` не поддерживает такую возможность, но аналогичного эффекта можно добиться, если выработать у себя привычку поступать следующим образом:

```
PUBLIC@ORA12CR1> connect / as sysdba
Connected.
Подключено.

SYS@ORA12CR1> create pfile='init_14_aug_2013_oral2crl.ora' from spfile;
File created.
Файл создан.
```

```

SYS@ORA12CR1> alter system set pga_aggregate_target=512m
  2 comment = 'Changed 14-aug-2013, AWR recommendation';
System altered.
Система изменена.

```

Таким способом хронология изменений будет сохраняться с течением времени в последовательности файлов параметров.

Исправление поврежденных файлов SPFILE

Последний вопрос, возникающий в отношении файлов SPFILE, формулируется так: “Поскольку файлы SPFILE являются двоичными, что произойдет, если один из них окажется поврежденным и запуск базы данных станет невозможным? Во всяком случае, файл `init.ora` был простым текстовым, и его всегда можно было отредактировать и исправить”. В действительности файлы SPFILE повреждаются не чаще файлов данных, журнальных файлов, управляющих файлов и т.д. Тем менее, если такая проблема все же возникла или в SPFILE установлено значение, которое не позволяет базе данных запуститься, то можно пойти двумя путями.

Прежде всего, объем двоичных данных в файле SPFILE очень мал. В среде UNIX/Linux простая команда `strings` позволяет извлечь все настройки:

```

[oral2cr1@dellpe dbs]$ strings $ORACLE_HOME/dbs/spfile$ORACLE_SID.ora
*.audit_file_dest='/home/oral2cr1/app/oral2cr1/admin
/oral2cr1/adump'
*.audit_trail='db'
*.compatible='12.1.0.0.0'
...

```

В среде Windows достаточно открыть файл с помощью программы `write.exe` (WordPad). Этот редактор отобразит весь обычный текст из файла, после чего посредством операции вырезания и вставки в файл `init<ORACLE_SID>.ora` можно создать файл PFILE, предназначенный для запуска экземпляра.

В случае “пропажи” файла SPFILE (по любой причине — хотя мне с подобными ситуациями сталкиваться не приходилось) информацию о файле параметров можно извлечь из сигнального журнала базы данных (более подробно сигнальный журнал рассматривается далее в главе). При каждом запуске базы данных в сигнальный журнал записывается раздел следующего вида:

```

Starting up:
Oracle Database 12c Enterprise Edition Release 12.1.0.1.0 - 64bit Production
With the Partitioning, OLAP, Advanced Analytics and Real Application
Testing options.
ORACLE_HOME = /home/oral2cr1/app/oral2cr1/product/12.1.0/dbhome_1
System name:   Linux
Node name:    dellpe
Release:      2.6.39-400.109.1.el6uek.x86_64
Version:      #1 SMP Tue Jun 4 23:21:51 PDT 2013
Machine:      x86_64
Using parameter settings in server-side spfile /home/oral2cr1/app/
oral2cr1/product/12.1.0/dbhome_1/dbs/spfileoral2cr1.ora
System parameters with non-default values:
  processes                = 300
  resource_limit            = TRUE
  sga_target                = 4800M

```

```

control_files          = "/home/oral2cr1/oradata/ORAl2CR1/controlfile/
❏ ol_mf_8wvv2pml_.ctl"
control_files          = "/home/oral2cr1/app/oral2cr1/fast_recovery_area/
❏ ORAl2CR1/controlfile/ol_mf_8wvv2ps2_.ctl"
db_block_size          = 8192
compatible             = "12.1.0.0.0"
db_create_file_dest    = "/home/oral2cr1/oradata"
db_recovery_file_dest  = "/home/oral2cr1/app/oral2cr1/fast_recovery_area"
db_recovery_file_dest_size= 4815M
undo_tablespace        = "UNDOTBS1"
remote_login_passwordfile= "EXCLUSIVE"
db_domain              = ""
dispatchers            = "(PROTOCOL=TCP) (SERVICE=oral2cr1XDB)"
local_listener         = "(ADDRESS=(PROTOCOL=tcp) (HOST=192.168.1.25)
❏ (PORT=1521))"
parallel_min_servers   = 0
parallel_max_servers   = 0
audit_file_dest        = "/home/oral2cr1/app/oral2cr1/admin/oral2cr1/adump"
audit_trail            = "DB"
db_name                = "oral2cr1"
open_cursors           = 300
_column_tracking_level = 1
pga_aggregate_target   = 1600M
diagnostic_dest        = "/home/oral2cr1/app/oral2cr1"
NOTE: remote asm mode is local (mode 0x1; from cluster type)
Starting background process PMON
Mon Sep 02 16:56:22 2013
PMON started with pid=2, OS id=21572

```

На основе этого раздела легко создать файл PFILE и затем посредством команды CREATE SPFILE преобразовать его в новый файл SPFILE.

Подключаемые базы данных

Подключаемые базы данных спроектированы как набор файлов, которые можно перемещать из одной корневой контейнерной базы данных в другую. То есть мы можем отсоединить подключаемую базу данных и впоследствии снова присоединить ее либо к той же самой, либо к какой-то другой корневой контейнерной базе данных, и получить исходную подключаемую базу данных обратно — со всеми схемами приложений, пользователями, метаданными, привилегиями, данными и даже настройками параметров (которые не были унаследованы от корневого контейнера). Это достигается за счет сохранения настроек параметров, специфичных для подключаемой базы данных, в таблице словаря данных SYS.PDB_SPFILE\$. Именно таким способом подключаемые базы данных могут переопределять значения некоторых параметров (не каждый параметр допускает установку на уровне подключаемой базы данных) в файле SPFILE и обеспечивать их перенос по мере перемещения между корневыми контейнерными базами данных.

Заключительные соображения по поводу файла параметров

В этом разделе были рассмотрены все основные особенности управления параметрами инициализации Oracle и файлами параметров. Вы узнали, как устанавливать параметры, просматривать их значения и обеспечивать сохранение настроек между

перезапусками базы данных. Мы исследовали два типа файлов параметров базы данных: унаследованные файлы PFILE (простые текстовые файлы) и более новые файлы SPFILE. Для всех существующих баз данных рекомендуется использовать файлы SPFILE из-за легкости администрирования и ясности, которую они привносят. Возможность применения единственного источника “истины” для базы данных в сочетании с возможностью выдачи команды ALTER SYSTEM для сохранения значений параметров делает файлы SPFILE привлекательным средством. Я начал их использовать сразу же, как только они стали доступны, и ни разу не пожалел об этом.

Трассировочные файлы

Трассировочные файлы являются источником отладочной информации. При возникновении проблем сервер генерирует трассировочный файл с огромным объемом диагностической информации. Когда разработчик запускает процедуру DBMS_MONITOR.SESSION_TRACE_ENABLE, сервер создает трассировочный файл с информацией, касающейся производительности. Трассировочные файлы доступны благодаря тому, что Oracle является хорошо инструментированным программным обеспечением. Под “инструментированным” я имею в виду, что программисты, написавшие ядро базы данных, поместили в него очень много отладочного кода и оставили его там на всякий случай.

Я встречал многих разработчиков, которые считали отладочный код накладными расходами — т.е. чем-то таким, от чего следует избавляться перед тем, как приложение поступит в производственную среду, пытаясь выжать максимум производительности из кода. Позже, конечно же, они обнаруживают, что код содержит ошибку и не функционирует настолько быстро, как должен (что конечные пользователи склонны расценивать как ошибку; для конечных пользователей низкая производительность является ошибкой). В этот момент они страстно желают, чтобы отладочный код оставался в коде приложения. Вы должны тестировать любой новый код перед его помещением в производственную среду, и это не должно делаться впопыхах.

СУБД Oracle, а также сервер приложений, приложения Oracle и разнообразные средства, такие как Application Express (APEX), хорошо инструментированы. Ниже перечислены признаки этого инструментария.

- **Представления V\$.** Большинство представлений V\$ содержат “отладочную” информацию. Представления V\$WAITSTAT, V\$SESSION_EVENT и многие другие предназначены исключительно для отображения информации о происходящем внутри ядра базы данных.
- **Команда AUDIT.** Эта команда позволяет указывать события, которые база данных должна записывать для последующего анализа.
- **Диспетчер ресурсов (DBMS_RESOURCE_MANAGER).** Это средство позволяет осуществлять микро-управление ресурсами (центральным процессором, устройствами ввода-вывода и т.п.) внутри базы данных. Возможность применения диспетчера ресурсов в базе данных обеспечивается тем, что база данных имеет доступ ко всем статистическим сведениям об использовании ресурсов во время выполнения.
- **События Oracle.** События позволяют запрашивать у Oracle генерацию необходимой трассировочной или диагностической информации.

- **DBMS_TRACE.** Эта внутренняя функция механизма PL/SQL создает подробное дерево вызовов хранимых процедур, возникших исключений и ошибок.
- **Триггеры событий базы данных.** Эти триггеры, например, ON SERVERERROR, позволяют отслеживать и регистрировать любые условия, которые можно считать “исключительными” или отклоняющимися от нормы. Например, в журнал можно записать SQL-запрос, который выполнялся при возникновении ошибки типа “переполнение временной области”.
- **SQL_TRACE/DBMS_MONITOR.** Служит для просмотра точного SQL-кода, событий ожидания и другой диагностической информации, связанной с производительностью/поведением, которая генерируется в результате выполнения вашего приложения. Средство SQL Trace также доступно в расширенном виде через событие Oracle 10046.

Инструментирование жизненно важно при проектировании и разработке приложений, поэтому с каждым новым выпуском СУБД Oracle становится все лучше и лучше инструментированной. На самом деле количество дополнительных инструментальных средств, появившихся между версиями Oracle9i Release 2 и Oracle 11g, а теперь Oracle 12c, феноменально велико. В версии Oracle 10g инструментирование кода ядра было поднято на совершенно новый уровень благодаря появлению средств AWR (Automatic Workload Repository — репозиторий автоматической рабочей нагрузки) и ASH (Active Session History — хронология активных сеансов). В Oracle 11g были введены новые средства — ADR (Automatic Diagnostic Repository — репозиторий автоматической диагностики) и SPA (SQL Performance Analyzer — анализатор производительности SQL). В версии Oracle 12c разработчики продвинулись еще дальше, добавив журнал DDL для отслеживания всех операций DDL в базе данных (того, что не должно происходить во многих типовых производственных базах данных изо дня в день) и журнал отладки для отслеживания исключительных условий в базе данных.

В этом разделе мы сосредоточим внимание на информации, которую можно найти в трассировочных файлах различных типов. Мы рассмотрим, что они собой представляют, где хранятся и что с ними можно делать.

Существуют два основных типа трассировочных файлов, и предоставляемые ими возможности значительно отличаются.

- **Ожидаемые и желательные трассировочные файлы.** Примером могут служить файлы, генерируемые в результате включения DBMS_MONITOR.SESSION_TRACE_ENABLE. Они содержат диагностическую информацию о сеансе и помогают в настройке приложения для оптимизации производительности и диагностирования любых имеющихся узких мест.
- **Трассировочные файлы, которые не ожидалось, но были сгенерированы сервером в результате ошибок ORA-00600 “Internal Error” (Внутренняя ошибка), ORA-03113 “End of file on communication channel” (Конец файла в канале передачи данных) или ORA-07445 “Exception Encountered” (Возникло исключение).** Эти трассировочные файлы содержат диагностическую информацию, которая предназначена главным образом для использования аналитиками службы поддержки Oracle и помогает выявить место возникновения внутренних ошибок в приложении. Применение этих файлов разработчиками приложений ограничено.

Запрошенные трассировочные файлы

Трассировочные файлы, которые ожидаются разработчиками, чаще всего генерируются в результате включения трассировки посредством DBMS_MONITOR (ALTER SESSION SET SQL_TRACE=TRUE в Oracle9i Release 2 и предшествующих версиях) либо использования расширенной трассировки через событие 10046, как показано ниже:

```
EODA@ORA12CR1> alter session set events
  2 '10046 trace name context forever, level 12'
  3 /
Session altered.
Сеанс изменен.
```

Эти трассировочные файлы содержат диагностическую информацию и информацию, относящуюся к производительности. Они открывают бесценные подробности внутренней работы вашего приложения базы данных. В нормально функционирующей базе данных вы будете видеть такие трассировочные файлы гораздо чаще, чем трассировочные файлы любого другого вида.

Местоположения файлов

Независимо от того, используется DBMS_MONITOR, SQL_TRACE или расширенная функция трассировки, Oracle начнет генерировать трассировочный файл в одном из следующих мест на компьютере сервера базы данных.

- Если вы применяете подключение посредством выделенного сервера, то трассировочный файл будет сгенерирован в каталоге, указанном в параметре user_dump_dest.
- Если вы используете подключение посредством разделяемого сервера, то трассировочный файл будет создан в каталоге, указанном в параметре background_dump_dest.

Чтобы посмотреть, куда попадут трассировочные файлы, можно выдать команду show parameter dump_dest в SQL*Plus, запросить представление V\$PARAMETER, воспользоваться созданной ранее процедурой (SCOTT.GET_PARAM) или запросить новое представление V\$DIAG_INFO. Ниже все перечисленные приемы демонстрируются по очереди.

```
EODA@ORA12CR1> show parameter dump_dest
```

| NAME | TYPE | VALUE |
|----------------------|--------|--|
| background_dump_dest | string | /home/oral2cr1/app/oral2cr1/diag/ rdbms/oral2cr1/oral2cr1/trace |
| core_dump_dest | string | /home/oral2cr1/app/oral2cr1/diag/ rdbms/oral2cr1/oral2cr1/cdump |
| user_dump_dest | string | /home/oral2cr1/app/oral2cr1/diag/ rdbms/oral2cr1/oral2cr1/trace |

Эта команда отображает три местоположения файлов дампа (трассировки). Каталог фонового дампа используется любым “серверным” процессом (полный список фоновых процессов Oracle и выполняемых ими функций вы найдете в главе 5). Параметр core_dump_dest (местоположение дампа ядра) применяется для

выгрузки “дампа ядра” (очень детальная диагностическая информация о процессе) при возникновении серьезной проблемы, такой как аварийный отказ процесса. Параметр `user_dump_dest` (местоположение пользовательского дампа) используется подключениями посредством выделенного и разделяемого сервера (описанными в главе 2), когда они генерируют трассировочный файл.

Чтобы продолжить исследование разнообразных методов просмотра этих местоположений для дампов, давайте взглянем на доступные таблицы `V$`:

```
EODA@ORA12CR1> select name, value
2   from v$parameter
3  where name like '%dump_dest%';
```

| NAME | VALUE |
|----------------------|--|
| background_dump_dest | /home/ora12cr1/app/ora12cr1/diag/rdbms/ora12cr1/ora12cr1/trace |
| user_dump_dest | /home/ora12cr1/app/ora12cr1/diag/rdbms/ora12cr1/ora12cr1/trace |
| core_dump_dest | /home/ora12cr1/app/ora12cr1/diag/rdbms/ora12cr1/ora12cr1/cdump |

Конечно, для запрашивания таблицы `V$PARAMETER` можно было бы также применить пакет `DBMS_UTILITY` с ранее созданной функцией `SCOTT.GET_PARAM`:

```
EODA@ORA12CR1> set serveroutput on
EODA@ORA12CR1> exec dbms_output.put_line( scott.get_param( 'user_dump_dest' ) )
/home/ora12cr1/app/ora12cr1/diag/rdbms/ora12cr1/ora12cr1/trace
```

PL/SQL procedure successfully completed.

Процедура PL/SQL успешно выполнена.

В версии Oracle 11g было добавлено новое средство — **ADR**. Частью его функциональности является новое представление `V$` по имени `V$DIAG_INFO`.

На заметку! `V$DIAG_INFO` — это представление, доступное в Oracle 11g (и последующих версиях), которое в старых выпусках отсутствует. Оно предлагает более легкий интерфейс для работы с трассировочной информацией, используемый новым средством **ADR**.

В целях улучшения читабельности в выводе следующего запроса к `V$DIAG_INFO` длинное путевое имя к домашнему каталогу **ADR** (**ADR Home**) заменено `$home$`. Это просто облегчает его восприятие в книге; поступать так не обязательно:

```
EODA@ORA12CR1> with home
2   as
3   (select value home
4     from v$diag_info
5     where name = 'ADR Home'
6   )
7   select name,
8          case when value <> home.home
9                then replace(value,home.home,'$home$')
10                 else value
11            end value
12   from v$diag_info, home
13  /
```

| NAME | VALUE |
|-----------------------|--|
| ----- | ----- |
| Diag Enabled | TRUE |
| ADR Base | /home/ora12cr1/app/ora12cr1 |
| ADR Home | /home/ora12cr1/app/ora12cr1/diag/ rdbms/ora12cr1/ora12cr1 |
| Diag Trace | \$home\$/trace |
| Diag Alert | \$home\$/alert |
| Diag Incident | \$home\$/incident |
| Diag Cdump | \$home\$/cdump |
| Health Monitor | \$home\$/hm |
| Default Trace File | \$home\$/trace/ora12cr1_ora_22319.trc |
| Active Problem Count | 0 |
| Active Incident Count | 0 |
| 11 rows selected. | |
| 11 строк выбрано. | |

Как видите, строки содержат пути к местоположениям разных трассировочных файлов. В Oracle 11g модернизированы соглашения о местах, где многие файлы хранятся по умолчанию, отличаясь чуть более удобной организацией; это упрощает процесс поддержки, когда вы протоколируете запрос службы в Oracle Support. Ниже описаны наиболее важные строки.

- **Diag Trace (Трассировка диагностики).** Место, куда попадают трассировочные файлы — как фоновые, так и пользовательского дампа — в Oracle 11g и последующих версиях.
- **Default Trace File (Файл трассировки по умолчанию).** Имя трассировочного файла вашего текущего сеанса. В ранних выпусках выяснить это имя было непросто (далее будет показано, каким образом). В Oracle Database 11g и последующих версиях простой запрос к V\$DIAG_INFO возвращает полностью определенное имя файла.

До выхода Oracle 11g и добавления информации Default Trace File вы должны были искать трассировочный файл вручную. Если применяется подключение к разделяемому серверу Oracle, то используется фоновый процесс, так что расположение трассировочных файлов определяется в `background_dump_dest`. В случае подключения к выделенному серверу для взаимодействия с Oracle применяется пользовательский процесс или процесс переднего плана, и трассировочные файлы попадают в каталог, указанный в параметре `user_dump_dest`.

Параметр `core_dump_dest` определяет, где будет сгенерирован файл “ядра” в случае возникновения серьезной внутренней ошибки Oracle (вроде ошибки сегментации в UNIX/Linux), или если в службе поддержки Oracle вас попросят его сгенерировать для получения дополнительной отладочной информации. Вообще говоря, нас интересуют местоположения фонового и пользовательского дампов. Если не указано иное, то на протяжении этой книги мы будем использовать подключения посредством выделенного сервера, так что все трассировочные файлы будут генерироваться в местоположении `user_dump_dest`.

Соглашение об именовании

Соглашение об именовании трассировочных файлов в Oracle время от времени меняется, но если вам известен пример имени трассировочного файла из имеющейся системы, то выяснить применяемый шаблон не составит особого труда. Например, на моих разнообразных серверах UNIX/Linux имена трассировочных файлов выглядят так, как показано в табл. 3.1.

Таблица 3.1. Примеры имен трассировочных файлов

| Имя трассировочного файла | Версия базы данных |
|---------------------------|--------------------|
| ora_10583.trc | 9i Release 1 |
| ora9ir2_ora_1905.trc | 9i Release 2 |
| ora10gr2_ora_6793.trc | 10g Release 2 |
| ora11gr2_ora_1990.trc | 11g Release 2 |
| ora12cr1_ora_2344.trc | 12c Release 1 |

На моих серверах имя трассировочного файла может быть разбито на части следующим образом.

- Первая часть имени файла — ORACLE_SID (за исключением версии Oracle9i Release 1, в которой решили эту часть опустить).
- Следующая часть имени файла — просто ora.
- Номер в имени трассировочного файла — это идентификатор процесса выделенного сервера, доступный через представление V\$PROCESS.

Таким образом, до выхода версии Oracle 11g, где было удобно использовать представление V\$DIAG_INFO, на практике (предполагая режим выделенного сервера) для определения имени трассировочного файла требовался доступ к четырем представлениям:

- V\$PARAMETER, применяемое для нахождения трассировочного файла user_dump_dest и необязательного идентификатора tracefile_identifier, который может быть использован в имени трассировочного файла;
- V\$PROCESS, применяемое для нахождения идентификатора процесса;
- V\$SESSION, используемое для корректной идентификации информации о сессии в других представлениях;
- V\$INSTANCE, применяемое для получения ORACLE_SID.

Как упоминалось ранее, для выяснения местоположения файлов можно использовать утилиту DBMS_UTILITY. Кроме того, часто идентификатор ORACLE_SID просто “известен”, поэтому формально требуется доступ только к представлениям V\$SESSION и V\$PROCESS, но для упрощения работы может понадобиться доступ ко всем четырем представлениям.

Следовательно, запрос для генерации имени трассировочного файла мог бы выглядеть так:

```

EODA@ORA12CR1> column trace new_val TRACE format a100
EODA@ORA12CR1> select c.value || '/' || d.instance_name || '_ora_' ||
a.spid || '.trc' trace
2   from v$process a, v$session b, v$parameter c, v$instance d
3   where a.addr = b.paddr
4     and b.audsid = userenv('sessionid')
5     and c.name = 'user_dump_dest'
6   /

TRACE
-----
/home/ora12cr1/app/ora12cr1/diag/rdbms/ora12cr1/ora12cr1/trace/
ora12cr1_ora_22319.trc

```

Вывод как раз показывает, что если файл существует, то к нему можно обратиться по такому имени (при наличии разрешений на чтение каталога трассировки).

В следующем примере генерируется трассировочный файл, демонстрируя создание файла после включения трассировки:

```

EODA@ORA12CR1> !ls &TRACE
ls: cannot access /home/ora12cr1/app/ora12cr1/diag/rdbms/ora12cr1/
ora12cr1/trace/ora12cr1_ora_22319.trc: No such file or directory
ls: невозможен доступ к /home/ora12cr1/app/ora12cr1/diag/rdbms/ora12cr1/
ora12cr1/trace/ora12cr1_ora_22319.trc: Нет такого файла или каталога
EODA@ORA12CR1> exec dbms_monitor.session_trace_enable
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно выполнена.
EODA@ORA12CR1> !ls &TRACE
/home/ora12cr1/app/ora12cr1/diag/rdbms/ora12cr1/ora12cr1/trace/
ora12cr1_ora_22319.trc

```

Как видите, перед включением трассировки в этом сеансе файл не существовал; однако, как только трассировка включена, мы можем его видеть.

Должно быть очевидным, что в среде Windows символы / понадобятся заменить символами \. В случае применения версии Oracle9i Release 1 вместо добавления имени экземпляра в имя трассировочного файла понадобится просто выдать следующий запрос:

```
select c.value || 'ora_' || a.spid || '.trc'
```

Маркирование трассировочных файлов

Существует способ “маркирования” трассировочных файлов, позволяющий находить их, даже когда доступ к представлениям V\$PROCESS и V\$SESSION не разрешен. При наличии доступа по чтению к каталогу user_dump_dest можно использовать параметр сеанса tracefile_identifier. С его помощью к имени трассировочного файла можно добавить уникально идентифицируемую строку, например:

```

EODA@ORA12CR1> alter session set tracefile_identifier = 'Look_For_Me';
Session altered.
Сеанс изменен.
EODA@ORA12CR1> !ls /home/ora12cr1/app/ora12cr1/diag/rdbms/ora12cr1/
ora12cr1/trace/*Look_For_Me*

```

```
ls: cannot access /home/oral2crl/app/oral2crl/diag/rdbms/oral2crl/
oral2crl/trace/*Look_For_Me*: No such file or directory
ls: невозможен доступ к /home/oral2crl/app/oral2crl/diag/rdbms/oral2crl/
oral2crl/trace/*Look_For_Me*: Нет такого файла или каталога
EODA@ORA12CR1> exec dbms_monitor.session_trace_enable
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно выполнена.
EODA@ORA12CR1> !ls
/home/oral2crl/app/oral2crl/diag/rdbms/oral2crl/oral2crl/trace/
*Look_For_Me*
```

Предшествующая строка кода не уместается в физические границы печатной страницы, поэтому в книге она разнесена на две строки, хотя должна вводиться как одна строка. Команда `ls` ищет файлы в следующем каталоге:

```
/home/oral2crl/app/oral2crl/diag/rdbms/oral2crl/oral2crl/trace
```

Символ `*` представляет собой групповой символ, указывающий команде `ls` на необходимость поиска любых файлов, в имени которых присутствует строка `Look_For_Me`. В этом примере предыдущая команда `ls` обнаруживает два таких файла:

```
/home/oral2crl/app/oral2crl/diag/rdbms/oral2crl/oral2crl/trace/
oral2crl_ora_22489_Look_For_Me.trc
/home/oral2crl/app/oral2crl/diag/rdbms/oral2crl/oral2crl/trace/
oral2crl_ora_22489_Look_For_Me.trm
```

Как видите, теперь трассировочному файлу назначено имя в стандартном формате `<ORACLE_SID>_ora_<PROCESS_ID>`, но с ним связана также заданная уникальная строка, что позволяет без особых затруднений найти имя “нашего” трассировочного файла. Трассировочный файл имеет расширение `.trc`. Имеется также соответствующий файл отображения трассировки (с расширением `.trm`), в котором находится структурная информация о трассировочном файле. Обычно вас будет интересовать только содержимое файла `.trc`.

Трассировочные файлы, генерируемые в ответ на внутренние ошибки

Этот раздел я хотел бы завершить рассмотрением несколько иного вида трассировочных файлов — файлов, генерируемых в результате ошибки `ORA-00600` или какой-то другой внутренней ошибки. Можно ли с ними что-то делать?

Если ответить на такой вопрос коротко, то эти файлы предназначены не для нас с вами — их используют сотрудники службы поддержки Oracle. Однако они могут пригодиться при заполнении запроса службы поддержки Oracle. Важно запомнить следующее: если вы столкнулись с внутренними ошибками, то единственный способ их устранения заключается в заполнении запроса службы поддержки. Если вы просто проигнорируете их, ошибки никогда не исчезнут сами по себе, разве что случайно.

Например, при создании следующей таблицы и выполнении запроса в Oracle 10g Release 1 вполне можно получить сообщение о внутренней ошибке (или же не получить; эта ошибка была выявлена и устранена в последующих исправлениях этого выпуска):

```
ops$tkyte@ORA10G> create table t ( x int primary key );
Table created.
Таблица создана.
```

```
ops$tkyte@ORA10G> insert into t values ( 1 );
1 row created.
1 строка создана.
```

```
ops$tkyte@ORA10G> exec dbms_stats.gather_table_stats( user, 'T' );
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно выполнена.
```

```
ops$tkyte@ORA10G> select count(x) over (
  2   from t;
  from t
  *
```

ERROR at line 2:

ORA-00600: internal error code, arguments: [12410], [], [], [], [], [], [], []

Ошибка в строке 2:

ORA-00600: код внутренней ошибки, аргументы: [12410], [], [], [], [], [], [], []

Теперь представьте, что вы — администратор базы данных, и этот трассировочный файл вдруг появляется в области трассировки. Или же вы — разработчик, и ваше приложение вызывает ошибку ORA-00600, поэтому нужно выяснить, что произошло. Этот трассировочный файл содержит огромный объем информации (около 35 000 строк), но в целом для вас и для меня он не особенно полезен. В общем случае нужно было бы просто сжать его и загрузить как часть запроса службы.

В Oracle 11g и последующих версиях процесс сбора трассировочной информации и ее загрузки в службу поддержки был модифицирован (и стал значительно проще). Новый инструмент командной строки в сочетании с пользовательским интерфейсом диспетчера предприятия позволяет просматривать трассировочную информацию в ADR, упаковывать ее и передавать в службу поддержки Oracle.

Инструмент ADRCI (Automatic Diagnostic Repository Command Interpreter — командный интерпретатор репозитория автоматической диагностики) позволяет просматривать “проблемы” (критические ошибки в базе данных) и инциденты (случаи возникновения этих критических ошибок), а затем упаковывать их для передачи в службу поддержки. Шаг упаковки включает извлечение не только трассировочной информации, но также деталей из сигнального журнала базы данных и прочей конфигурационной/тестовой информации. Например, я создал в своей базе данных условия для возникновения критической ошибки. (Я не собираюсь сообщать какие-либо подробности. Вы должны генерировать собственные критические ошибки.) Я узнал, что в базе данных присутствует “проблема”, т.к. инструмент ADRCI сообщил следующее:

```
[oral2crl@dellpe ~]$ adrci
ADRCI: Release 12.1.0.1.0 - Production on Mon Sep 2 17:45:38 2013
Copyright (c) 1982, 2013, Oracle and/or its affiliates. All rights reserved.
ADR base = "/home/oral2crl/app/oral2crl"
adrci> show problem
```



```
ADR Home = /home/ora12crl/app/ora12crl/diag/rdbms/ora12crl/ora12crl:
*****
PROBLEM_ID PROBLEM_KEY          LAST_INCIDENT LASTINC_TIME
-----
1          ORA 7445 [qctcopn] 36281          2013-09-02 17:52:11.438000 -04:00
```

2 сентября 2013 года я вызвал в базе данных ошибку ORA-4031 — серьезную проблему (ошибка была зафиксирована и исправлена). Теперь я могу посмотреть, что затронула эта ошибка, выдав команду `show incident`:

```
adrci> show incident

ADR Home = /home/ora12crl/app/ora12crl/diag/rdbms/ora12crl/ora12crl:
*****
INCIDENT_ID PROBLEM_KEY          CREATE_TIME
-----
36249       ORA 7445 [qctcopn]      2013-09-02 17:45:06.791000 -04:00
36250       ORA 7445 [qctcopn]      2013-09-02 17:51:58.469000 -04:00
36281       ORA 7445 [qctcopn]      2013-09-02 17:52:11.438000 -04:00
```

Я отмечаю наличие трех инцидентов и с помощью команды `show tracefile` могу вывести информацию, связанную с каждым из них:

```
adrci> show tracefile -I 36250
diag/rdbms/ora12crl/ora12crl/incident/incdir_36250/
ora12crl_ora_22682_i36250.trc
```

Это показывает местоположение трассировочного файла для инцидента номер 36250. Более того, при желании я могу просмотреть массу деталей об инциденте:

```
adrci> show incident -mode detail -p "incident_id=36250"
ADR Home = /home/ora12crl/app/ora12crl/diag/rdbms/ora12crl/ora12crl:
*****

*****
INCIDENT INFO RECORD 1
*****

INCIDENT_ID          36250
STATUS               ready
CREATE_TIME          2013-09-02 17:51:58.469000 -04:00
PROBLEM_ID           1
CLOSE_TIME           <NULL>
FLOOD_CONTROLLED     none
ERROR_FACILITY       ORA
ERROR_NUMBER         7445
ERROR_ARG1           qctcopn
ERROR_ARG2           SIGSEGV
ERROR_ARG3           ADDR:0x18
ERROR_ARG4           PC:0xB859512
ERROR_ARG5           Address not mapped to object
ERROR_ARG6           <NULL>
ERROR_ARG7           <NULL>
ERROR_ARG8           <NULL>
ERROR_ARG9           <NULL>
ERROR_ARG10          <NULL>
ERROR_ARG11          <NULL>
ERROR_ARG12          <NULL>
```

```

SIGNALLING_COMPONENT      <NULL>
SIGNALLING_SUBCOMPONENT   <NULL>
SUSPECT_COMPONENT         <NULL>
SUSPECT_SUBCOMPONENT      <NULL>
ECID                      <NULL>
IMPACTS                   0
PROBLEM_KEY               ORA 7445 [qctcopn]
FIRST_INCIDENT            36249
FIRSTINC_TIME             2013-09-02 17:45:06.791000 -04:00
LAST_INCIDENT             36281
LASTINC_TIME              2013-09-02 17:52:11.438000 -04:00
IMPACT1                   0
IMPACT2                   0
IMPACT3                   0
IMPACT4                   0
KEY_NAME                  Client ProcId
KEY_VALUE                 oracle@dellpe
                           (TNS V1-V3).22682_140239502662112
KEY_NAME                  SID
KEY_VALUE                 416.23
KEY_NAME                  ProcId
KEY_VALUE                 31.41
KEY_NAME                  PQ
KEY_VALUE                 (0, 1378158717)
OWNER_ID                  1
INCIDENT_FILE             /home/ora12crl/app/ora12crl/diag/rdbms/
ora12crl/ora12crl/trace/ora12crl_ora_22682.trc
OWNER_ID                  1
INCIDENT_FILE             /home/ora12crl/app/ora12crl/diag/rdbms/
ora12crl/ora12crl/incident/incdir_36250/ora12crl_ora_22682_i36250.trc

```

И, наконец, я могу создать “пакет”, включающий сведения об инциденте, который будет полезен для службы поддержки. Пакет будет содержать все, что необходимо аналитикам службы поддержки для начала работы над проблемой.

Настоящий раздел не претендует на то, чтобы быть исчерпывающим обзором или введением в утилиту ADRCI, которая полностью документирована в руководстве *Oracle Database Utilities* (Утилиты базы данных Oracle). Вместо этого я просто хотел напомнить о существовании этого инструмента — инструмента, облегчающего использование трассировочных файлов.

До появления утилиты ADRCI в версии Oracle 11g могли ли вы делать с незапланированными трассировочными файлами что-то помимо отправки их в службу поддержки? Да, в них присутствует информация, которая может помочь ответить на вопросы о том, какая произошла ошибка и где именно. Трассировочный файл может также помочь при возникновении ряда других проблем.

Предыдущий пример продемонстрировал, что ADRCI предлагает удобный способ анализа трассировочных файлов в Oracle 12c (здесь была показана лишь малая часть доступных команд). В Oracle 10g и предшествующих версиях вы могли делать то же самое, хотя с большей долей ручного труда. Например, даже беглый взгляд на начало трассировочного файла предоставляет некоторую полезную информацию:

```

/home/ora10gr1/admin/ora10gr1/udump/ora10gr1_ora_2578.trc
Oracle Database 10g Enterprise Edition Release 10.1.0.4.0 - Production

```

```

With the Partitioning, OLAP and Data Mining options
ORACLE_HOME = /home/ora10gr1
System name:   Linux
Node name:    dellpe
Release:      2.6.9-11.ELsmp
Version:      #1 SMP Fri May 20 18:26:27 EDT 2005
Machine:      i686
Instance name: ora10gr1
Redo thread mounted by this instance: 1
Oracle process number: 16
Unix process pid: 2578, image: oracle@dellpe (TNS V1-V3)

```

Информацию о базе данных важно иметь при заполнении запроса службы поддержки на сайте <http://support.oracle.com> или при выяснении, является ли то, с чем вы столкнулись, известной проблемой. Вдобавок вы можете просмотреть экземпляр Oracle, в котором произошла ошибка. Довольно часто параллельно выполняются несколько экземпляров, так что изоляция проблемы в конкретном экземпляре весьма полезна.

А вот другой раздел трассировочного файла, о котором следует знать:

```

*** 2010-01-20 14:32:40.007
*** ACTION NAME: () 2010-01-20 14:32:39.988
*** MODULE NAME: (SQL*Plus) 2010-01-20 14:32:39.988
*** SERVICE NAME: (SYS$USERS) 2010-01-20 14:32:39.988

```

Эта часть трассировочного файла присутствует только в Oracle 10g и последующих версиях; в Oracle9i и предшествующих версиях ее нет. Она показывает доступную информацию о сеансе в столбцах ACTION (действие) и MODULE (модуль) представления V\$SESSION. Здесь можно выяснить, какой сеанс SQL*Plus привел к возникновению ошибки (вы и ваши разработчики можете и должны установить информацию ACTION и MODULE; некоторые среды, такие как Oracle Forms и APEX, делают это автоматически).

Кроме того, здесь есть строка SERVICE NAME (имя службы). Это имя реальной службы, которая применялась для подключения к базе данных (SYS\$USERS в рассматриваемом случае); легко заметить, что подключение осуществлялось не через службу TNS. Если бы мы вошли в базу данных, используя user/pass@ora10g.localdomain, то могли бы видеть такую строку:

```

*** SERVICE NAME: (ORA10G) 2010-01-20 14:32:39.988

```

Здесь ora10g — имя службы (не строка подключения TNS; точнее, оно представляет собой имя конечной службы, зарегистрированной в прослушивающем процессе TNS, к которому произведено подключение). Эта информация полезна также при отслеживании процесса/модуля, на который повлияла возникшая ошибка.

В заключение, прежде чем переходить к действительной ошибке, можно выявить идентификатор сеанса (19 в этом примере), порядковый номер сеанса (27995 в этом примере) и связанную с ним информацию даты/времени (все выпуски):

```

*** SESSION ID: (19.27995) 2010-01-20 14:32:39.988

```

Теперь мы готовы к анализу самой ошибки:

```

ksedmp: internal or fatal error
ksedmp: внутренняя или неисправимая ошибка

```

```

ORA-00600: internal error code, arguments: [12410], [], [], [], [], [], [], []
ORA-00600: код внутренней ошибки, аргументы: [12410], [], [], [], [], [], [], []
Current SQL statement for this session:
Текущий SQL-оператор для этого сеанса:
select count(x) over ()
  from t
----- Call Stack Trace -----
----- Трассировка стека вызовов -----
_ksedmp+524
_ksfdmp.160+14
_kgeriv+139
_kgesiv+78
_ksesic0+59
_qerixAllocate+4155
_qknRwsAllocateTree+281
_qknRwsAllocateTree+252
_qknRwsAllocateTree+252
_qknRwsAllocateTree+252
_qknDoRwsAllocate+9
...

```

Здесь мы видим пару важных фрагментов информации. Первый из них — SQL-оператор, который выполнялся в момент возникновения внутренней ошибки, что очень полезно для отслеживания приложения или приложений, подвергшихся влиянию ошибки. Кроме того, поскольку мы располагаем информацией об SQL-операторе, мы можем начать поиск возможных способов обхода ошибки — опробовать разные способы кодирования SQL-запроса, чтобы выяснить возможность быстрого обхода проблемы на время, пока ведутся работы по устранению внутренней ошибки. Более того, можно вырезать и вставить подозрительный SQL-оператор в среду интерфейса SQL*Plus и посмотреть, удастся ли получить воспроизводимый тестовый сценарий для службы поддержки Oracle (естественно, это наиболее предпочтительный подход).

Вторым важным фрагментом информации является код ошибки (обычно 600, 3113 или 7445) и другие аргументы, связанные с кодом ошибки. Эти сведения в сочетании с информацией трассировки стека, которая показывает набор и порядок вызова внутренних подпрограмм Oracle, могут помочь в нахождении существующей программной ошибки (и способов ее обхода, заплат и т.п.). Например, вот как может выглядеть строка поиска:

```
ora-00600 12410 ksesic0 qerixAllocate qknRwsAllocateTree
```

С помощью расширенного поиска (поиска в базе данных программных ошибок с применением всех слов) на сайте My Oracle Support (Моя поддержка Oracle) мы немедленно находим ошибку 3800614, “ORA-600 [12410] ON SIMPLE QUERY WITH ANALYTIC FUNCTION” (ORA-600 [12410] на простом запросе с аналитической функцией). Если вы перейдете на сайт <http://support.oracle.com> и поищите приведенный текст, то обнаружите сведения об этой ошибке, выясните, что она была исправлена в следующем выпуске, и узнаете, что для нее имеются заплатки — вся эта информация доступна. Я многократно убеждался в том, что возникающая в моей системе ошибка уже случалась ранее, и для нее есть исправления или способы обхода.

Заключительные соображения по поводу трассировочных файлов

Итак, вы теперь знаете два типа основных трассировочных файлов, места их хранения и способы их нахождения. Будем надеяться, что вам придется использовать трассировочные файлы главным образом для настройки и повышения производительности приложений, а не для заполнения запросов в службу поддержки. В заключение хочу отметить, что сотрудники службы поддержки Oracle имеют доступ ко многим недокументированным “событиям”, которые очень полезны для извлечения огромного объема диагностической информации при возникновении любой ошибки в базе данных. Например, если ошибка ORA-01555 Snapshot Too Old (слишком старый снимок) возникла в ситуации, в которой, по вашему мнению, это было совершенно невозможно, то сотрудники службы поддержки помогут выполнить настройку таких диагностических событий для отслеживания точной причины ошибки путем создания трассировочного файла при каждом ее возникновении.

Сигнальный файл

Сигнальный файл (также известный как сигнальный журнал (alert log)) — это своего рода ежедневник базы данных. Он представляет собой простой текстовый файл, заполняемый с момента “рождения” (создания) базы данных и до ее последнего дня (когда она будет удалена). В этом файле вы найдете записанные в хронологическом порядке исторические сведения о базе данных — переключения журнальных файлов, возможные внутренние ошибки, время создания табличных пространств, их отключения и повторного подключения и т.п. Этот файл чрезвычайно полезен для просмотра хронологии базы данных. Я совершенно убежден, что чем больше информации содержит этот файл, тем лучше, поэтому позволяю ему становиться достаточно большими, прежде чем его архивировать.

Я не стану описывать абсолютно все, что попадает в сигнальный журнал — эта тема достаточно обширна. Однако я рекомендую просмотреть его у себя и оценить обилие содержащейся в нем информации. В этом разделе мы рассмотрим конкретный пример извлечения информации из сигнального журнала — в данном случае для создания отчета о работе базы данных.

В прошлом я применял файл сигнального журнала для веб-сайта <http://asktom.oracle.com> и для генерации отчета о работоспособности своей базы данных. Вместо того чтобы вручную проходить по файлу и выяснять, что происходило (в частности, время запуска и остановки базы данных), я решил воспользоваться преимуществами, предлагаемыми СУБД и языком SQL, для автоматизации таких работ. Результатом стала методика создания динамического отчета о работоспособности непосредственно из сигнального журнала.

Команда EXTERNAL TABLE (которая подробно рассматривается в главах 10 и 15) позволяет запрашивать сигнальный журнал и просматривать сохраненную в нем информацию. Я обнаружил, что при каждом запуске базы данных в сигнальный журнал вносилась пара записей:

```
Thu May 6 14:24:42 2004
Starting ORACLE instance (normal)
```

То есть всегда имеется запись метки времени в постоянном формате с фиксированной шириной и сообщение Starting ORACLE instance (Запуск экземпляра

ORACLE). Кроме того, я обратил внимание, что этим записям должно было предшествовать сообщение ALTER DATABASE CLOSE (в случае чистого завершения), сообщение об отмене остановки или вообще ничего — сообщения отсутствовали, свидетельствуя об аварийном отказе системы. Но любое сообщение должно было сопровождаться связанной с ним меткой времени. Таким образом, до тех пор, пока не возник аварийный отказ системы, в сигнальный журнал должна была записываться какая-то значащая метка времени (в случае отказа системы метка времени должна была фиксироваться незадолго до аварии, т.к. запись в журнал производится довольно часто).

Я заметил, что мог бы легко сгенерировать отчет о работоспособности, если бы:

- собрал все записи вроде Staring ORACLE instance %;
- собрал все записи, которые соответствуют формату даты (фактически, представляющие значения даты);
- связал с каждой записью Staring ORACLE instance две предшествующие ей записи (которые были бы записями с датой).

В приведенном ниже коде создается внешняя таблица, позволяющая запрашивать сигнальный журнал. (Обратите внимание, что строка /background/dump/dest должна быть заменена реальным каталогом фонового дампа, а в операторе CREATE TABLE необходимо указать свое имя сигнального журнала.)

```
EODA@ORA12CR1> create or replace
2 directory data_dir
3 as
4 '/home/ora12cr1/app/ora12cr1/diag/rdbms/ora12cr1/ora12cr1/trace/'
5 /
```

Directory created.

Каталог создан.

```
EODA@ORA12CR1> CREATE TABLE alert_log
2 (
3     text_line varchar2(4000)
4 )
5 ORGANIZATION EXTERNAL
6 (
7     TYPE ORACLE_LOADER
8     DEFAULT DIRECTORY data_dir
9     ACCESS PARAMETERS
10    (
11        records delimited by newline
12        fields
13    )
14    LOCATION
15    (
16        'alert_oral2cr1.log'
17    )
18 )
19 reject limit unlimited
20 /
```

Table created.

Таблица создана.

Теперь эту информацию можно запросить в любое время:

```

EODA@ORA12CR1> select to_char(last_time,'dd-mon-yyyy hh24:mi') shutdown,
2      to_char(start_time,'dd-mon-yyyy hh24:mi') startup,
3      round((start_time-last_time)*24*60,2) mins_down,
4      round((last_time-lag(start_time) over (order by r)),2) days_up,
5      case when (lead(r) over (order by r) is null )
6          then round((sysdate-start_time),2)
7          end days_still_up
8  from (
9  select r,
10     to_date(last_time, 'Dy Mon DD HH24:MI:SS YYYY') last_time,
11     to_date(start_time,'Dy Mon DD HH24:MI:SS YYYY') start_time
12  from (
13  select r,
14     text_line,
15     lag(text_line,1) over (order by r) start_time,
16     lag(text_line,2) over (order by r) last_time
17  from (
18  select rownum r, text_line
19  from alert_log
20  where text_line like '_____:__:__ 20__'
21     or text_line like 'Starting ORACLE instance %'
22     )
23  )
24  where text_line like 'Starting ORACLE instance %'
25  )
26  /
```

| SHUTDOWN | STARTUP | MINS_DOWN | DAYS_UP | DAYS_STILL_UP |
|-------------------|-------------------|-----------|---------|---------------|
| ----- | | | | |
| | 28-jun-2013 16:04 | | | |
| 28-jun-2013 17:02 | 28-jun-2013 17:02 | .03 | .04 | |
| 29-jun-2013 06:00 | 01-jul-2013 09:42 | 3102.53 | .54 | |
| 02-jul-2013 14:59 | 02-jul-2013 14:59 | .03 | 1.22 | |
| 02-jul-2013 15:00 | 02-jul-2013 15:00 | .03 | 0 | |
| 02-jul-2013 15:10 | 02-jul-2013 15:10 | .03 | .01 | |
| 02-jul-2013 17:01 | 02-jul-2013 17:02 | 1.55 | .08 | |
| 18-jul-2013 02:00 | 18-jul-2013 11:31 | 571.37 | 15.37 | |
| 05-aug-2013 09:00 | 06-aug-2013 09:06 | 1445.62 | 17.9 | |
| 14-aug-2013 09:09 | 14-aug-2013 09:58 | 49.42 | 8 | |
| 31-aug-2013 14:08 | 02-sep-2013 10:51 | 2683.15 | 17.17 | |
| 02-sep-2013 14:32 | 02-sep-2013 14:51 | 18.93 | .15 | |
| 02-sep-2013 15:13 | 02-sep-2013 15:13 | .03 | .02 | |
| 02-sep-2013 15:15 | 02-sep-2013 15:15 | .05 | 0 | |
| 02-sep-2013 16:53 | 02-sep-2013 16:54 | .03 | .07 | |
| 02-sep-2013 16:56 | 02-sep-2013 16:56 | .03 | 0 | .07 |

16 rows selected.
16 строк выбрано.

Я не стану здесь останавливаться на нюансах SQL-запроса, а отмечу лишь, что самый внутренний запрос (строки 18–21) собирает строки, содержащие “Starting” и дату (вспомните, что использование подчеркивания (_) в конструкции LIKE обеспечивает сопоставление в точности с одним символом — не больше и не меньше).

Запрос также нумерует строки с применением `rownum`. На следующем уровне запроса используется встроенная аналитическая функция `LAG()` для ухода на одну и на две строки для каждой строки и сдвига соответствующих данных вверх, чтобы третья строка запроса имела данные строк 1, 2 и 3. Строка 4 содержит данные строк 2, 3 и 4 и т.д. В итоге мы сохраняем только строки, которые были подобны `Starting ORACLE instance %` и теперь имеют две предшествующих метки времени, ассоциированные с ними. С этого места вычисление времени простоя базы данных не составляет труда: достаточно вычесть одно значение даты из другого. Вычисление времени работы экземпляра не намного сложнее: достаточно уйти на предыдущую строку, получить ее время запуска и вычесть это значение из значения времени остановки, относящегося к этой строке.

Моя база данных Oracle 12c начала свое существование 28 июня 2013 года и останавливалась много раз (к моменту получения этого вывода она функционировала на протяжении 0,07 дня без перерыва).

Если вас интересует другой пример извлечения полезной информации из сигнального журнала, пройдите по ссылке <http://tinyurl.com/y8wkhjt>. На этой странице приведена демонстрация вычисления среднего времени, необходимого для архивирования заданного оперативного журнального файла. Когда вы понимаете, что находится в сигнальном журнале, генерация подобных запросов становится простой.

В дополнение к применению внешней таблицы для опроса сигнальных журналов в Oracle 12c вы можете легко просматривать их с помощью инструмента ADRCL. Он позволяет осуществлять поиск, редактирование (пересмотр) и мониторинг (интерактивное отображение новых записей по мере их появления в журнале). К тому же сигнальный журнал в Oracle 11g и последующих выпусках доступен в двух версиях — в старой, которая только что использовалась, и в версии XML:

```
EODA@ORA12CR1> column value new_val V
EODA@ORA12CR1> select value from v$diag_info where name = 'Diag Alert';
VALUE
-----
/home/oral2cr1/app/oral2cr1/diag/rdbms/oral2cr1/oral2cr1/alert
EODA@ORA12CR1> !ls &V/log.xml
/home/oral2cr1/app/oral2cr1/diag/rdbms/oral2cr1/oral2cr1/alert/log.xml
EODA@ORA12CR1> !head &V/log.xml
<msg time='2013-06-28T16:04:25.378-04:00' org_id='oracle' comp_id='rdbms'
msg_id='dbkrlCheckSuppressAlert:332:7003611' type='NOTIFICATION'
group='startup'
level='16' host_id='localhost.localdomain' host_addr=':::1'
pid='32628' version='1'>
<txt>Adjusting the default value of parameter parallel_max_servers
</txt>
</msg>
<msg time='2013-06-28T16:04:25.378-04:00' org_id='oracle' comp_id='rdbms'
msg_id='dbkrlCheckSuppressAlert:332:2000778772' type='NOTIFICATION'
group='startup'
level='16' host_id='localhost.localdomain' host_addr=':::1'
```


Располагая утилитами либо инструментами для генерации отчетов из XML (например, такими как база данных Oracle, применяющая XDB — XML DB), вы также можете запрашивать и строить отчеты в этом формате.

Естественно, диспетчер предприятия также отображает важную информацию из сигнального журнала.

Файлы данных

Файлы данных наряду с журнальными файлами повторения действий представляют собой наиболее важный набор файлов в базе данных. Именно в этих файлах в конечном итоге будут храниться данные. Каждая база данных имеет, по меньшей мере, один связанный с нею файл данных, и обычно таких файлов будет значительно больше. Только самая простая “тестовая” база данных будет содержать лишь один файл данных. На самом деле, как вы видели в главе 2, даже простейшая команда CREATE DATABASE по умолчанию создает базу данных с тремя файлами данных, которые перечислены ниже:

NAME

```
-----
/home/oral2crl/app/oral2crl/product/12.1.0/dbhome_1/dbs/dbs1ora12c.dbf
/home/oral2crl/app/oral2crl/product/12.1.0/dbhome_1/dbs/dbx1ora12c.dbf
/home/oral2crl/app/oral2crl/product/12.1.0/dbhome_1/dbs/dbu1ora12c.dbf
```

Один файл данных предназначен для табличного пространства SYSTEM (которое вмещает настоящий словарь данных Oracle), один — для табличного пространства SYSAUX (где хранятся другие несловарные объекты в Oracle 10g и последующих версиях) и еще один — для табличного пространства USER (понятие табличного пространства будет раскрыто в разделе “Табличные пространства” далее в главе). Любая реальная база данных будет иметь, по крайней мере, эти три файла данных.

После краткого обзора типов файловой системы мы обсудим организацию этих файлов в базе данных Oracle и организацию данных внутри них. Для понимания этого вы должны знать, что собой представляют табличные пространства, сегменты, экстенды и блоки. Все они являются единицами распределения, которые Oracle использует для хранения объектов в базе данных, и вскоре они будут описаны более подробно.

Краткий обзор механизмов файловой системы

Существуют четыре механизма файловой системы (в версии Oracle 12c их только три), с применением которых хранятся ваши данные в среде Oracle. Под вашими данными понимается словарь данных, данные журнала, данные отката, таблицы, индексы, LOB-объекты и т.д. — данные, о которых вам придется лично заботиться в конце дня. Ниже приведено краткое описание этих типов файлов.

- **“Готовые” файловые системы операционной системы (ОС).** Это файлы, появляющиеся в файловой системе подобно документам текстовых процессоров. Их можно видеть в окне проводника Windows или в среде UNIX/Linux как результат запуска команды `ls`. Для перемещения этих файлов с места на место можно использовать простые утилиты ОС, такие как `xcopy` в Windows или `cp` в UNIX/Linux. Готовые файловые системы ОС исторически являются наиболее

популярным методом хранения данных в Oracle, но я вижу, что с появлением ASM (речь об этом пойдет ниже) ситуация меняется. Готовые файловые системы обычно буферизируются, т.е. ОС будет кешировать информацию при чтении и в ряде случаев при записи на диск.

- **Чистые (raw) разделы.** Это не файлы — это диски без файловой системы. Их нельзя просматривать с помощью команды `ls` или проводника Windows. Они представляют собой всего лишь большие области диска без какой-либо файловой системы внутри. Для Oracle весь чистый раздел выглядит как одиночный крупный файл. Этим он отличается от готовой файловой системы, которая может содержать многие десятки или даже сотни файлов данных. В настоящее время чистые разделы применяются лишь в небольшом проценте развернутых баз данных Oracle из-за присущих им накладных расходов по администрированию. Чистые разделы не являются устройствами с буферизацией — все операции ввода-вывода выполняются напрямую без какой-либо буферизации данных с участием ОС (что применительно к базе данных в целом считается положительным качеством).

На заметку! В версии Oracle 11g чистые разделы были объявлены устаревшими и в Oracle 12c они больше не поддерживаются. При наличии существующей базы данных, в которой используются чистые разделы, вам понадобится извлечь данные с помощью помпы данных или прибегнуть к услугам другого инструмента репликации, такого как Golden Gate, чтобы переместить свои данные в новую базу данных, в которой применяется одна из поддерживаемых файловых систем. В качестве альтернативы в существующую базу данных можно добавить новые пространства имен с поддерживаемой файловой системой и переместить в них нужные данные из чистых разделов. Такой подход работает, только если табличное пространство `SYSTEM` не является чистым разделом.

- **Автоматическое управление памятью (Automatic Storage Management — ASM).** Это новая функциональная возможность Oracle 10g Release 1 (в обеих редакциях Standard и Enterprise). В выпусках, предшествующих Oracle 11g Release 2, ASM представляет собой файловую систему, предназначенную для использования исключительно базой данных. Проще всего думать о ней как о файловой системе базы данных. В этой файловой системе вы не будете хранить, скажем, текстовый файл со списком покупок, а только информацию, связанную с базой данных: таблицы, индексы, резервные копии, управляющие файлы, файлы параметров, журнальные файлы, архивы и т.д. Но даже в ASM существует эквивалент файла данных; концептуально данные по-прежнему хранятся в файлах, но файловой системой является ASM. Система ASM рассчитана на работу либо на отдельном компьютере, либо в кластеризованной среде. Начиная с Oracle 11g Release 2, ASM предоставляет не только файловую систему для базы данных, но также дополнительно кластеризованную файловую систему, описанную ниже.
- **Кластеризованная файловая система.** Эта файловая система предназначена специально для (кластеризованной) среды RAC и предлагает то, что выглядит подобно готовой файловой системе, которая совместно используется многими узлами (компьютерами) в кластеризованной среде. Таким образом,

хотя и можно было бы смонтировать устройство NFS или создать общий ресурс Samba (метод совместного использования дисков в среде Windows/UNIX/Linux, похожий на NFS) для готовой файловой системы с целью доступа из множества узлов кластера, это породило бы одиночную точку отказа. Если узел, который владеет файловой системой и обеспечивает совместную работу с ней, откажет, то файловая система станет недоступной. В версиях Oracle, предшествующих 11g Release 2, в этой области предлагалась кластерная файловая система Oracle (Oracle Cluster File System — OCFS), которая в настоящее время доступна только для Windows и UNIX/Linux. Другие поставщики предлагают сертифицированные кластеризованные файловые системы, также работающие с Oracle. В версии Oracle 11g Release 2 предоставляется еще один вариант в форме кластерной файловой системы для автоматического управления памятью Oracle (Oracle Automatic Storage Management Cluster File System — ACFS). Кластеризованная файловая система привносит комфорт готовой файловой системы в кластеризованную среду.

Интересно отметить, что база данных может включать файлы из любой или всех описанных выше файловых систем — вовсе не обязательно выбирать какую-то одну. Вполне может существовать база данных, внутри которой часть данных хранится в традиционных готовых файловых системах, часть — в чистых разделах, часть — в ASM и часть — в кластеризованной файловой системе. Такой подход облегчает переход с одной технологии на другую или дает представление о возможностях нового типа файловой системы без переноса туда всей базы данных. А теперь, поскольку подробное обсуждение файловых систем и всех их характеристик выходит за рамки настоящей книги, мы возвращаемся к типам файлов Oracle. Независимо от того, где хранится файл — в готовой файловой системе, в чистом разделе, в системе ASM или в кластеризованной файловой системе, — описанные ниже концепции всегда будут применимы к нему.

Иерархия хранения в базе данных Oracle

База данных состоит из одного или более *табличных пространств* (tablespace). Табличное пространство — это логический контейнер хранения в Oracle, который находится на верхушке иерархии хранения и образован из одного или нескольких файлов данных. Файлы могут быть готовыми файлами файловой системы, чистыми разделами, файлами базы данных, управляемой посредством ASM, или файлами в кластеризованной файловой системе. Табличное пространство содержит сегменты, описанные в следующем разделе.

Сегменты

Сегменты (segment) представляют собой основную организационную структуру внутри табличного пространства. Сегменты — это просто объекты базы данных, которые потребляют хранилище, т.е. объекты наподобие таблиц, индексов, сегментов отката и т.д. При создании секционированной таблицы мы создаем не один сегмент для всей таблицы, а по сегменту на секцию. При создании индекса обычно создается сегмент индекса и т.д. Каждый объект, занимающий место в хранилище, в итоге сохраняется в одиночном сегменте. Существуют сегменты отката, временные сегменты, сегменты кластеров, сегменты индексов и прочие.

На заметку! Утверждение о том, что каждый объект, занимающий место в хранилище, в итоге сохраняется в одиночном сегменте, может сбивать с толку. Вы встретите многочисленные операторы CREATE, которые создают многосегментные объекты. Путаница возникает вследствие того факта, что отдельный оператор CREATE может, в конечном счете, создавать объекты, состоящие из нуля, одного или большего числа сегментов! Например, оператор CREATE TABLE T (x int primary key, y clob) создаст четыре сегмента: один для таблицы T, один для индекса, который будет создан для поддержки первичного ключа, и два для объекта CLOB (один сегмент для индекса LOB и один для самих данных LOB). С другой стороны, оператор CREATE TABLE T (x int, y date) cluster by_cluster создаст *ноль* сегментов (в этом случае сегментом является кластер). Мы дополнительно исследуем эту концепцию в главе 10.

Экстененты

Сегменты состоят из одного или более *экстенентов* (extent). Экстенент — это логически непрерывное распределение пространства в файле. (В общем случае сами файлы не занимают непрерывное место на диске, иначе нам никогда не пришлось бы пользоваться инструментом дефрагментации дисков. Кроме того, в случае применения таких дисковых технологий, как RAID (Redundant Array of Independent Disks — избыточный массив независимых дисков), отдельный файл может распространяться на несколько физических дисков.) По традиции каждый сегмент начинается, по крайней мере, с одного экстенента. В версии Oracle 11g Release 2 была введена концепция “отложенного” (deferred) сегмента — сегмента, который не распределяет экстенент немедленно, так что в этой и последующих версиях СУБД сегмент может отложить выделение своего начального экстенента до момента вставки в него данных. Когда объект вырастает до размеров, выходящих за пределы начального экстенента, он запрашивает распределение следующего экстенента. Этот второй экстенент не обязательно будет расположен на диске по соседству с первым — он может даже не попасть в тот же файл, где находится первый экстенент. Второй экстенент может располагаться очень далеко от первого, но пространство внутри экстенента всегда считается логически непрерывным в файле. Размеры экстенентов варьируются от одного блока данных Oracle (рассматривается следующим) до 2 Гбайт.

Блоки

Экстененты, в свою очередь, состоят из *блоков* (block). Блок представляет собой наименьшую единицу распределения пространства в Oracle. Именно в блоках хранятся строки данных, записи индексов или временные результаты сортировки. Блок — это то, что Oracle обычно считывает с диска или записывает на диск. Как правило, блоки в Oracle имеют один из четырех общих размеров: 2 Кбайт, 4 Кбайт, 8 Кбайт или 16 Кбайт (хотя в некоторых случаях допускается также размер 32 Кбайт; ограничения на максимальный размер блоков накладываются операционной системой).

На заметку! Есть один малоизвестный факт: стандартный размер блока для базы данных не обязательно должен быть степенью двойки. Степени двойки — всего лишь общепринятое соглашение. На самом деле можно создать базу данных с размером блока 5 Кбайт, 7 Кбайт или n Кбайт, где n — любое значение из диапазона 2–32 Кбайт. Однако я не рекомендую поступать так в реальных базах данных — придерживайтесь обычных размеров блоков. Использование нестандартных размеров блоков может легко превратиться в проблему при поддержке — если только лично вы применяете размер блока 5 Кбайт, то можете столкнуться со сложностями, которые другие пользователи просто никогда не увидят.

Взаимосвязь между сегментами, экстендами и блоками демонстрируется на рис. 3.1.

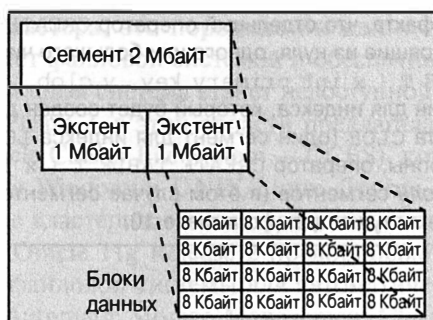


Рис. 3.1. Сегменты, экстенды и блоки

Сегмент образован из одного или более экстендов, а экстенд — это непрерывная область блока. Начиная с версии Oracle9i Release 1, база данных может содержать блоки до шести разных размеров.

На заметку! Поддержка нескольких размеров блоков была введена для того, чтобы переносимые табличные пространства можно было использовать в большем числе ситуаций. Возможность переноса табличного пространства позволяет администратору базы данных перемещать или копировать уже сформатированные файлы данных из одной базы данных и присоединять их к другой — например, чтобы немедленно скопировать все таблицы и индексы из базы данных OLTP (Online Transaction Processing — оперативная обработка транзакций) в базу данных DW (Data Warehouse — информационное хранилище). Однако во многих случаях в базе данных OLTP могут применяться блоки небольшого размера (2 Кбайт или 4 Кбайт), тогда как в базе данных DW будут использоваться блоки значительно большего размера (8 Кбайт или 16 Кбайт). Без поддержки нескольких размеров блоков в одной базе данных перенос такой информации оказался бы невозможным. Табличные пространства с несколькими размерами блоков должны применяться для содействия переносу табличных пространств; обычно они не используются для чего-то другого.

Для базы данных будет определен стандартный размер блока, представляющий собой размер, который указан в файле инициализации, задействованном при выполнении команды CREATE DATABASE. Табличное пространство SYSTEM всегда будет использовать стандартный размер блока, но вы можете создать другие табличные пространства с нестандартными размерами блоков 2 Кбайт, 4 Кбайт, 8 Кбайт, 16 Кбайт и 32 Кбайт (в зависимости от операционной системы). Общее количество размеров блоков равно шести, если и только если во время создания базы данных был указан нестандартный размер блока (не степень двойки). Следовательно, фактически база данных будет иметь максимум пять размеров блоков: один стандартный и четыре нестандартных размера.

Каждое табличное пространство будет обладать согласованным размером блока, т.е. все блоки в этом табличном пространстве будут иметь один и тот же размер. В многосегментном объекте, таком как таблица со столбцом LOB, каждый сегмент в табличном пространстве может иметь отличающийся размер блока, но любой отде-

льно взятый сегмент (который содержится в табличном пространстве) будет состоять из блоков в точности одинаковых размеров. Большинство блоков, независимо от их размеров, имеют один и тот же общий формат, который выглядит подобно показанному на рис. 3.2.



Рис. 3.2. Структура блока

Исключениями из этого формата являются блоки сегментов LOB и блоки гибридных сжатых столбцов в хранилище Exadata, например, но подавляющее большинство блоков в базе данных будут иметь формат, сходный с представленным на рис. 3.2. Заголовок блока содержит информацию о типе блока (блок таблицы, блок индекса и т.д.), информацию об активных и прошедших транзакциях (эту информацию имеют только блоки, управляемые транзакциями — скажем, во временном блоке сортировки она отсутствует) и адрес (местоположение) блока на диске.

Следующие два компонента блока встречаются в наиболее распространенных типах блоков базы данных — блокам традиционных (heap-organized) таблиц. Типы таблиц баз данных подробно рассматриваются в главе 10, а пока достаточно знать, что к этому типу относится большинство таблиц.

Список указателей таблиц, если присутствует, содержит информацию о таблицах, хранящих строки в этом блоке (в одном блоке могут храниться данные из более чем одной таблицы). Список указателей строк содержит информацию, которая описывает строки, находящиеся в блоке. Он представляет собой массив указателей на места, где расположены строки в части данных блока. Эти три компонента блока вместе называются дополнительной памятью блока, являющейся пространством блока, которое не доступно для размещения данных, но используется Oracle для управления самим блоком.

Оставшиеся два фрагмента блока просты: в блоке может присутствовать свободное пространство, и обычно в нем будет находиться область, где в текущий момент хранятся данные.

Теперь, когда вы в общих чертах ознакомились с сегментами, которые состоят из экстенгов, в свою очередь состоящих из блоков, давайте подробнее рассмотрим табличные пространства и то, каким образом файлы вписываются в общую картину.

Табличные пространства

Как отмечалось ранее, табличное пространство представляет собой контейнер — оно хранит сегменты. Каждый сегмент принадлежит только одному табличному пространству. Табличное пространство может содержать внутри себя множество сегментов. Все экстенги отдельного сегмента будут находиться в табличном пространстве, которое связано с этим сегментом. Сегменты никогда не пересекают границы таб-

личных пространств. Само табличное пространство связано с одним или большим числом файлов данных. Экстент для любого заданного сегмента в табличном пространстве будет полностью содержаться внутри одного файла данных. Тем не менее, сегмент может иметь экстенты из множества разных файлов данных. Графически табличное пространство может выглядеть так, как показано на рис. 3.3.

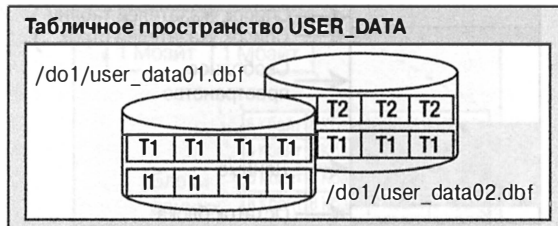


Рис. 3.3. Табличное пространство, содержащее два файла данных, три сегмента и четыре экстента

На рис. 3.3 изображено табличное пространство `USER_DATA`. Оно состоит из двух файлов данных — `user_data01.dbf` и `user_data02.dbf`. Для него выделено три сегмента: `T1`, `T2` и `I1` (вероятно, две таблицы и индекс). В табличном пространстве распределены четыре экстента, каждый из которых представлен как логически непрерывный набор блоков базы данных. Сегмент `T1` содержит два экстента — по одному в каждом файле. Как показано, сегменты `T2` и `I1` имеют по одному экстенту. Если этому табличному пространству понадобится дополнительное место, можно либо изменить размеры уже распределенных для него файлов данных, либо добавить к нему третий файл данных.

Табличное пространство — логический контейнер для хранения данных в Oracle. Как разработчики, мы будем создавать в табличных пространствах сегменты. Нам никогда не придется опускаться до уровня файлов — мы не указываем, что хотим, чтобы наши экстенты были распределены в конкретном файле (делать это можно, но обычно так не поступают). Взамен мы создаем в табличных пространствах объекты, а Oracle заботится обо всем остальном. Если когда-либо в будущем администратор базы данных решит переместить файлы на диске для обеспечения более равномерного распределения операций ввода-вывода, это нас вполне устроит. Оно никак не скажется на обработке данных.

Резюме по иерархии хранения

В качестве резюме по иерархии хранения в Oracle можно привести следующие утверждения.

1. *База данных* состоит из одного или большего числа табличных пространств.
2. *Табличное пространство* образовано одним или несколькими файлами данных. Эти файлы могут быть готовыми файлами в файловой системе, чистыми разделами, файлами базы данных, управляемой ASM, или файлами в кластеризованной файловой системе. Табличное пространство содержит сегменты.
3. *Сегмент* (`TABLE`, `INDEX` и т.д.) формируется одним или большим количеством экстентов. Сегмент существует в табличном пространстве, но его данные могут располагаться во многих файлах данных внутри этого табличного пространства.

4. *Экстент* — логически непрерывный набор блоков на диске. Экстент находится в одном табличном пространстве и, более того, всегда в единственном файле внутри этого табличного пространства.
5. *Блок* — это наименьшая единица распределения пространства в базе данных. Блок является наименьшим элементом ввода-вывода, используемым базой данных при взаимодействии с файлами данных.

Табличные пространства, управляемые словарем и управляемые локально

Прежде чем двигаться дальше, давайте рассмотрим еще одну тему, связанную с табличными пространствами: управление экстенентами в табличном пространстве. До выхода версии Oracle 8.1.5 существовал только один метод управления выделением экстенентов внутри табличного пространства: табличное пространство, управляемое словарем данных (dictionary-managed tablespace). То есть управление памятью внутри табличного пространства производилось в таблицах словаря данных, во многом подобно управлению данными бухгалтерского учета, возможно, с участием таблиц DEBIT (дебет) и CREDIT (кредит). На стороне дебета представлены все экстененты, выделенные объектам. На стороне кредита имеются все свободные экстененты, доступные для использования. Когда объект нуждается в еще одном экстененте, он должен запросить его у системы. Затем Oracle обратится к своим таблицам словаря данных, выполнит ряд запросов, найдет свободное пространство (или нет), после чего обновит строку в одной из таблиц (или удалит ее) и вставит строку в другую таблицу. СУБД Oracle управляет пространством хранения очень похоже на то, как мы поступаем в своих приложениях: изменяя и перемещая данные.

Этот SQL-запрос, выполняемый от нашего имени в фоновом режиме для получения дополнительного пространства, называют рекурсивным SQL-запросом. Наш SQL-оператор INSERT вызывает выполнение еще одного рекурсивного SQL-оператора для получения дополнительного места. Частое выполнение этого рекурсивного SQL-запроса может обходиться достаточно дорого с точки зрения ресурсов. Такие обновления словаря данных должны быть сериализованы; они не могут делаться одновременно. Это то, чего следует избегать.

В ранних версиях Oracle с упомянутой проблемой управления пространством — накладные расходы в виде рекурсивного SQL-запроса — мы сталкивались наиболее часто при работе с временными табличными пространствами (это было до появления “подлинных” временных табличных пространств, создаваемых посредством команды CREATE TEMPORARY TABLESPACE). Пространство должно было регулярно распределяться (приходилось удалять данные из одной таблицы словаря и вставлять их в другую) и освобождаться (только что перемещенные строки приходилось помещать обратно туда, где они располагались первоначально). Как правило, эти операции должны были выполняться последовательно, что вело к резкому снижению степени параллелизма и увеличению времени ожидания. В версии Oracle 7.3 (давным-давно, еще в 1996 году) была введена концепция подлинного временного табличного пространства — нового типа табличного пространства, выделенного специально для хранения временных данных и призванного смягчить описанную проблему.

До появления этого типа табличного пространства временные данные управлялись в тех же самых табличных пространствах, что и постоянные данные и трактовались во многом так же, как постоянные данные.

Во временном табличном пространстве нельзя было создавать собственные постоянные объекты. По существу это было единственным отличием; пространство хранения по-прежнему управлялось посредством таблиц словаря данных. Тем не менее, как только экстенст распределялся во временном табличном пространстве, система должна была удерживать его (т.е. не возвращать пространство обратно). Когда в следующий раз кто-то запрашивал выделение части временного табличного пространства для любых целей, СУБД Oracle должна была искать уже выделенный экстенст в своем внутреннем списке распределенных экстенстов. При обнаружении такого экстенста она должна была просто использовать его повторно или же выделить его старым способом. В результате после запуска базы данных и ее работы в течение определенного периода временный сегмент выглядел бы полным, хотя в действительности он был бы всего лишь “распределенным”. Свободные экстенсты имелись в наличии; они просто управлялись по-другому. Когда кто-то нуждался во временном пространстве, СУБД Oracle должна была искать это пространство в структуре данных, расположенной в памяти, а не выполнять затратный в плане ресурсов рекурсивный SQL-запрос.

В Oracle 8.1.5 и последующих версиях был совершен следующий шаг в направлении сокращения накладных расходов, связанных с управлением пространством. Была введена концепция локально управляемых табличных пространств как противоположность табличным пространствам, управляемым словарем. Локальное управление пространством фактически позволило делать для всех табличных пространств то, что в Oracle 7.3 делалось для временных табличных пространств: оно устранило необходимость применения словаря данных для управления памятью в табличном пространстве.

В локально управляемых табличных пространствах для управления экстенстами используется битовая карта (bitmap), хранящаяся в каждом файле данных. Для получения экстенста системе достаточно лишь установить в 1 бит в битовой карте. Для освобождения пространства система устанавливает бит в 0. По сравнению с применением табличных пространств, управляемых словарем, это происходит невероятно быстро. Больше не нужно сериализовать длительно выполняющиеся операции на уровне базы данных для запросов памяти по всем табличным пространствам. Вместо этого последовательность выполнения обеспечивается на уровне табличного пространства, что гарантирует очень быструю работу. Локально управляемые табличные пространства обладают и другими замечательными характеристиками, такими как принудительное использование унифицированного размера экстенста, но это уже относится к компетенции администратора базы данных.

Забегая наперед, следует отметить, что единственным методом управления пространством хранения, который вы должны применять, является локально управляемое табличное пространство. На самом деле в Oracle9i и последующих версиях при создании базы данных с использованием утилиты DBCA (Database Configuration Assistant — помощник по конфигурированию сервера баз данных) табличное пространство SYSTEM будет создано как локально управляемое, а если это так, то все остальные табличные пространства в базе данных также будут управляться локально,

и унаследованный метод управления посредством словаря работать не будет. Это не означает, что табличные пространства, управляемые словарем, не поддерживаются в базе данных, в которой табличное пространство SYSTEM управляется локально — их просто невозможно создать:

```

EODA@ORA12CR1> create tablespace dmt
  2 datafile '/tmp/dmt.dbf' size 2m
  3 extent management dictionary;
create tablespace dmt
*
ERROR at line 1:
ORA-12913: Cannot create dictionary managed tablespace
ОШИБКА в строке 1:
ORA-12913: невозможно создать табличное пространство, управляемое словарем

EODA@ORA12CR1> !oerr ora 12913
12913, 00000, "Cannot create dictionary managed tablespace"
// *Cause: Attemp to create dictionary managed tablespace in database
//           which has system tablespace as locally managed
// *Action: Create a locally managed tablespace.
12913, 00000, "Невозможно создать табличное пространство, управляемое словарем"
// *Причина: попытка создания табличного пространства, управляемого словарем,
//           в базе данных с локально управляемым табличным пространством SYSTEM
// *Действие: создайте локально управляемое табличное пространство.

```

Обратите внимание, что `oerr` — это утилита, доступная только в UNIX/Linux; на других платформах за описанием полученной ошибки необходимо обращаться в руководство *Oracle Database Error Messages* (Сообщения об ошибках базы данных Oracle).

На заметку! Вас может заинтересовать, по какой причине я утверждаю: “Это не означает, что табличные пространства, управляемые словарем, не поддерживаются в базе данных, в которой табличное пространство SYSTEM управляется локально — их просто невозможно создать”. Если их нельзя создать, то почему нам может понадобиться их поддерживать? Ответ кроется в средстве переносимых табличных пространств. Вы можете переносить табличное пространство, управляемое словарем, в базу данных с табличным пространством SYSTEM, которое управляется локально. Вы можете подключить это табличное пространство и иметь в своей базе данных табличное пространство, управляемое с помощью словаря, но вы не можете *создать* его с нуля в этой базе данных.

Невозможность создания табличных пространств, управляемых словарем, является положительным побочным эффектом, т.к. запрещает пользоваться унаследованным механизмом хранения, который был менее эффективным и опасно предрасположенным к фрагментации пространства. Локально управляемые табличные пространства, вдобавок к более высокой эффективности в плане выделения и освобождения пространства, также предотвращают фрагментацию табличного пространства. Мы рассмотрим их детально в главе 10.

Временные файлы

Временные файлы данных (или просто временные файлы) в Oracle представляют собой специальный тип файла данных. Они применяются для хранения промежуточных результатов масштабных операций сортировки и хеширования, а также для хранения данных глобальной временной таблицы или результирующего набора, когда объема ОЗУ для этого оказывается недостаточно. В Oracle 12c и последующих версиях временные табличные пространства также могут хранить данные UNDO, сгенерированные операциями, которые выполняются в глобальных временных таблицах. В ранних выпусках такие данные направлялись в табличное пространство UNDO, поэтому приводили к генерации данных REDO; теперь это больше не происходит. Во временные файлы никогда не попадут постоянные объекты данных, такие как таблица или индекс, а только содержимое временных таблиц вместе с их индексами. Таким образом, вы никогда не будете создавать во временном файле таблицы своего приложения, но можете хранить в нем данные при использовании временных таблиц.

Временные файлы в Oracle трактуются специальным образом. Обычно каждое изменение, вносимое в объект, будет записываться в журналы повторения. Позже эти журналы транзакций можно воспроизводить для “повторения транзакции”, например, во время восстановления после сбоя. Временные файлы из этого процесса исключены. В частности, для транзакций во временных таблицах (расположенных во временных файлах) данные REDO никогда не генерируются, хотя они могут иметь сгенерированные данные UNDO. Следовательно, *могут существовать сгенерированные данные REDO*, работающие с временными таблицами, т.к. данные UNDO всегда защищены данными REDO, как будет показано в главе 9. Данные UNDO, сгенерированные для глобальных временных таблиц, предназначены для поддержки отката работы, выполненной в текущем сеансе, либо из-за ошибочной обработки данных, либо вследствие какого-то общего отказа транзакции. У администратора базы данных никогда не возникает потребность в резервном копировании временного файла данных; фактически такая попытка была бы напрасной тратой времени, поскольку восстановление временного файла данных совершенно невозможно.

На заметку! В Oracle 12c и последующих версиях данные UNDO, сгенерированные для глобальных временных таблиц, могут быть сохранены во временном табличном пространстве. По умолчанию данные UNDO будут генерироваться в постоянном табличном пространстве UNDO, как это было в предыдущих выпусках. Чтобы разрешить сохранение данных UNDO, сгенерированных для глобальных временных таблиц, во временном файле, можно установить в TRUE настройку `init.ora` уровня системы или устанавливаемый параметр `TEMP_UNDO_ENABLED` уровня сеанса. В таком случае для этих операций данные REDO генерироваться не будут. Мы возвратимся к этой теме в главе 9.

Рекомендуется сконфигурировать базу данных на применение локально управляемых временных табличных пространств. Необходимо гарантировать использование администратором базы данных команды `CREATE TEMPORARY TABLESPACE`. Совершенно нежелательно просто заменить постоянное табличное пространство временным, т.к. при этом невозможно получить все преимущества, связанные с временными файлами.

Одна из особенностей подлинных временных файлов заключается в том, что если операционная система разрешает это, то временные файлы будут создаваться разреженными — т.е. в действительности они не будут потреблять дисковое пространство до тех пор, пока в этом не возникнет необходимость. В этом легко убедиться с помощью следующего примера (в Oracle Linux):

```
EODA@ORA12CR1> !df -h /tmp
Filesystem              Size  Used Avail Use% Mounted on
/dev/mapper/VolGroup-lv_root
                        50G   6.5G   41G   14%  /

EODA@ORA12CR1> create temporary tablespace temp_huge
  2 tempfile '/tmp/temp_huge.dbf' size 2g;
Tablespace created.
Табличное пространство создано.

EODA@ORA12CR1> !df -h /tmp
Filesystem              Size  Used Avail Use% Mounted on
/dev/mapper/VolGroup-lv_root
                        50G   6.5G   41G   14%  /

EODA@ORA12CR1> !ls -l /tmp/temp_huge.dbf
-rw-rw----. 1 oral2crl oral2crl 2147491840 Sep  3 13:28 /tmp/temp_huge.dbf
```

На заметку! Команда UNIX/Linux под названием `df` показывает свободное пространство на диске (от “disk free”). В ее выводе видно, что перед добавлением в базу данных временного файла размером 2 Гбайт в файловой системе, содержащей `/tmp`, было свободно 41 Гбайт. После добавления временного файла объем свободного дискового пространства файловой системы *по-прежнему* остался равным 41 Гбайт.

Как видите, для этого файла не понадобилось много места. Если мы взглянем на вывод `ls`, то увидим нормальный файл размером 2 Гбайт, но фактически в текущий момент он потребляет всего несколько килобайт дискового пространства. Так что на самом деле мы могли бы создать сотни таких временных файлов по 2 Гбайт, несмотря на то, что на диске свободно примерно 41 Гбайт. Звучит неплохо — бесплатное хранилище для всех! Однако проблема в том, что как только мы приступим к использованию этих временных файлов, и они начнут увеличиваться в размерах, то очень скоро мы получим сообщение об ошибке типа “свободного места больше нет”. Поскольку пространство выделяется или физически назначается файлу операционной системой по мере необходимости, есть реальный шанс нехватки свободного места (особенно если после создания временных файлов кто-то другой заполнит файловую систему своими данными).

Способ решения такой проблемы варьируется в зависимости от операционной системы. В среде UNIX/Linux для этого можно выдать команду `dd`, заполнив файл данными, что вынудит ОС физически назначить файлу дисковое пространство, или применить команду `cp`, чтобы создать неразрезанный файл, например:

```
EODA@ORA12CR1> !cp --sparse=never /tmp/temp_huge.dbf /tmp/temp_huge_not_sparse.dbf
EODA@ORA12CR1> !df -h /tmp
Filesystem              Size  Used Avail Use% Mounted on
/dev/mapper/VolGroup-lv_root
                        50G   8.5G   39G   19%  /
```

```

EODA@ORA12CR1> drop tablespace temp_huge including contents and datafiles;
Tablespace dropped.
Табличное пространство удалено.
EODA@ORA12CR1> create temporary tablespace temp_huge
2 tempfile '/tmp/temp_huge_not_sparse.dbf' reuse;
Tablespace created.

```

После копирования разреженного файла размером 2 Гбайт в /tmp/temp_huge_not_sparse.dbf и создания с помощью опции REUSE временного табличного пространства, использующего этот временный файл, мы гарантируем, что временному файлу выделено все нужное ему пространство файловой системы, а база данных действительно имеет 2 Гбайт временного пространства для работы с ним.

На заметку! Согласно моему опыту, файловая система NTFS в Windows не создает разреженные файлы, и это применимо к разновидностям UNIX/Linux. В качестве положительной стороны, если вы должны создать временное табличное пространство размером 15 Гбайт в системе UNIX/Linux и располагаете поддержкой временных файлов, то обнаружите, что это происходит очень быстро (мгновенно); просто удостоверьтесь в наличии 15 Гбайт свободного пространства и держите их в уме.

Управляющие файлы

Управляющие файлы — это довольно небольшие файлы (в предельных случаях они могут вырастать до 64 Мбайт), которые содержат каталог других файлов, необходимых Oracle. Файл параметров сообщает экземпляру, где находятся управляющие файлы, а управляющие файлы указывают экземпляру местоположения файлов базы данных и оперативных журналов повторения транзакций.

Управляющие файлы также сообщают Oracle другие сведения, такие как информация об имеющихся контрольных точках, имя базы данных (которое должно соответствовать указанному в параметре db_name внутри файла параметров), метка времени создания базы данных, хронология архивного журнала повторения транзакций (в некоторых случаях это может приводить к получению управляющего файла большего размера), информация RMAN и т.д.

Управляющие файлы должны мультиплексироваться либо оборудованием (RAID), либо самой базой данных Oracle, когда поддержка RAID или зеркального отображения недоступна. Для каждого управляющего файла должно существовать более одной копии, которые во избежание их утери в случае отказа диска должны храниться на разных дисках. Утеря управляющего файла не является фатальной — она лишь значительно затрудняет процесс восстановления.

Скорее всего, разработчику никогда не придется иметь дело с управляющими файлами. Для администратора базы данных они являются важной частью базы данных, но для разработчика приложений они не представляют особого интереса.

Файлы журналов повторения транзакций

Файлы журналов повторения транзакций критически важны для базы данных Oracle. Они представляют собой журналы транзакций для базы данных. В основном

эти файлы используются только при восстановлении, но могут также применяться в перечисленных ниже ситуациях:

- восстановление экземпляра после аварийного отказа системы;
- восстановление носителей информации после восстановления файла данных из резервной копии;
- обслуживание резервной базы данных;
- подключение к Streams, или процессам извлечения информации из журналов повторения транзакций Golden Gate для совместного использования информации (это всего лишь необычное название репликации);
- предоставление администраторам возможности инспектирования транзакций хронологической базы данных посредством утилиты Oracle LogMiner.

Главным предназначением этих файлов является их применение в случае отказа экземпляра или носителя либо в качестве метода поддержания резервной базы данных для обхода отказа. Если на компьютере базы данных отключилось электропитание, в результате чего произошел сбой в работе экземпляра, то Oracle будет использовать оперативные журналы повторения для восстановления системы в состоянии, которое непосредственно предшествовало моменту исчезновения питания. Если диск, содержащий файлы данных, серьезно пострадал, то для восстановления резервной копии этого диска на соответствующий момент времени вдобавок к оперативным журналам Oracle будет применять архивные журналы повторения. Кроме того, если вы случайно удалили таблицу или какую-то важную информацию и зафиксировали операцию, то можете взять резервную копию и запросить у Oracle восстановление данных на момент прямо перед этим происшествием с использованием оперативных и архивных журналов повторения.

Практически каждая операция, выполняемая в Oracle, генерирует какой-то объем данных, которые записываются в журнальные файлы повторения транзакций. При вставке строки в таблицу конечный результат этой вставки записывается в журналы повторений. При удалении строки факт ее удаления также фиксируется. Результаты удаления таблицы аналогичным образом записываются в журналы повторений. Данные из удаленной таблицы не фиксируются, но рекурсивный SQL-запрос, выполняемый Oracle для удаления таблицы, генерирует информацию повторения. Например, Oracle удалит строку из таблицы SYS.OBJ\$ (и других внутренних объектов словаря); это приведет к генерированию данных повторения транзакции, и если различные режимы дополнительной протоколирования включены, то действительный оператор DROP TABLE будет записан в поток журнала повторения транзакций.

Некоторые операции могут выполняться в режиме генерации минимально необходимого объема данных повторения. Например, я могу создать индекс с атрибутом NOLOGGING. Это означает, что сведения о первоначальном создании данных индекса не будут записаны в журнал, но любые рекурсивные SQL-запросы, выполняемые Oracle от моего имени, протоколироваться будут. Например, вставка в таблицу SYS.OBJ\$ строки, свидетельствующей о существовании индекса, будет зафиксирована в журнале, как и все последующие изменения индекса, вызванные выполнением SQL-операторов вставки, обновления и удаления. Но факт первоначального записывания структуры индекса на диск в журнале не фиксируется.

Я упоминал о двух типах файлов журналов повторения — оперативных и архивных. Мы подробно рассмотрим каждый из них в последующих разделах. В главе 9 мы еще раз взглянем на журналы повторения в сочетании с сегментами отката, чтобы выяснить, какое воздействие они оказывают на разработчиков. А пока мы просто сосредоточим внимание на их сущности и назначении.

Оперативный журнал повторения транзакций

Каждая база данных Oracle имеет, по крайней мере, две группы файлов оперативных журналов повторения. Каждая группа журналов повторения состоит из одного или нескольких элементов (управление журналами повторения осуществляется на уровне групп элементов). Отдельные элементы этих групп журнальных файлов являются точными зеркальными образами друг друга. Файлы оперативного журнала повторения фиксированы по размеру и используются циклически. СУБД Oracle будет выполнять запись в группу 1 журнальных файлов, а когда дойдет до конца этого набора файлов, перейдет к группе 2 и перепишет все содержимое этих файлов от начала до конца. Заполнив группу 2 журнальных файлов, она снова вернется к группе 1 (при наличии только двух групп журнальных файлов; естественно, при наличии трех групп журнальных файлов программа перейдет к группе 3). Этот процесс показан на рис. 3.4.

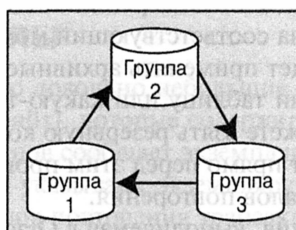


Рис. 3.4. Запись в группы журнальных файлов

Действие по переходу от одной группы журнальных файлов к другой называют *переключением журнальных файлов*. Важно отметить, что в неудачно сконфигурированной базе данных переключение журнальных файлов может вызывать временную “паузу”. Поскольку журналы повторения служат для восстановления транзакций в случае отказа, мы должны быть уверены в том, что не нуждаемся в содержимом журнального файла повторения до того, как появится возможность его использования. Если такой уверенности нет, Oracle будет на мгновение приостанавливать операции в базе данных и удостоверяться в том, что данные в кеше, “защищаемые” этим журналом повторения, надежно записаны на диск (создана контрольная точка). После этого обработка будет возобновлена и файл журнала повторения будет использован повторно.

Мы только что затронули одну из основных концепций баз данных: *создание контрольной точки*. Чтобы понять, как применяются оперативные журналы повторения, необходимо знать кое-что о контрольных точках, о работе кеша буфера базы данных, а также о том, что делает процесс под названием Database Block Writer (процесс записи блока базы данных), или DBWn. Кеш буфера базы данных и процесс DBWn подробно раскрываются позже, но мы забежим немного вперед и кратко рассмотрим их сейчас.

Кеш буфера базы данных — это место временного хранения блоков базы данных. Он представляет собой структуру в области SGA. По мере чтения блоки сохраняются в этом кеше в надежде, что позже не придется повторять процесс их физического считывания. Кеш буфера — самое первое устройство, на которое следует обратить внимание при настройке производительности. Оно существует единственно ради того, чтобы очень медленный процесс физического ввода-вывода выглядел как выполняющийся намного быстрее, чем есть на самом деле. При модификации блока путем обновления в нем строки изменения производятся в памяти — в блоках кеша буфера. При этом объем информации, достаточный для повторения данного изменения, сохраняется в буфере журнала повторения транзакций — еще одной структуре данных области SGA. При фиксации изменений, делающей их постоянными, Oracle не проходит по всем измененным блокам в SGA с целью их записи на диск, а только записывает содержимое буфера журнала повторения в оперативные журналы повторения. До тех пор, пока модифицированный блок находится в кеше буфера, а не на диске, мы нуждаемся в содержимом оперативного журнала на случай отказа базы данных. Если немедленно после фиксации изменений отключить электропитание, кеш буфера базы данных будет очищен.

Когда такое произойдет, единственным свидетельством выполненного изменения является запись в файле журнала повторения. После перезапуска базы данных Oracle действительно воспроизведет транзакцию, модифицируя блок снова, как мы это делали, и зафиксирует изменение. Таким образом, пока измененный блок хранится в кеше и не записан на диск, мы не можем повторно использовать (перезаписывать) файл журнала изменений.

Именно здесь в игру вступает процесс DBWn. Этот фоновый процесс Oracle отвечает за освобождение места в кеше буфера, когда он заполняется, и, что еще важнее, за создание контрольных точек. Создание контрольной точки — это запись “грязных” (измененных) блоков из кеша буфера на диск. Это делается Oracle автоматически в фоновом режиме. Запись контрольной точки может быть вызвана многими причинами, наиболее распространенной из которых является переключение журнальных файлов.

После заполнения журнального файла 1 и переключения на журнальный файл 2 база данных Oracle инициирует создание контрольной точки. В этот момент процесс DBWn начинает записывать на диск все “грязные” блоки, защищенные группой журнальных файлов 1. До тех пор, пока DBWn не сбросит на диск все эти блоки, защищенные упомянутым журнальным файлом, Oracle не может применять его повторно (перезаписывать). При попытке воспользоваться журнальным файлом до завершения создания контрольной точки процессом DBWn в сигнальном журнале базы данных появится сообщение, подобное показанному ниже:

...

Thread 1 cannot allocate new log, sequence 66
Checkpoint not complete

Current log# 2 seq# 65 mem# 0:
/home/oral2crl/app/oral2crl/oradata/orcl/redo01.log

*Поток 1 не может разместить новый журнал, последовательность 66
Создание контрольной точки не завершено*

*Номер текущего журнала 2 последовательности 65 области памяти 0:
/home/oral2crl/app/oral2crl/oradata/orcl/redo01.log*

Итак, когда появилось такое сообщение, обработка в базе данных была приостановлена на время, пока процесс DBWn торопливо создавал контрольную точку. В этот момент Oracle предоставляет процессу DBWn всю доступную вычислительную мощь, чтобы он мог как можно быстрее завершить свою работу.

Это сообщение никогда не должно появляться в хорошо настроенном экземпляре базы данных. Если вы видите его, то знайте, что вы ввели искусственные и нежелательные ожидания при работе конечных пользователей. Этого всегда можно избежать. Цель (причем администратора базы данных, а не обязательно разработчика) заключается в обеспечении достаточного количества журнальных файлов, чтобы не приходилось пытаться повторно использовать один из них до завершения создания контрольной точки (инициированного переключением журнальных файлов). Частое появление такого сообщения свидетельствует о том, что администратор базы данных не выделил приложению достаточное число оперативных журналов повторения транзакций, или о том, что процесс DBWn нуждается в настройке для более эффективной работы.

Разные приложения будут генерировать разное количество журналов повторения. Система поддержки принятия решений (Decision Support System — DSS) или система информационного хранилища (DW) по своей природе ежедневно будет генерировать значительно меньше данных для оперативных журналов повторения, чем система OLTP. Система, которая выполняет множество манипуляций с изображениями в больших двоичных объектах (Binary Large Objects — BLOB) базы данных, может генерировать радикально больше данных повторения, чем простая система обработки заказов. Система обработки заказов, имеющая сто пользователей, скорее всего, будет генерировать десятую часть объема данных повторения, которые генерировала бы система с тысячей пользователей. Таким образом, не существует какого-то “правильного” размера для журналов повторения, хотя вы должны стремиться обеспечить им такой размер, который достаточен для конкретной рабочей нагрузки.

При установке размеров и количества оперативных журналов повторения вы должны принимать во внимание множество факторов. Многие факторы выходят за рамки тематики этой книги, но некоторые из них перечислены ниже, чтобы вы могли получить общее представление.

- *Пиковые рабочие нагрузки.* Желательно, чтобы системе не приходилось дожидаться завершения создания контрольных точек и чтобы в ней не возникали узкие места во время пиковых нагрузок. Вы должны определять размер журналов повторения в расчете не на среднюю часовую пропускную способность, а на пиковую нагрузку. Если на протяжении дня генерируется 24 Гбайт журнальных данных, но 10 Гбайт из них попадают в период с 9:00 до 11:00, то размеры журнальных файлов должны быть достаточно большими, чтобы выдерживать эту двухчасовую пиковую нагрузку. Установка их размеров, исходя из средней величины 1 Гбайт в час, вероятно, окажется недостаточной.
- *Наличие множества пользователей, изменяющих одни и те же блоки.* Это также может требовать применения больших журнальных файлов. После того, как каждый из них модифицирует те же самые блоки, желательно обеспечить максимально возможное количество их обновлений перед записью на диск. Каждое переключение журнальных файлов будет инициировать создание контрольной точки, поэтому неплохо сделать так, чтобы переключение происходило как можно реже. Тем не менее, это может повлиять на время восстановления.

- *Среднее время восстановления.* Если вы должны обеспечить минимальное время восстановления, то можете склониться в сторону использования меньших по размеру журнальных файлов даже с учетом предыдущего обстоятельства. При восстановлении обработка одного или двух небольших журнальных файлов потребует меньшего времени, чем обработка одного гигантского файла. Возможно, это приведет к общему замедлению повседневной работы системы (из-за создания дополнительных контрольных точек), но время, уходящее на восстановление, сократится. Впрочем, существуют также и другие параметры базы данных, которые можно применять для сокращения времени восстановления в качестве альтернативы использованию небольших журнальных файлов.

Архивный журнал повторения транзакций

База данных Oracle может функционировать в одном из двух режимов: ARCHIVELOG и NOARCHIVELOG. Отличие между ними связано с тем, что происходит с файлом журнала, когда Oracle собирается использовать его повторно. При этом необходимо ответить на важный вопрос: “Сохраним ли мы копию этого журнала повторения или пусть Oracle просто перезапишет его, утратив навсегда?” Если вы не сохраните этот файл, то восстановление данных на текущий момент времени из резервной копии окажется невозможным.

Предположим, что резервное копирование выполняется раз в неделю по субботам. В полдень пятницы, после того как за неделю были сгенерированы сотни журнальных файлов, жесткий диск отказывает. Если база данных не работает в режиме ARCHIVELOG, то у вас есть только две возможности.

- Удалить табличное пространство или пространства, ассоциированные с отказавшим диском. Любое табличное пространство, которое имеет файл на этом диске, должно быть удалено вместе со всем его содержимым. Если отказ диска затронул табличное пространство SYSTEM (словарь данных Oracle) или какое-то другое табличное пространство, связанное с системой, например, UNDO, то сделать это не получится. В таком случае придется воспользоваться следующей возможностью.
- Восстановить данные, сохраненные в предыдущую субботу, и потерять результаты работы, выполненной за неделю.

Ни один из вариантов не выглядит привлекательным. Оба они предполагают утрату данных. С другой стороны, если бы база данных работала в режиме ARCHIVELOG, то достаточно было бы найти другой диск и восстановить на него испорченные файлы из субботней резервной копии. Затем нужно было бы применить архивные журнальные файлы и, в конце концов, оперативные журналы повторения (в сущности, воспроизводя в ускоренном режиме выполненные за неделю транзакции). При этом ничего не теряется. Данные восстанавливаются на момент отказа диска.

Люди часто говорят мне, что не нуждаются в режиме ARCHIVELOG для их производственных систем. Я еще ни разу не встретил того, кто был бы корректен при таком заявлении. Я уверен в том, что система не может считаться производственной, если она не функционирует в режиме ARCHIVELOG. База данных, которая не работает в режиме ARCHIVELOG, однажды потеряет данные. Это неизбежно; *вы потеряете данные* (не возможно, а *обязательно*), если ваша база данных не выполняется в режиме ARCHIVELOG.

Обычно в качестве оправдания можно услышать: “Мы используем RAID-5, поэтому полностью защищены”. Мне приходилось сталкиваться с ситуациями, когда из-за производственных дефектов все диски в массиве RAID останавливались, причем почти одновременно. В некоторых случаях аппаратный контроллер приводил к повреждению файлов данных, в результате чего устройства RAID надежно защищали уже поврежденные данные. Кроме того, система RAID не защищает от ошибок операторов — одной из наиболее распространенных причин потери данных. Применение RAID не означает, что данные находятся в безопасности; данные могут быть доступнее, они могут быть *более защищены*, но данные, хранящиеся только на устройстве RAID, когда-нибудь будут утеряны — это лишь вопрос времени.

Очень часто можно услышать: “Если бы у нас были резервные копии базы данных на момент, предшествующий аппаратному сбою или ошибке оператора, то мы смогли бы восстановить данные”. По этому поводу можно лишь сказать, что отказ от использования режима ARCHIVELOG в системе, содержащей сколько-нибудь ценные данные, нельзя ничем оправдать. Стремление к высокой производительности оправданием не является; правильно сконфигурированная архивация увеличивает накладные расходы лишь незначительно либо вовсе не влияет на них. Указанное обстоятельство и тот факт, что быстрая система, которая теряет данные, совершенно бесполезна, обуславливает необходимость применения архивации даже в том случае, если это приведет к стопроцентному росту накладных расходов. Функциональное средство можно отнести к накладным расходам, если от него можно отказаться без утери чего-либо важного; накладные расходы подобны сахарной глазури на пирожном. Сохранение данных и обеспечение невозможности их потери не является накладными расходами — наоборот, это первоочередная задача администратора базы данных!

В режиме NOARCHIVELOG могут находиться только тестовые системы или *возможно* системы для разработки приложений. Большинство систем разработки должны функционировать в режиме ARCHIVELOG по двум причинам.

- Это способ обработки данных в производственной среде; вы хотите, чтобы система разработки действовала и реагировала так, как будет себя вести реальная рабочая система.
- Во многих случаях разработчики извлекают свой код из словаря данных, модифицируют его и компилируют обратно в базу. База разработки хранит текущую версию кода. Если база разработки подвергается отказу диска после обеда, что произойдет со всем кодом, который вы компилировали и повторно компилировали все утро? Он будет утерян.

Не позволяйте никому отговаривать вас от использования режима ARCHIVELOG. На разработку приложения было потрачено много времени, поэтому вы должны стремиться к тому, чтобы пользователи ему доверяли. Потеря их данных никоим образом не способствует повышению доверия к системе.

На заметку! В некоторых случаях работа крупного информационного хранилища в режиме NOARCHIVELOG может оказаться оправданной, если в нем продумано применяются табличные пространства READ ONLY и если любое табличное пространство READ WRITE, пострадавшее от сбоя, должно быть полностью перестроено путем повторной загрузки данных.

Файлы паролей

Файл паролей — это необязательный файл, который позволяет удаленному администратору или пользователю SYSDBA получать доступ к базе данных.

Когда вы пытаетесь запустить Oracle, не существует какой-то доступной базы данных, к которой можно было бы обратиться для проверки подлинности паролей. При запуске Oracle в локальной системе (т.е. не по сети, а с компьютера, на котором будет располагаться экземпляр базы данных) для выполнения действий по аутентификации Oracle будет использовать средства ОС.

При установке Oracle необходимо указать группу ОС для администраторов. Обычно по умолчанию в системе UNIX/Linux этой группой будет DBA, а в системе Windows — ORA_DBA. Тем не менее, это может быть любое имя группы, допустимое на текущей платформе. Такая группа обладает “особым” статусом в том смысле, что любой входящий в нее пользователь может подключаться к базе данных “как SYSDBA”, не указывая имя пользователя или пароль. Например, в своей установленной копии Oracle 12c Release 1 я задал группу `oral2crl`. Любой пользователь, входящий в эту группу, может подключаться без указания имени пользователя/пароля:

```
[tkyte@dellpe ~]$ groups
tkyte oral2crl orallgr2 oral0gr2
[tkyte@dellpe ~]$ sqlplus / as sysdba

SQL*Plus: Release 12.1.0.1.0 Production on Tue Sep 3 14:15:31 2013
Copyright (c) 1982, 2013, Oracle. All rights reserved.
```

```
Connected to:
Oracle Database 12c Enterprise Edition Release 12.1.0.1.0 - 64bit Production
With the Partitioning, OLAP, Advanced Analytics and Real Application
Testing options
SYS@ORA12CR1> show user
USER is "SYS"
```

Как видите, все работает. Я подключился и теперь могу запускать базу данных, останавливать ее или выполнять любые действия по администрированию. Однако предположим, что эти операции необходимо выполнять с другого компьютера через сеть. В этом случае я пробую подключиться с применением строки соединения TNS. Однако попытка оказывается неудачной:

```
[tkyte@dellpe ~]$ sqlplus /@oral2crl as sysdba

SQL*Plus: Release 12.1.0.1.0 Production on Tue Sep 3 14:16:22 2013
Copyright (c) 1982, 2013, Oracle. All rights reserved.

ERROR:
ORA-01017: invalid username/password; logon denied
ОШИБКА:
ORA-01017: недопустимое имя пользователя/пароль; вход запрещен
```

Аутентификация со стороны ОС пользователя SYSDBA через сеть не будет работать даже при установке в TRUE значения очень нежелательного (с точки зрения безопасности) параметра `REMOTE_OS_AUTHENT`. Таким образом, аутентификация со стороны ОС не работает и, как обсуждалось ранее, когда вы пытаетесь запустить эк-

земплар для монтирования и открытия базы данных, то по определению пока еще нет какой-либо базы данных, где бы можно было найти сведения по аутентификации. Это один из примеров пресловутой проблемы с выяснением, что первично — курица или яйцо.

На выручку приходит файл паролей. В этом файле хранится список имен пользователей и паролей, которым разрешена удаленная аутентификация через сеть как SYSDBA. Для их аутентификации СУБД Oracle должна использовать этот файл, а не обычный список паролей, хранящийся в базе данных.

Итак, давайте исправим ситуацию. Для начала проверьте, что параметр `REMOTE_LOGIN_PASSWORDFILE` установлен в свое стандартное значение `EXCLUSIVE`, которое предусматривает использование предоставленного файла паролей только одной базой данных:

```
EODA@ORA12CR1> show parameter remote_login_passwordfile
```

| NAME | TYPE | VALUE |
|---------------------------|--------|-----------|
| remote_login_passwordfile | string | EXCLUSIVE |

На заметку! Другими допустимыми значениями для этого параметра являются `NONE`, указывающее на отсутствие файла паролей (нет удаленных подключений SYSDBA) и `SHARED` (один файл паролей может быть задействован более, чем одной базой данных).

Затем с помощью утилиты командной строки (в UNIX/Linux и Windows) под названием `orapwd` необходимо создать и заполнить файл паролей:

```
[oral2crl@dellpe ~]$ orapwd
Usage: orapwd file=<fname> entries=<users> force=<y/n> asm=<y/n>
dbuniquename=<dbname> format=<legacy/12> sysbackup=<y/n> sysdg=<y/n>
syskm=<y/n> delete=<y/n> input_file=<input-fname>

Использование: orapwd file=<имя_файла> entries=<пользователи> force=<y/n>
asm=<y/n> dbuniquename=<имя_базы_данных> format=<legacy/12>
sysbackup=<y/n> sysdg=<y/n> syskm=<y/n> delete=<y/n>
input_file=<имя_входного_файла>
```

```
Usage: orapwd describe file=<fname>
Использование: orapwd describe file=<имя_файла>
```

```
where
```

```
...
```

```
There must be no spaces around the equal-to (=) character.
```

```
где
```

```
...
```

Пробелы возле знака равенства (=) не разрешены.

Команда, которую мы будем применять при входе с учетной записью операционной системы, владеющей программным обеспечением Oracle, выглядит следующим образом:

```
[oral2crl@dellpe dbs]$ orapwd file=orapw$ORACLE_SID password=bar
entries=20
```

В моем случае создается файл паролей по имени `orapworal2crl` (идентификатором `ORACLE_SID` у меня является `oral2crl`). Это соответствует соглашению

об именовании данного файла, используемому на большинстве платформ UNIX/Linux (детали, связанные с именованием этого файла для конкретной платформы можно найти в руководстве по установке и администрированию ОС). Файл размещен в каталоге `$ORACLE_HOME/dbs`. В системе Windows этот файл имеет имя `PW$ORACLE_SID%.ora` и находится в каталоге `%ORACLE_HOME%\database`. Прежде чем запускать команду для создания данного файла, необходимо перейти в правильный каталог или переместить файл в этот каталог впоследствии.

В настоящее время в файле присутствует только один пользователь SYS, несмотря на то, что база данных содержит и другие учетные записи SYSDBA (они пока еще не записаны в файл паролей). Однако, располагая этой информацией, можно выполнить первое подключение к базе данных через сеть в качестве пользователя SYSDBA:

```
[tkyte@dellpe ~]$ sqlplus sys/bar@ora12crl as sysdba
SQL*Plus: Release 12.1.0.1.0 Production on Tue Sep 3 14:21:08 2013
Copyright (c) 1982, 2013, Oracle. All rights reserved.

Connected to:
Oracle Database 12c Enterprise Edition Release 12.1.0.1.0 - 64bit Production
With the Partitioning, OLAP, Advanced Analytics and Real Application
Testing options
SYS@ORA12CR1>
```

На заметку! Если на этом шаге вы столкнетесь с ошибкой ORA-12505 “TNS:listener does not currently know of SID given in connect Descriptor” (ORA-12505 “Прослушивающему процессу TNS в настоящий момент не известен идентификатор SID, предоставленный в дескрипторе соединения”), значит, прослушивающий процесс базы данных не сконфигурирован со статической регистрационной записью для этого сервера. Администратору базы данных не разрешено удаленно подключаться как SYSDBA, когда экземпляр базы не запущен. Именно так обстоят дела в большинстве установок Oracle9i и последующих версий. Необходимо настроить статическую регистрацию сервера в конфигурационном файле `listener.ora`. Поищите на странице поиска в документации OTN (Oracle Technology Network) для своей версии СУБД подробную информацию о конфигурации этой статической службы, указав в качестве искомой строку “Configuring Static Service Information”.

Мы прошли аутентификацию и, следовательно, вошли в систему. Теперь мы можем успешно выполнять удаленный запуск, останов и администрирование базы данных с применением учетной записи SYSDBA. В системе присутствует еще один пользователь, OPS\$TKYTE, которому были выданы права SYSDBA, но он пока не может удаленно подключаться к базе данных:

```
[tkyte@dellpe ~]$ sqlplus 'ops$tkyte/foobar'@ora12crl as sysdba
SQL*Plus: Release 12.1.0.1.0 Production on Tue Sep 3 14:22:21 2013
Copyright (c) 1982, 2013, Oracle. All rights reserved.

ERROR:
ORA-01017: invalid username/password; logon denied
ОШИБКА:
ORA-01017: недопустимое имя пользователя/пароль; вход запрещен
```

Причина неудачи подключения в том, что пользователь ORS\$TKYTE еще не занесен в файл паролей. Чтобы записать его в файл паролей, необходимо “заново выдать” этой учетной записи привилегию SYSDBA:

```
[tkyte@dellpe ~]$ sqlplus / as sysdba
SQL*Plus: Release 12.1.0.1.0 Production on Tue Sep 3 14:23:11 2013
Copyright (c) 1982, 2013, Oracle. All rights reserved.

Connected to:
Oracle Database 12c Enterprise Edition Release 12.1.0.1.0 - 64bit Production
With the Partitioning, OLAP, Advanced Analytics and Real Application
Testing options
SYS@ORA12CR1> grant sysdba to ops$tkyte;

Grant succeeded.
Выдано успешно.
SYS@ORA12CR1> exit
Disconnected from Oracle Database 12c Enterprise Edition Release
12.1.0.1.0 - 64bit Production
With the Partitioning, OLAP, Advanced Analytics and Real Application
Testing options
```

```
[tkyte@dellpe ~]$ sqlplus 'ops$tkyte/foobar'@ora12cr1 as sysdba
SQL*Plus: Release 12.1.0.1.0 Production on Tue Sep 3 14:23:25 2013
Copyright (c) 1982, 2013, Oracle. All rights reserved.
```

```
Connected to:
Oracle Database 12c Enterprise Edition Release 12.1.0.1.0 - 64bit Production
With the Partitioning, OLAP, Advanced Analytics and Real Application
Testing options
SYS@ORA12CR1>
```

Этот запрос создал запись в файле паролей, и теперь Oracle будет поддерживать пароль в синхронизированном состоянии. Если пользователь OPS\$TKYTE изменит свой пароль, старый пароль перестанет действовать для удаленных подключений SYSDBA, а новый вступит в силу.

Аналогичный процесс повторяется для любого пользователя, который имел привилегию SYSDBA, но не был записан в файл паролей.

Файл отслеживания изменений

Файл отслеживания изменений — это необязательный файл, предназначенный для использования с Oracle 10g Enterprise Edition и последующими версиями. Его единственное назначение заключается в отслеживании блоков, которые были изменены с момента последнего инкрементного резервного копирования. Это позволяет утилите RMAN (Recovery Manager — диспетчер восстановления) выполнять резервное копирование только тех блоков базы данных, которые действительно были изменены, без необходимости в считывании всей базы данных.

В версиях, предшествующих Oracle 10g, при инкрементном резервном копировании для выявления блоков, которые были изменены с момента последнего инкрементного копирования, приходилось считывать весь набор файлов базы данных. Таким образом, если в базу данных объемом 1 Тбайт было добавлено всего 500 Мбайт новых данных (например, производилась загрузка информационного хранилища), то для нахождения этих 500 Мбайт новой информации с целью резервного копирования пришлось бы прочитать 1 Тбайт данных. В итоге процесс инкрементного резервного копирования сохранял значительно меньший объем данных, но ему по-прежнему требовалось считывать всю базу данных.

В Oracle 10g Enterprise Edition и последующих версиях такая проблема больше не возникает. В процессе работы базы данных и по мере того, как блоки модифицируются, Oracle дополнительно поддерживает файл, который указывает утилите RMAN, какие блоки были изменены. Создается этот файл отслеживания изменений довольно просто с помощью команды ALTER DATABASE:

```
SYS@ORA12CR1> alter database enable block change tracking using file
  2  '/home/ora12cr1/oradata/ORA12CR1/changed_blocks.bct';
Database altered.
База данных изменена.
```

Внимание! В этой книге я время от времени повторяю: не забывайте, что команды, которые устанавливают параметры, модифицируют базу данных или вносят существенные изменения, не должны выдаваться необдуманно и, безусловно, их следует протестировать, прежде чем выполнять в “реальной” системе. Предыдущая команда в действительности заставит базу данных делать больший объем работ и приведет к росту потребления ресурсов.

Чтобы отключить и удалить файл отслеживания измененных блоков, нужно снова воспользоваться командой ALTER DATABASE:

```
SYS@ORA12CR1> alter database disable block change tracking;
Database altered.
База данных изменена.
```

Обратите внимание, что эта команда удалит файл отслеживания измененных блоков. Она не просто отключает функцию отслеживания — она также удаляет сам файл.

На заметку! В определенных операционных системах, таких как Windows, вы можете обнаружить, что после запуска приведенного примера — создание файла отслеживания измененных блоков и затем отключение функции отслеживания — файл по-прежнему существует. Это проблема, специфичная для ОС: во многих операционных системах она не проявляется. Она возникает, только если файл отслеживания измененных блоков создается и отключается в одном сеансе. Сеанс, создавший файл отслеживания измененных блоков, оставляет его открытым, а в ряде ОС не разрешено удалять файл, который открыт предыдущим процессом (например, процессом сеанса, создавшего файл). В этом нет ничего страшного — просто позже понадобится удалить файл самостоятельно.

Функция отслеживания измененных блоков может быть включена в режиме ARCHIVELOG или NOARCHIVELOG. Но помните, что база данных, работающая в режи-

ме NOARCHIVELOG, в котором ежедневно генерируемый журнал повторения не сохраняется, не имеет возможности восстановить все изменения в случае отказа носителя (диска либо устройства)! Однажды такая база данных неизбежно утратит данные. Мы рассмотрим эти два режима работы базы данных более подробно в главе 9.

Ретроспективные журналы

Ретроспективные журналы (flashback log) были введены в Oracle 10g для поддержки команды FLASHBACK DATABASE — нового средства этого выпуска базы данных в редакции Enterprise Edition. Ретроспективные журналы содержат “предшествующие образы” модифицированных блоков базы данных, которые можно применять для восстановления базы данных в состоянии на определенный момент времени в прошлом.

Команда FLASHBACK DATABASE

Команда FLASHBACK DATABASE была добавлена для ускорения достаточно медленного процесса восстановления базы данных на определенный момент времени. Ее можно использовать вместо полного восстановления базы данных и наката с применением архивных журналов, и она предназначена главным образом для повышения скорости восстановления после “катастрофы”. Например, давайте посмотрим, что администратор базы данных может сделать для восстановления после случайного удаления схемы не той базы данных (удалять нужно было версию в тестовой среде). Администратор сразу же заметил свою ошибку и немедленно остановил СУБД. Но что делать дальше?

До появления возможности FLASHBACK DATABASE администратору, скорее всего, пришлось бы выполнить следующие действия.

1. Остановить базу данных.
2. Восстановить последнюю полную резервную копию базы данных с магнитной ленты (как правило), что в общем случае представляет собой длительный процесс. Обычно это инициируется в RMAN посредством команды `RESTORE DATABASE UNTIL <момент_времени>`.
3. Восстановить все архивные журналы повторений, сгенерированные с момента резервного копирования.
4. Используя архивные журналы повторения (и возможно информацию из оперативных журналов повторения), произвести накат (последовательное выполнение всех ранее выполненных транзакций) базы данных до момента, непосредственно предшествующего ошибочному применению команды `DROP USER`. Шаги 3 и 4 в этом списке обычно инициируются в RMAN через команду `RECOVER DATABASE UNTIL <момент_времени>`.
5. Открыть базу данных с опцией `RESETLOGS`.

Это был весьма непростой процесс с многочисленными шагами, который в целом занимал длительное время (в течение которого, естественно, никто не мог получать доступ к базе данных). Причин для подобного восстановления базы данных на определенный момент времени множество: некорректно сработавший сценарий модернизации; неудавшийся процесс модернизации; невнимательная выдача команды

кем-то с достаточными полномочиями (пожалуй, самый часто встречающийся случай); создание каким-то процессом проблем с целостностью данных в крупной базе данных (опять-таки, случайное событие; возможно, процесс был запущен два раза вместо одного или содержал в себе ошибку). Какой бы ни была причина, конечным результатом станет длительный простой.

Ниже описаны шаги, которые администратору базы данных понадобится выполнить в Oracle 10g Enterprise Edition и последующих версиях для восстановления базы данных при наличии сконфигурированной возможности ретроспективы.

1. Остановить базу данных.
2. Выполнить начальное монтирование базы данных и выдать команду ретроспективного восстановления, используя номер SCN (внутренние часы Oracle), точку восстановления (указатель на SCN) либо метку времени (показание обычных часов) с точностью до пары секунд.
3. Открыть базу данных с указанием опции RESETLOGS.

Для применения этой функциональности база данных должна работать в режиме ARCHIVELOG и быть сконфигурированной для включения команды FLASHBACK DATABASE. Иначе говоря, возможность ретроспективы должна быть настроена до того, как в ней возникнет потребность. Ее нельзя включить после того, как база данных окажется поврежденной, поэтому решение о ее использовании следует принимать заранее.

Область для быстрого восстановления

Область для быстрого восстановления (Fast Recovery Area) — новая концепция в Oracle 10g и последующих версиях. Начиная с Oracle 10g, впервые за многие годы (свыше 25 лет) концепция резервных копий в Oracle изменилась. В прошлом в основе резервного копирования и восстановления лежала концепция применения устройств с последовательным доступом, таких как накопители на магнитных лентах. Устройства произвольного доступа (диски) всегда считались слишком дорогостоящими, чтобы тратить их на обычные резервные копии. В итоге использовались сравнительно дешевые ленточные устройства большой емкости.

Однако в наши дни можно приобрести терабайты дискового хранилища по очень низкой цене. Например, мой сын Алан стал первым ребенком в квартале, обладающим сетевым устройством хранения (network attached storage — NAS) объемом 1 Тбайт. Оно обошлось мне в \$125. Я вспоминаю первый жесткий диск, установленный на моем персональном компьютере: он имел колоссальный на тот момент объем 40 Мбайт. На самом деле мне пришлось его еще и разбить на два логических диска, поскольку ОС, с которой я имел дело тогда (MS-DOS), была не в состоянии распознавать диски емкостью больше 32 Мбайт. За последние 25 лет или около того ситуация коренным образом изменилась.

Область Fast Recovery Area в Oracle — это местоположение, где СУБД будет управлять множеством файлов, связанных с резервным копированием и восстановлением баз данных. В этой области (*области*, которая является частью набора дисков и предназначена для этой цели — например, каталог) можно найти:

- фрагменты резервных копий RMAN (полных и/или инкрементных резервных копий);

- копии образов RMAN (побайтовые копии файлов данных и управляющих файлов);
- оперативные журналы повторения;
- архивные журналы повторения;
- мультиплексированные управляющие файлы;
- ретроспективные журналы.

Новая область Fast Recovery Area применяется Oracle для управления перечисленными файлами, поэтому серверу известно, что присутствует на диске, а что отсутствует (и находится, возможно, где-то на магнитной ленте). Используя эту информацию, Oracle может выполнять такие операции, как восстановление поврежденных файлов данных с диска на диск или ретроспективный откат (операция “перемотки в начало”) базы данных для отмены операции, которая не должна была выполняться. Например, команду FLASHBACK DATABASE можно было бы применить для возвращения базы данных в состояние, в котором она пребывала пять минут назад (не делая полное восстановление базы данных и ее восстановление на определенный момент времени), что позволяет отменить, к примеру, случайное удаление учетной записи пользователя.

Область Fast Recovery Area — нечто большее, чем логическая концепция. Это область, содержащая типы файлов, которые были рассмотрены в главе. Ее использование не является обязательным, но если вы намерены применять некоторые расширенные возможности вроде FLASHBACK DATABASE, то должны использовать эту область для хранения информации.

Файлы DMP (файлы экспорта/импорта)

Утилиты Export (Экспорт) и Import (Импорт) — уважаемые инструменты извлечения и загрузки данных Oracle, которые входили в состав множества версий. Задача утилиты Export заключается в создании независимого от платформы файла DMP, который содержит все обязательные метаданные (в форме операторов CREATE и ALTER) и дополнительно сами данные для воссоздания таблиц, схем и даже целых баз данных. Единственное назначение утилиты Import состоит в чтении этих файлов DMP, выполнении операторов DDL и загрузке любых найденных данных.

На заметку! В версии Oracle Database 11g Release 2 утилита Export была официально объявлена устаревшей. Она поставляется только для применения с унаследованными структурами баз данных. Новые типы данных, новые структуры, новые функциональные возможности баз данных этот инструмент не поддерживает. Я настоятельно рекомендую пользоваться Data Pump — пришедшим на замену инструментом экспорта/импорта, который несколько лет тому назад появился в Oracle 10g.

Файлы DMP спроектированы так, чтобы обеспечивать обратную совместимость, т.е. более новые версии могут успешно читать и обрабатывать файлы DMP из предшествующих выпусков. Я слышал о людях, которые экспортировали базу данных версии 5 и благополучно импортировали ее в Oracle 10g (просто ради проверки). Таким образом, процесс Import может считывать старые версии файлов DMP и об-

рабатывать содержащиеся в них данные. Однако обратное, скорее всего, окажется невозможным: процесс Import, входящий в состав Oracle9i Release 1, не может, да и не должен, считывать файл DMP, который создан версией Oracle9i Release 2 или Oracle 10g Release 1. Например, я экспортировал простую таблицу из Oracle 11g Release 2. Когда я попробовал воспользоваться этими файлами DMP в Oracle9i Release 1, то быстро обнаружил, что процесс Import этой версии даже не пытается обрабатывать файл DMP версии Oracle 11g Release 2:

```
$ imp userid=/ full=y
```

```
Import: Release 10.2.0.4.0 - Production on Wed Jan 20 18:21:03 2010
```

```
Copyright (c) 1982, 2007, Oracle. All rights reserved.
```

```
Connected to: Oracle Database 10g Enterprise Edition Release 10.2.0.4.0 -  
Production
```

```
With the Partitioning, OLAP, Data Mining and Real Application Testing options
```

```
IMP-00010: not a valid export file, header failed verification
```

```
IMP-00000: Import terminated unsuccessfully
```

```
IMP-00010: недопустимый экспортный файл, заголовок не прошел верификацию
```

```
IMP-00000: процесс Import завершен неудачно.
```

Даже в случае, если бы инструмент Import мог распознать файл, высоки шансы на то, что DDL-код, сгенерированный инструментом Export из версии Oracle 11g Release 2, не сможет быть распознан в более ранних версиях Oracle. Например, предположим, что вы экспортируете из любой выпуска Oracle версии 9i Release 2 или выше. В экспортном файле вы обнаружите, что к каждому оператору CREATE TABLE добавлена опция COMPRESS или NOCOMPRESS. Средство базового сжатия таблиц появилось в версии Oracle9i Release 2. Если попробовать прочитать этот экспортный файл в любой версии базы, предшествовавшей 9i Release 2, то выяснится, что содержащийся в нем DDL-код не выполняется, причем в 100% случаев. Ни один оператор CREATE TABLE не будет работать, поскольку ключевые слова NOCOMPRESS/COMPRESS в старых выпусках не распознаются.

Файлы DMP являются независимыми от платформы, поэтому их можно безопасно экспортировать из любой платформы, перенести на другую платформу и затем импортировать (при условии, что установленная версия Oracle разрешает это). Однако при переносе файлов с помощью FTP в среде Windows возникает проблема, обусловленная тем, что по умолчанию Windows посчитает файл DMP текстовым и будет стремиться преобразовать символы перевода строки (маркеры конца строки в UNIX/Linux) в пары символов возврата каретки/перевода строки, полностью разрушая файл DMP. При передаче файла DMP по FTP в среде Windows необходимо выбирать двоичный режим. Если окажется, что процесс импорта не работает, необходимо удостовериться, что размеры исходного и целевого файлов совпадают. Я даже не могу вспомнить, сколько раз эта проблема стопорила процесс переноса файлов.

Файлы DMP являются двоичными, т.е. их нельзя редактировать для внесения изменений. Из них можно извлечь большой объем информации (DDL-команды CREATE и многое другое), но их нельзя редактировать в текстовом (или в любом другом) редакторе. В своей книге *Expert One-on-One Oracle* я посвятил немало страниц описанию утилит Import и Export и файлов DMP. Популярность этих утилит снижается в связи с появлением предлагающих большие возможности утилит Data Pump.

Файлы Data Pump

Data Pump (помпа данных) — это файловый формат, который используется, по меньшей мере, двумя инструментами в Oracle. Внешние таблицы могут загружать и выгружать данные в этом формате, а новые программные средства импорта/экспорта IMPDP и EXPDP используют его во многом подобно тому, как утилиты Import и Export работают с файловым форматом DMP.

На заметку! Формат Data Pump ориентирован на применение исключительно в Oracle 10g Release 1 и последующих версиях — в выпуске Oracle9i он не существовал и не может там использоваться.

Значительная часть проблем, связанных с описанными ранее файлами DMP, со временем станут актуальными и для файлов Data Pump. Эти файлы представляют собой межплатформенные (переносимые) двоичные файлы, которые содержат метаданные (хранящиеся не в операторах CREATE/ALTER, а в XML-коде) и возможно данные. Для нас, как конечных пользователей инструментов, важно то, что в качестве структуры представления метаданных применяется XML. Утилиты IMPDP и EXPDP обладают рядом сложных возможностей фильтрации и преобразования, которых никогда ранее не было в старых утилитах Import/Export. Частично это обусловлено использованием языка XML и тем, что оператор создания таблицы хранится не в форме CREATE TABLE, а в виде документа разметки. Это позволяет легко реализовать запрос типа “заменить все ссылки на табличное пространство FOO ссылками на табличное пространство BAR”. Когда метаданные хранились в файле DMP в виде операторов CREATE/ALTER, для решения этой задачи утилите Import приходилось подробно анализировать каждый SQL-оператор перед его выполнением (т.е. заниматься несвойственной ей задачей). Тем не менее, утилите IMPDP для этого достаточно применить простую XML-трансформацию — если строка FOO относится к табличному пространству, то она должна быть заключена в дескрипторы <TABLESPACE>FOO</TABLESPACE> (или представлена как-то иначе).

Использование XML позволило утилитам EXPDP и IMPDP оставить старые утилиты Import и Export далеко позади с точки зрения предлагаемых возможностей. В главе 15 мы подробно проанализируем основные особенности этих инструментов. А пока давайте посмотрим, как применять формат Data Pump для быстрого извлечения определенных данных из базы данных A и их перемещения в базу данных B. Здесь мы будем использовать “внешнюю таблицу в обратном направлении”.

Внешние таблицы, впервые появившиеся в Oracle9i Release 1, предоставили возможность считывать плоские файлы (простые текстовые файлы) так, как если бы они были таблицами базы данных. Это позволяет полностью задействовать для их обработки возможности, предлагаемые языком SQL. Они допускали только чтение и были предназначены для загрузки данных из внешнего мира в Oracle. Внешние таблицы в Oracle 10g Release 1 и последующих версиях поддерживают и другое направление: их можно применять для извлечения данных из базы данных в формате Data Pump с целью облегчения переноса данных на другой компьютер или на другую платформу. Перед началом упражнения понадобится создать объект DIRECTORY, который указывает Oracle место выгрузки данных:

```
EODA@ORA12CR1> create or replace directory tmp as '/tmp';
Directory created.
Каталог создан.
```

```
EODA@ORA12CR1> create table all_objects_unload
2 organization external
3 ( type oracle_datapump
4   default directory TMP
5   location( 'allobjects.dat' )
6 )
7 as
8 select * from all_objects
9 /
Table created.
```

И это буквально все, что необходимо: в каталоге /tmp мы имеем файл по имени allobjects.dat, в котором находится содержимое запроса select * from all_objects. Теперь можно заглянуть в него:

```
EODA@ORA12CR1> !strings /tmp/allobjects.dat | head
"EODA"."U"
x86_64/Linux 2.4.xx
AL32UTF8
12.00.00.00.00
001:001:000001:000001
i<?xml version="1.0"?>
<ROWSET>
<ROW>
  <STRMTABLE_T>
```

Это только заголовок, или верхняя часть, файла. Теперь с помощью двоичного режима передачи FTP (аналогично файлам DMP) этот файл можно переместить на любую другую платформу, где установлена версия Oracle 12c, после чего выдать оператор CREATE DIRECTORY (чтобы сообщить базе данных, где расположен файл) и оператор CREATE TABLE следующего вида:

```
create table t
( OWNER          VARCHAR2(30),
  OBJECT_NAME    VARCHAR2(30),
  SUBOBJECT_NAME VARCHAR2(30),
  OBJECT_ID      NUMBER,
  DATA_OBJECT_ID NUMBER,
  OBJECT_TYPE    VARCHAR2(19),
  CREATED        DATE,
  LAST_DDL_TIME  DATE,
  TIMESTAMP      VARCHAR2(19),
  STATUS         VARCHAR2(7),
  TEMPORARY      VARCHAR2(1),
  GENERATED     VARCHAR2(1),
  SECONDARY      VARCHAR2(1)
)
organization external
( type oracle_datapump
  default directory TMP
  location( 'allobjects.dat' )
);
```

Теперь все готово для чтения выгруженных данных, используя SQL-код непосредственно. В этом и заключается основное преимущество файлового формата Data Pump: он позволяет немедленно перемещать данные из одной системы в другую посредством сети “на своих двоих” по мере необходимости. Вспомните об этом, когда в следующий раз потребуется захватить подмножество данных домой на выходные дни для проведения тестирования.

Даже если кодировка баз данных отличается (в рассмотренном примере это было не так), теперь благодаря файловому формату Data Pump имеется возможность распознавания разных символьных наборов, а также их обработки. Преобразование символьных наборов может выполняться на лету, когда это нужно, обеспечивая “корректность” представления данных для каждой базы данных.

Мы вернемся к рассмотрению файлового формата Data Pump в главе 15, но настоящий раздел должен был содействовать в общем понимании содержимого файлов этого формата.

Плоские файлы

Плоские файлы существуют с момента зарождения электронной обработки данных. Мы встречаемся с ними буквально на каждом шагу. Например, описанный ранее сигнальный журнал является плоским файлом.

Я считаю довольно точным следующее определение плоского файла.

Электронная запись, лишенная всех специфичных для приложения (программы) форматов. Это позволяет перемещать элементы данных в другие приложения с целью манипулирования ими. Такой режим отбрасывания форматов электронных данных предотвращает потерю данных из-за устаревания оборудования и патентованного программного обеспечения.

Плоский файл представляет собой просто файл, в котором каждая “строка” является “записью” и содержит определенный текст, как правило, разделенный символами запятой или вертикальной черты. Плоские файлы легко читаются в Oracle либо с помощью унаследованной утилиты загрузки данных SQL*Loader, либо с применением внешних таблиц. Все эти вопросы будут подробно обсуждаться в главе 15 (внешние таблицы также рассматриваются в главе 10).

Тем не менее, создавать плоские файлы в Oracle не настолько легко. По ряду причин не существует какого-то простого инструмента командной строки для экспорта информации в плоский файл. Инструменты вроде APEX, SQL Developer и Enterprise Manager упрощают этот процесс, но официальные средства командной строки, которые можно было бы легко использовать в сценариях или как-то по-другому для выполнения этой операции, отсутствуют.

Одна из причин, по которой я решил упомянуть плоские файлы в этой главе, связана со стремлением предложить набор инструментов, способных создавать простые плоские файлы. За долгие годы практики я разработал три метода решения этой задачи, каждый из которых соответствует своей ситуации. Для выполнения работы первая утилита применяет комбинацию PL/SQL и пакета UTL_FILE с динамическим SQL. При небольших объемах данных (сотни или тысячи строк) она

обеспечивает достаточную гибкость и скорость. Однако эта утилита должна создавать свои файлы на компьютере сервера базы данных, что временами нежелательно. Для таких ситуаций я написал утилиту SQL*Plus, которая создает плоские файлы на компьютере, где функционирует SQL*Plus. Поскольку SQL*Plus может подключаться к серверу Oracle в любой точке сети, это предоставляет возможность выгружать в плоский файл любые данные из любой базы данных в сети. Наконец, в случаях, когда требуется максимальная скорость, ничего другого кроме написания программы на языке С не остается (по моему мнению). Таким образом, я разработал на Pro*C утилиту командной строки для выгрузки данных в плоские файлы. Все три утилиты свободно доступны по адресу <http://tkyte.blogspot.com/2009/10/httpasktomoraclecomtkyteflat.html>, и там же будут появляться любые новые инструменты, созданные для выгрузки в плоские файлы.

Резюме

В этой главе мы исследовали важные типы файлов, используемые базой данных Oracle — от небольших файлов параметров (без которых даже не удастся приступить к работе) до крайне важных журналов повторения и файлов данных. Мы рассмотрели структуры хранения данных Oracle, начиная с табличных пространств и сегментов и заканчивая экстенентами и блоками базы данных, представляющими собой наименьшую единицу хранения. Мы кратко объяснили работу контрольных точек и даже забежали немного наперед, указав действия, которые выполняют физические процессы или потоки Oracle. Мы также раскрыли многие дополнительные файловые типы, такие как файлы паролей, файлы отслеживания изменений, файлы Data Pump и т.д. Теперь вы готовы приступить к исследованию структур памяти Oracle в следующей главе.

ГЛАВА 4

Структуры памяти

В этой главе мы рассмотрим три основные структуры памяти Oracle.

- **Глобальная область системы (System Global Area — SGA).** Это большой совместно используемый сегмент памяти, к которому в тот или иной момент времени будут обращаться практически все процессы Oracle.
- **Глобальная область процесса (или программы) (Process (Prorgam) Global Area — PGA).** Это закрытая область отдельного процесса или потока; из других процессов или потоков она не доступна.
- **Глобальная область пользователя (User Global Area — UGA).** Эта область памяти связана с конкретным сеансом. Она располагается либо в области SGA (если подключение к базе данных выполнено посредством разделяемого сервера), либо в области PGA (если подключение к базе данных осуществлено через выделенный сервер).

На заметку! В ранних версиях Oracle разделяемый сервер назывался *многопоточным сервером (Multi-Threaded Server — MTS)*. В этой книге мы будем всегда использовать термин “разделяемый сервер”.

При обсуждении управления памятью в Oracle нам предстоит исследовать пять режимов.

- Автоматическое управление памятью (automatic memory management — AMM), применяемое к областям SGA и PGA, доступно только в Oracle 11g и последующих версиях. Администратор базы данных устанавливает всего лишь один параметр `MEMORY_TARGET`, чтобы позволить базе данных определить, каким образом назначать размеры всем областям памяти.
- Автоматическое управление разделяемой памятью (automatic shared memory management — ASMM), используемое в отношении к SGA. Администратор базы данных устанавливает целевой размер для области SGA (через `SGA_TARGET`).
- Ручное управление разделяемой памятью, применяемое к SGA. Администратор базы данных вручную устанавливает размеры индивидуальных областей памяти SGA (посредством `DB_CACHE_SIZE`, `SHARED_POOL_SIZE` и т.д.).
- Автоматическое управление памятью PGA, используемое в отношении к PGA. Администратор базы данных устанавливает целевой размер для области PGA (через `PGAAggregateTarget`).

- Ручное управление памятью PGA, применяемое к PGA. Администратор базы данных вручную устанавливает размеры индивидуальных областей памяти PGA (посредством `SORT_AREA_SIZE`, `HASH_AREA_SIZE` и т.д.). В Oracle настоятельно рекомендуют не пользоваться этим методом, но мы обсудим его, чтобы предоставить фундамент для понимания других концепций управления памятью.

Какой метод управления памятью из перечисленных выше должен применяться? В идеальной среде все использовали бы автоматическое управление памятью, не так ли? Достаточно установить один параметр (`MEMORY_TARGET`) — и все готово. Однако реальный мир оказывается не настолько черно-белым. Временами существуют такие аспекты среды, когда вы лучше знаете, как должны выглядеть целевые размеры областей памяти, поэтому применяете один из менее автоматизированных методов управления памятью. Цель настоящей главы заключается в том, чтобы помочь вам освоить все аспекты памяти Oracle и принимать обоснованные решения относительно того, каким образом включать управление памятью.

На заметку! Некоторые конфигурации операционных систем несовместимы с автоматическим управлением памятью (например, Linux HugePages). За подробными сведениями по этому поводу обращайтесь к справочному руководству администратора баз данных Oracle для операционных систем, основанных на Linux и UNIX.

Мы займемся исследованием всех упомянутых методов после того, как обсудим управление памятью PGA и UGA, сначала ручное, а затем автоматическое. После этого мы перейдем к области SGA, снова рассмотрев методы ручного и автоматического управления памятью. Закончим мы тем, что взглянем на то, каким образом управлять памятью SGA и PGA, используя единственный параметр.

Глобальная область процесса и глобальная область пользователя

Область PGA — это фрагмент памяти, специфичный для процесса. Другими словами, это память, связанная с отдельным процессом или потоком операционной системы. Такая память не является доступной любому другому процессу или потоку в системе. Обычно она выделяется посредством вызовов функции `C` под названием `malloc()` или `memmap()`, а ее объем увеличивается (и уменьшается) во время выполнения. Область PGA никогда не выделяется в области SGA базы данных Oracle; она всегда распределяется процессом или потоком только локально — *P* в PGA означает *process* (процесс) или *program* (программа), и это не разделяемая область.

Область UGA, по сути, представляет состояние процесса. Сеанс должен всегда иметь доступ к этой памяти. Расположение UGA зависит от способа подключения к базе данных Oracle. Если подключение было выполнено через разделяемый сервер, область UGA должна храниться в структуре памяти, к которой имеет доступ каждый процесс разделяемого сервера — областью подобного рода будет SGA. Таким образом, сеансом может пользоваться любой из разделяемых серверов, поскольку любой из них может считывать и записывать данные сеанса. С другой стороны, в случае применения подключения посредством выделенного сервера потребность в универ-

сальном доступе к состоянию сеанса отпадает, и область UGA становится практически синонимом PGA; фактически она будет содержаться внутри области PGA выделенного сервера. Если вы просмотрите статистическую информацию о системе, то заметите, что в режиме выделенного сервера область UGA хранится в области PGA выделенного сервера (размер PGA будет больше или равен размеру области UGA; размер PGA будет включать также и размер UGA).

Итак, область PGA содержит память процесса и может включать область UGA. Остальные области памяти PGA в основном служат для выполнения сортировки внутри памяти, слияния битовых индексов и хеширования. Можно без опаски утверждать, что наряду с UGA эти процессы вносят наибольший вклад в формирование области PGA.

Начиная с версии Oracle9i Release 1, существуют два способа управления этой отличной от UGA памятью в области PGA.

1. *Ручное управление памятью PGA*, при котором вы указываете Oracle, какой объем памяти можно задействовать для выполнения сортировки и хеширования в любой момент, когда это необходимо в конкретном процессе.
2. *Автоматическое управление памятью PGA*, при котором вы указываете базе данных Oracle, какой объем памяти она должна пытаться использовать на уровне всей системы.

Начиная с версии Oracle 11g Release 1, автоматическое управление памятью PGA может быть реализовано с применением одного из двух приемов.

1. Установка параметра инициализации PGA_AGGREGATE_TARGET и сообщение Oracle объема памяти PGA, который следует попытаться использовать на уровне экземпляра.
2. Установка параметра инициализации MEMORY_TARGET и уведомление Oracle об объеме общей памяти, которую экземпляр базы данных должен применять для SGA и PGA; размер PGA будет выведен из этого параметра самим экземпляром.

Способы выделения и использования памяти в каждом из этих случаев заметно отличаются, так что давайте обсудим их по очереди.

На заметку! Следует отметить, что в среде Oracle9i при использовании подключения посредством разделяемого сервера можно применять только ручное управление памятью PGA. Такое ограничение было устранено в Oracle 10g Release 1 (и последующих версиях). Начиная с этого выпуска, с подключениями через разделяемый сервер можно использовать либо автоматическое, либо ручное управление памятью PGA. Если вы работаете с Oracle 10g Release 1 (и последующими версиями), то должны применять автоматическое управление памятью PGA.

Режимы управления памятью PGA контролируются параметром инициализации базы данных WORKAREA_SIZE_POLICY и могут изменяться на уровне сеанса. В Oracle9i Release 2 и последующих версиях этот параметр инициализации по умолчанию получает значение AUTO, что устанавливает автоматическое управление памятью PGA, когда это возможно. В Oracle9i Release 1 его стандартным значением было MANUAL (вручную).

В последующих разделах мы поочередно рассмотрим каждый из этих подходов.

Ручное управление памятью PGA

При ручном управлении памятью PGA наибольшее влияние на размер области PGA без учета объема памяти, выделенного сеансом для таблиц PL/SQL и других переменных, будут оказывать следующие параметры.

- **SORT_AREA_SIZE.** Общий объем памяти, который будет использоваться для сортировки информации перед ее сбросом на диск (с применением дискового пространства во временном табличном пространстве, назначенном пользователю).
- **SORT_AREA_RETAINED_SIZE.** Объем памяти, который будет использоваться для хранения отсортированных данных после завершения сортировки. То есть, если значением параметра **SORT_AREA_SIZE** является 512 Кбайт, а параметра **SORT_AREA_RETAINED_SIZE** — 256 Кбайт, то процесс сервера задействует до 512 Кбайт памяти для сортировки данных во время начальной обработки запроса. По завершении сортировки область сортировки будет “сокращена” до 256 Кбайт, а любые данные сортировки, которые не уместились в этот объем, будут записаны во временное табличное пространство.
- **HASH_AREA_SIZE.** Объем памяти, который процесс сервера может применять для хранения хеш-таблиц. Эти структуры используются во время хеш-соединений, обычно при соединении крупных наборов данных друг с другом. Меньший из двух наборов будет хешироваться в памяти, а все данные, не уместившиеся в хеш-область памяти, будут сохранены во временных табличных пространствах по ключу соединения.

Описанные параметры управляют объемом памяти, который Oracle будет использовать для выполнения сортировки или хеширования данных в памяти перед применением временного табличного пространства на диске, и размером той части этого сегмента памяти, которая будет сохранена после завершения сортировки. В общем случае память с объемом, вычисленным как **SORT_AREA_SIZE-SORT_AREA_RETAINED_SIZE**, выделится из области PGA, а память с объемом, заданным параметром **SORT_AREA_RETAINED_SIZE**, будет находиться в области UGA. Сведения о текущем использовании памяти PGA и UGA можно получать и отслеживать размеры этих областей, запрашивая специальные *V\$-представления Oracle*, которые также называются *динамическими представлениями производительности*.

Например, давайте проведем небольшой тест, при котором один сеанс будет сортировать большой объем данных, а из второго сеанса мы будем отслеживать информацию об использовании памяти UGA/PGA в первом сеансе. Чтобы решить эту задачу предсказуемым образом, мы создадим копию таблицы **ALL_OBJECTS**, содержащей около 72 000 строк безо всяких индексов (таким образом, мы знаем, что в случае применения конструкции **ORDER BY** в этой таблице будет происходить сортировка):

```
EODA@ORA12CR1> create table t as select * from all_objects;
Table created.
Таблица создана.

EODA@ORA12CR1> exec dbms_stats.gather_table_stats( user, 'T' );
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.
```

Чтобы предотвратить любые побочные эффекты начального полного разбора запросов, мы запустим приведенный ниже сценарий, но пока проигнорируем его вывод. Мы выполним его еще раз в новом сеансе, чтобы посмотреть на результаты использования памяти в управляемой среде. Мы будем применять размеры области сортировки 64 Кбайт, 1 Мбайт и 1 Гбайт по очереди, поэтому сохраним этот сценарий как `run_query.sql` (в сценарии предполагается, что именем пользователя и паролем являются `eoda/foo`, так что измените их в соответствии с имеющейся средой):

```
connect eoda/foo
set serveroutput off
set echo on
column sid new_val SID
select sid from v$mystat where rownum = 1;
alter session set workarea_size_policy=manual;
alter session set sort_area_size = &1;
prompt run @reset_stat &SID and @watch_stat in another session here!
pause
set termout off
select * from t order by 1, 2, 3, 4;
set termout on
prompt run @watch_stat in another session here!
pause
```

А теперь запустим готовый сценарий:

```
@run_query 65536
@run_query 1048576
@run_query 1073741820
```

Пока что не обращайте внимания на вывод; мы просто “разогреваем” разделяемый пул и делаем все “однородным”.

На заметку! При обработке SQL-запроса в базе данных сначала должен быть произведен его разбор. Различают разборы двух типов. Первый — это *полный разбор (hard parse)*, выполняемый при первом разборе запроса экземпляром базы данных и включающий в себя генерацию плана запроса и оптимизацию. Второй тип — *частичный разбор (soft parse)*, в ходе которого могут быть пропущены многие действия, обязательные во время полного разбора, поскольку при этом могут быть задействованы результаты работы, сделанной в течение полного разбора. Мы осуществляем полный разбор приведенных запросов, чтобы в следующем разделе не пришлось оценивать объем работы, которая выполнена данной операцией.

Прежде чем продолжить, я рекомендую выйти из сеанса SQL*Plus и снова войти, чтобы получить согласованную среду, или среду, в которой никакая работа пока еще не делалась. Теперь нам нужна возможность измерения показателей памяти сеанса, выполняющего крупные запросы ORDER BY из второго отдельного сеанса. Если мы будем использовать тот же самый сеанс, то запрос, предназначенный для просмотра объема памяти, задействованной под сортировку, сам может повлиять на интересующие нас показатели. Для замера памяти из второго сеанса применяется небольшой разработанный мною сценарий SQL*Plus. В действительности это пара сценариев:

вы узнаете, когда запускать их, из `run_query.sql`. Сценарий, выполняющий сброс небольшой таблицы и установку переменной SQL*Plus в идентификатор SID, который требуется отслеживать, называется `reset_stat.sql`:

```
drop table sess_stats;

create table sess_stats
( name varchar2(64), value number, diff number );

variable sid number
exec :sid := &1
```

На заметку! Перед использованием этого сценария (или любого сценария, если уж на то пошло) удостоверьтесь в наличии понимания того, что он делает. Сценарий удаляет и воссоздает таблицу по имени `SESS_STATS`. Если в вашей схеме уже имеется такая таблица, вам придется выбрать для нее другое имя!

Второй сценарий называется `watch_stat.sql`, и в этом учебном примере в нем используется SQL-оператор `MERGE`, так что можно изначально вставить статистические значения для сеанса, а потом возвратиться и обновить их — без необходимости в отдельном сценарии `INSERT/UPDATE`:

```
merge into sess_stats
using
(
select a.name, b.value
  from v$statname a, v$sesstat b
 where a.statistic# = b.statistic#
    and b.sid = :sid
    and (a.name like '%ga %'
        or a.name like '%direct temp%')
) curr_stats
on (sess_stats.name = curr_stats.name)
when matched then
  update set diff = curr_stats.value - sess_stats.value,
           value = curr_stats.value
when not matched then
  insert ( name, value, diff )
  values
    ( curr_stats.name, curr_stats.value, null )
/

select name,
       case when name like '%ga %'
         then round(value/1024,0)
         else value
       end kbytes_writes,
       case when name like '%ga %'
         then round(diff /1024,0)
         else value
       end diff_kbytes_writes
  from sess_stats
 order by name;
```

Я акцентирую внимание на формулировке “в этом учебном примере”, т.к. строки, выделенные полужирным, являются именами интересующих нас статистических показателей, которые изменяются от примера к примеру. В этом конкретном случае нас интересует все, содержащее в имени ga (pga и uga), либо все, включающее direct temp, что в Oracle 10g и последующих версиях покажет прямые чтения и записи во временное пространство (количество операций ввода-вывода при чтении и записи).

На заметку! В Oracle9i прямой ввод-вывод во временное табличное пространство подобным образом не помечался, поэтому должна применяться конструкция WHERE, включающая в себя and (a.name like '%ga %' or a.name like '%physical % direct%').

Когда этот сценарий watch_stat.sql выполняется в командной строке SQL*Plus, мы видим список статистических показателей, касающихся памяти PGA и UGA для данного сеанса, а также ввода-вывода во временную область. Если теперь запустить в сеансе сценарий @run_query 65536, появится вывод, подобный следующему:

```
EODA@ORA12CR1> @run_query 65536
Connected.
Подключено.
EODA@ORA12CR1> set serveroutput off
EODA@ORA12CR1> set echo on
EODA@ORA12CR1> column sid new_val SID
EODA@ORA12CR1> select sid from v$mystat where rownum = 1;

      SID
-----
      23

EODA@ORA12CR1> alter session set workarea_size_policy=manual;
Session altered.
Сеанс изменен.
EODA@ORA12CR1> alter session set sort_area_size = &1;
old  1: alter session set sort_area_size = &1
new  1: alter session set sort_area_size = 65536
Session altered.
Сеанс изменен.
EODA@ORA12CR1> prompt run @reset_stat &SID and @watch_stat in another
session here!
run @reset_stat      23 and @watch_stat in another session here!
EODA@ORA12CR1> pause
```

Здесь мы можем видеть идентификатор SID этого нового сеанса (23). Вдобавок мы устанавливаем ручной режим управления памятью PGA и значение SORT_AREA_SIZE в 65536 (64 Кбайт). Далее сценарий предлагает запустить два других сценария в другом сеансе, что мы и сделаем:

```
EODA@ORA12CR1> @reset_stat 23
Table dropped.
Таблица удалена.
```


Table created.

Таблица создана.

PL/SQL procedure successfully completed.

Процедура PL/SQL успешно завершена.

EODA@ORA12CR1> @watch_stat

6 rows merged.

6 строк объединено.

| NAME | KBYTES_WRITES | DIFF_KBYTES_WRITES |
|---|---------------|--------------------|
| physical reads direct temporary tablespace | 0 | 0 |
| physical writes direct temporary tablespace | 0 | 0 |
| session pga memory | 876 | 0 |
| session pga memory max | 876 | 0 |
| session uga memory | 334 | 0 |
| session uga memory max | 334 | 0 |

6 rows selected.

6 строк выбрано.

На заметку! Сценарий `watch_stat.sql` должен запускаться в том же сеансе, что и `reset_stat.sql`; сценарий `reset_stat.sql` устанавливает переменную привязки `:sid`, необходимую для оператора `MERGE`. В качестве альтернативы можно выполнить `exec :sid := <sid>`, где `<sid>` — идентификатор SID, полученный при тестировании.

Таким образом, как видите, до начала работы мы располагаем около 334 Кбайт данных в области UGA и 876 Кбайт данных в области PGA. Первый возникающий вопрос: сколько памяти мы задействуем из областей PGA и UGA? Другими словами, мы используем 334 + 876 Кбайт памяти или какой-то другой объем? Это сложный вопрос, на который невозможно ответить, не зная, был ли наблюдаемый сеанс с SID 23 подключен к базе данных посредством выделенного или разделяемого сервера. Но, даже располагая этой информацией, ответить на поставленный вопрос может быть нелегко. В режиме выделенного сервера область UGA полностью содержится внутри области PGA, и в этом случае наш процесс или поток потреблял бы 876 Кбайт памяти.

В режиме разделяемого сервера область UGA выделяется из области SGA, а область PGA распределена в памяти операционной системы, которая является закрытой для процесса разделяемого сервера. Итак, в режиме разделяемого сервера в момент получения последней строки приведенного запроса процесс разделяемого сервера может эксплуатироваться кем-то еще. Эта область PGA больше не является «нашей», поэтому формально мы используем 334 Кбайт памяти (за исключением момента действительного запуска запроса, когда мы применяем 1210 Кбайт памяти из объединенных областей PGA и UGA). В этом случае я использую выделенный сервер (иначе было бы невозможно точно выполнить тест) и конечный объем 876 Кбайт памяти из объединенных областей PGA и UGA. Давайте теперь запустим первый крупный запрос в сеансе с идентификатором 23, в котором применяется ручное управление памятью PGA в режиме выделенного сервера. Мы просто должны возвратиться в сеанс, где был запущен сценарий `run_query.sql` и нажать <Enter>, чтобы начать выполнение запроса.

На заметку! Поскольку мы не устанавливали параметр `SORT_AREA_RETAINED_SIZE`, его отображаемое значение будет равно нулю, но действительно используемое значение будет совпадать со значением `SORT_AREA_SIZE`.

```
EODA@ORA12CR1> set termout off
EODA@ORA12CR1> prompt run @watch_stat in another session here!
run @watch_stat in another session here!
EODA@ORA12CR1> pause
```

Точка, где вы видите `set termout off` — это место, где запускается крупный запрос; здесь мы указываем SQL*Plus о необходимости выполнить запрос, ничего не выводя на экран (вывод свыше 70 000 строк займет немало времени). Если мы снова запустим сценарий `watch_stat.sql` во втором сеансе, то увидим примерно такой вывод:

```
EODA@ORA12CR1> @watch_stat
6 rows merged.
6 строк объединено.
```

| NAME | KBYTES_WRITES | DIFF_KBYTES_WRITES |
|---|---------------|--------------------|
| physical reads direct temporary tablespace | 3000 | 3000 |
| physical writes direct temporary tablespace | 3000 | 3000 |
| session pga memory | 1196 | 320 |
| session pga memory max | 1260 | 384 |
| session uga memory | 654 | 320 |
| session uga memory max | 718 | 384 |
| 6 rows selected. | | |

6 строк выбрано.

Обратите внимание, что на этот раз значения `session xxx memory` и `session xxx memory max` не совпадают. Значение `session xxx memory` представляет объем памяти, применяемый в текущий момент времени.

Значение `session xxx memory max` отражает максимальный объем памяти, использованный в ходе сеанса при обработке запроса.

На заметку! В этих примерах не стоит ожидать, что вы получите в точности такие же цифры, как приведенные выше. На объем задействованной памяти влияют многие аспекты, такие как версия Oracle, операционная система, ее функциональные средства и возможности, объем данных, помещенных в таблицу T, и т.д. Вам следует ожидать вариаций в объемах памяти, однако общая картина должна быть аналогичной.

Как видите, объем используемой памяти возрос — мы выполнили определенную сортировку данных. Объем памяти UGA увеличился с 334 Кбайт до 718 Кбайт (максимальная величина) во время обработки запроса. Чтобы выполнить запрос и сортировку, СУБД Oracle выделила область сортировки для нашего сеанса. Кроме того, объем памяти PGA увеличился с 876 Кбайт до 1196 Кбайт. Вдобавок можно заметить, что было сделано 3000 операций записи и чтения в и из временного пространства (т.к. сортируемые данные не уместились в 64 Кбайт — значение параметра `SORT_AREA_SIZE`).

К моменту завершения запроса и исчерпания результирующего набора объем PGA также несколько уменьшился (следует отметить, что в Oracle8i и предшествующих версиях объем PGA вообще не уменьшился бы; эта особенность появилась в Oracle9i и сохранилась в последующих версиях).

Давайте повторим эту же операцию, но изменим размер SORT_AREA_SIZE, увеличив его до 1 Мбайт. Мы выйдем из наблюдаемого сеанса, снова войдем в него и выполним указания по увеличению SORT_AREA_SIZE до 1 Мбайт. Помните, что в другом сеансе, из которого осуществляется мониторинг, понадобится снова запустить сценарий reset_stat.sql. Поскольку начальные значения остаются постоянными (вывод после первого запуска watch_stat.sql должен быть таким же в новом сеансе), я их не повторяю, а привожу только конечный результат:

| NAME | KBYTES_WRITES | DIFF_KBYTES_WRITES |
|---|---------------|--------------------|
| physical reads direct temporary tablespace | 1043 | 1043 |
| physical writes direct temporary tablespace | 1043 | 1043 |
| session pga memory | 1196 | 320 |
| session pga memory max | 2732 | 1856 |
| session uga memory | 718 | 384 |
| session uga memory max | 1756 | 1422 |
| 6 rows selected. | | |
| 6 строк выбрано. | | |

Как видите, на этот раз во время обработки запроса объем памяти UGA значительно возрос. Он временно увеличился почти до 1700 Кбайт (немного больше 1 Мбайт — значения параметра SORT_AREA_SIZE), но при этом количество физических операций ввода-вывода, которые пришлось выполнить для сортировки данных, существенно уменьшилось (в случае работы с памятью большего объема обмен данными с диском происходит реже). Мы можем также избежать многопроходной сортировки — ситуации, когда количество небольших наборов отсортированных данных, которые требуется объединить вместе, настолько велико, что СУБД Oracle вынуждена записывать данные во временное пространство более одного раза. Теперь давайте перейдем к пределу и используем для SORT_AREA_SIZE значение 1 Гбайт:

| NAME | KBYTES_WRITES | DIFF_KBYTES_WRITES |
|---|---------------|--------------------|
| physical reads direct temporary tablespace | 0 | 0 |
| physical writes direct temporary tablespace | 0 | 0 |
| session pga memory | 1132 | 256 |
| session pga memory max | 11372 | 10496 |
| session uga memory | 654 | 320 |
| session uga memory max | 10631 | 10296 |
| 6 rows selected. | | |
| 6 строк выбрано. | | |

Легко заметить, что хотя мы разрешили SORT_AREA_SIZE принимать величины вплоть до 1 Гбайт памяти, в действительности было использовано только около 10 Мбайт. Это говорит о том, что параметр SORT_AREA_SIZE определяет верхнюю границу, а не стандартное значение, и представляет собой только размер выделенной памяти. Кроме того, мы снова выполнили только одну сортировку, но на этот раз она происходила полностью в памяти; никакое временное пространство на диске не было задействовано, о чем свидетельствует отсутствие операций физического ввода-вывода.

При запуске того же самого теста в средах с разными версиями Oracle или даже с другими операционными системами вы можете столкнуться с отличающимся поведением, а значения во всех случаях будут несколько разниться от приведенных. Но общее поведение базы данных должно быть таким же. Другими словами, при увеличении размера области сортировки и выполнении крупных сортировок объемом памяти, используемой сеансом, будет расти. Может оказаться, что объем памяти PGA увеличивается и уменьшается или же остается постоянным во времени, как только что было показано. Например, я уверен, что если вы выполните предыдущий тест в среде Oracle8i, то заметите, что объем памяти PGA не уменьшился до первоначального (т.е. во всех случаях значение `SESSION_PGA_MEMORY` равно значению `SESSION_PGA_MEMORY_MAX`). Этого и следовало ожидать, поскольку в выпусках Oracle8i память PGA управляется как куча и формируется вызовами `malloc()`. В Oracle9i и последующих выпусках новые методы присоединяют и освобождают рабочие области по мере необходимости с применением функций выделения памяти, специфичных для операционной системы.

Ниже перечислены важные моменты, о которых необходимо помнить во время использования параметров `*_AREA_SIZE`.

- Эти параметры управляют максимальным объемом памяти, потребляемой операцией `SORT`, `HASH` или `BITMAP MERGE`.
- В рамках одиночного запроса может выполняться множество операций, которые используют эту память, и может быть создано несколько областей хеширования/сортировки. Помните, что одновременно может быть открыто много курсоров, каждый из которых предъявляет собственные требования к `SORT_AREA_RETAINED`. Таким образом, если вы установите размер области сортировки в 10 Мбайт, то сможете применять в сеансе 10, 100, 1000 или более мегабайтов ОЗУ. Эти значения не являются предельными значениями сеанса, а определяют лимиты для одиночной операции; внутри сеанса могут выполняться многочисленные операции сортировки в одном запросе или открываться много запросов, требующих одной сортировки.
- Память для этих областей выделяется “по мере необходимости”. Если вы установите размер области сортировки в 1 Гбайт, как было сделано в рассмотренном примере, то это вовсе не означает, что вы распределяете 1 Гбайт ОЗУ. Это значит лишь то, что вы предоставляете данному процессу Oracle право выделения такого объема памяти для выполнения операции сортировки/хеширования.

Автоматическое управление памятью PGA

Начиная с Oracle9i Release 1, был введен новый способ управления памятью PGA, исключающий использование параметров `SORT_AREA_SIZE`, `BITMAP_MERGE_AREA_SIZE` и `HASH_AREA_SIZE`. Он был предназначен для решения следующих проблем.

- **Простота применения.** Способ правильной установки параметров `*_AREA_SIZE` вызывал немалую путаницу. Вдобавок возникали сложности с тем, как эти параметры в действительности работали, и каким образом распределялась память.

- **Ручное управление памятью было методом “один размер на все случаи”.** Обычно по мере роста числа пользователей, одновременно выполняющих похожие приложения в отношении базы данных, объем памяти, задействованной под сортировку и хеширование, увеличивался по линейному закону. Если 10 параллельно работающих пользователей с областью сортировки размером 1 Мбайт потребляли 10 Мбайт памяти, то 100 одновременно работающих пользователей, скорее всего, будут использовать 100 Мбайт, 1000 пользователей — 1000 Мбайт и т.д. Если только администратор базы данных не сидит за консолью, постоянно изменяя размеры областей сортировки/хеширования, то весьма вероятно, что в течение всего рабочего дня все пользователи будут иметь дело с одними и теми же значениями. Вспомните, как в предыдущем примере количество операций физического ввода-вывода во временном пространстве снижалось с увеличением разрешенного для применения объема ОЗУ. Если вы самостоятельно выполните это упражнение, то почти наверняка заметите уменьшение времени отклика с увеличением объема ОЗУ, доступного для сортировки. Выделение памяти вручную фиксирует объем памяти, который будет использоваться для сортировки, на уровне более или менее постоянного числа, независимо от действительно доступного объема памяти. Автоматическое управление памятью позволяет пользоваться этой памятью, когда она доступна; такой метод управления динамически регулирует объем применяемой памяти на основе рабочей нагрузки.
- **Контроль расхода памяти.** Предыдущая проблема делала трудным, если не невозможным, удержание экземпляра Oracle в рамках разумного распределения памяти. Контролировать объем памяти, который экземпляр собирался использовать, было невозможно, т.к. не существовало никакого реального средства управления количеством одновременно выполняющихся операций сортировки и хеширования. Слишком легко можно было попытаться задействовать больше реальной памяти (действительно свободной физической памяти), чем доступно на компьютере.

Давайте приступим к исследованию автоматического управления памятью PGA. Сначала нужно просто установить размер области SGA. Эта область является фрагментом памяти фиксированного размера, поэтому ее размер можно определить очень точно, и он будет ее полным размером (до тех пор, пока вы его не измените). Затем вы сообщаете СУБД Oracle о том, что этим объемом памяти нужно попробовать ограничить все *рабочие области* (новое обобщающее понятие для применяемых областей сортировки и хеширования). После этого на компьютере с 2 Гбайт физической памяти теоретически можно выделить 768 Мбайт памяти под область SGA и 768 Мбайт под область PGA, оставив 512 Мбайт для операционной системы и других процессов. Формулировка “теоретически” используется потому, что на самом деле все работает не настолько точно, хотя и близко к этому. Прежде чем приступить к обсуждению причин, почему это так, мы рассмотрим настройку автоматического управления памятью PGA и его включение.

Процесс настройки этого управления памятью предусматривает выбор подходящих значений для двух параметров инициализации экземпляра.

- `WORKAREA_SIZE_POLICY`. Этот параметр может быть установлен либо в `MANUAL`, в результате чего для управления объемом выделенной памяти будут применяться параметры размеров областей сортировки и хеширования, либо в `AUTO`, когда объем выделенной памяти будет варьироваться на основе памяти `PGA`, в текущий момент распределенной для экземпляра. Стандартным и рекомендуемым значением является `AUTO`.
- `PGAAggregateTarget`. Этот параметр управляет суммарным объемом памяти, который экземпляр должен выделять для всех рабочих областей, используемых для выполнения сортировки или хеширования данных. Его стандартное значение зависит от версии и может быть установлено посредством разнообразных инструментов, таких как `DBCA` (Database Configuration Assistant — помощник по конфигурированию сервера баз данных). В общем случае, если вы применяете автоматическое управление памятью `PGA`, то должны устанавливать этот параметр явным образом.

На заметку! В Oracle 11g Release 1 и последующих версиях вместо установки `PGAAggregateTarget` можно установить параметр `MemoryTarget`. Когда экземпляр использует параметр `MemoryTarget`, он решает, сколько памяти выделить для `SGA` и `PGA` соответственно. Экземпляр может также решить повторно выделить эти объемы памяти прямо во время функционирования базы данных. Однако этот факт не затрагивает работу автоматического управления памятью `PGA` (описанную далее в главе); взамен экземпляр просто определяет установку `PGAAggregateTarget`.

Итак, предполагая, что параметр `WORKAREA_SIZE_POLICY` установлен в `AUTO`, а параметр `PGAAggregateTarget` имеет ненулевое значение, будет применяться новый метод автоматического управления памятью `PGA`. Его можно “включить” в сеансе посредством команды `ALTER SESSION` или на уровне системы с помощью команды `ALTER SYSTEM`.

На заметку! Не забывайте о предыдущем предостережении относительно того, что в Oracle9i подключения посредством разделяемого сервера не будут использовать автоматическое управление памятью. Вместо этого для определения объема ОЗУ, распределяемого под различные операции, они будут применять параметры `SortAreaSize` и `HashAreaSize`. В Oracle 10g и последующих версиях автоматическое управление памятью `PGA` доступно для обоих типов подключений. В случае использования подключений через разделяемый сервер в Oracle9i важно правильно установить параметры `SortAreaSize` и `HashAreaSize`.

Таким образом, основная цель применения автоматического управления памятью `PGA` — максимизация использования ОЗУ с одновременным предотвращением применения большего объема ОЗУ, чем требуется. В условиях ручного управления памятью эта цель была практически недостижимой. При установке `SortAreaSize` в 10 Мбайт один пользователь, выполняющий операцию сортировки, задействовал бы до 10 Мбайт памяти под область сортировки. Если бы 100 пользователей делали то же самое, они заняли бы до 1000 Мбайт памяти. Имея в системе 500 Мбайт свободной памяти, единственный пользователь, выполняющий сортировку, мог бы занять значительно больший объем памяти, а 100 пользователей вынуждены были

бы довольствоваться существенно меньшим объемом. Именно для решения таких проблем и было разработано автоматическое управление памятью PGA. При незначительной рабочей нагрузке использование памяти может быть максимальным, а по мере увеличения нагрузки и роста числа пользователей, выполняющих операции сортировки или хеширования, объем распределяемой для них памяти будет уменьшаться, обеспечивая эксплуатацию всего доступного ОЗУ, но не допуская попытки запрашивать больше памяти, чем существует физически.

Определение способа выделения памяти

Очень часто задают вопросы наподобие “Каким образом эта память выделена?” и “Какой объем ОЗУ используется сеансом?”. На них трудно ответить по той простой причине, что алгоритмы обслуживания памяти при автоматическом управлении не документированы и могут (и будут) изменяться от выпуска к выпуску. Когда применяется автоматическое управление, вы частично утрачиваете контроль над процессами, поскольку лежащие в основе алгоритмы самостоятельно принимают решения о том, что следует делать, и как управлять компонентами.

На основе примечаний MOS (My Oracle Support — Моя поддержка Oracle) под номерами 147806.1 и 223730.1 можно сделать ряд наблюдений.

- Параметр `PGA_AGGREGATE_TARGET` является целью верхнего предела. Это не объем, который заранее выделяется при запуске базы данных. В этом легко убедиться, установив `PGA_AGGREGATE_TARGET` в значение, которое намного больше, чем доступный объем физической памяти сервера. В результате вы не увидите никаких крупных распределений памяти (запомните, что если вы установите `MEMORY_TARGET`, а затем `PGA_AGGREGATE_TARGET` в значение, превышающее `MEMORY_TARGET`, то при начальном запуске экземпляра Oracle будет сгенерирована ошибка `ORA-00838`, которая не позволит запустить экземпляр).
- Объем памяти PGA, доступной заданному сеансу, выводится из установки `PGA_AGGREGATE_TARGET`. Алгоритм для определения максимального размера, используемый процессом, зависит от версии базы данных. Объем памяти PGA, распределяемой процессом, обычно является функцией от доступного объема памяти и количества процессов, соперничающих за пространство.
- По мере возрастания рабочей нагрузки экземпляра (при увеличении числа параллельно выполняющихся запросов и одновременно работающих пользователей) объем памяти PGA, выделенной рабочим областям, будет уменьшаться. База данных будет пытаться сохранять суммарный объем памяти, распределенной всем процессам, в пределах, указанных в параметре `PGA_AGGREGATE_TARGET`. Это аналогично тому, как если бы администратор базы данных в течение всего дня сидел за консолью и устанавливал параметры `SORT_AREA_SIZE` и `HASH_AREA_SIZE` в соответствии с объемом работы, выполняемой в базе данных. Вскоре мы посредством теста исследуем такое поведение.

Хорошо, а как выяснить размеры различных рабочих областей, выделенные сеансу? Это можно сделать с помощью той же методики, которая применялась ранее в сеансе с ручным управлением памятью PGA для выяснения объема памяти, используемой сеансом, и количества выполненных операций ввода-вывода во временном пространстве. Я выполнил следующий тест на четырехпроцессорной машине Oracle

Linux с установленной версией Oracle 12.1.0.1 и подключениями посредством выделенного сервера. Мы начнем с создания таблицы для хранения показателей, которые собираемся отслеживать (приведенный ниже код помещен в файл по имени stats.sql):

```
create table sess_stats
as
select name, value, 0 active
  from
  (
select a.name, b.value
  from v$statname a, v$sesstat b
 where a.statistic# = b.statistic#
    and b.sid = (select sid from v$mystat where rownum=1)
    and (a.name like '%ga %'
        or a.name like '%direct temp%')
 union all
select 'total: ' || a.name, sum(b.value)
  from v$statname a, v$sesstat b, v$session c
 where a.statistic# = b.statistic#
    and (a.name like '%ga %'
        or a.name like '%direct temp%')
    and b.sid = c.sid
    and c.username is not null
 group by 'total: ' || a.name
 );
```

Для представления показателей в этой таблице мы будем применять следующие столбцы.

- NAME — название статистической информации, которую мы собираем (информация PGA и UGA из V\$SESSTAT для текущего сеанса плюс вся информация о памяти экземпляра базы данных, а также о записи во временное табличное пространство).
- VALUE — значение данного показателя.
- ACTIVE — количество сеансов, выполняющих работу в экземпляре. Перед началом мы предполагаем, что экземпляр находится в состоянии ожидания; наш сеанс в данный момент является единственным пользовательским сеансом, поэтому значение равно 0.

Затем в интерактивном сеансе запускается показанный ниже сценарий SQL*Plus (сохраненный в файле single_load.sql). Таблица T была создана заранее и содержит около 70 000 строк.

```
connect eoda/foo
set echo on
declare
  l_first_time boolean default true;
begin
  for x in ( select * from t order by 1, 2, 3, 4 )
  loop
    if ( l_first_time )
```



```

then
    insert into sess_stats
    ( name, value, active )
    select name, value,
           (select count(*)
            from v$session
            where status = 'ACTIVE'
            and username is not null)
    from
    (
    select a.name, b.value
    from v$statname a, v$sesstat b
    where a.statistic# = b.statistic#
    and b.sid = (select sid from v$mystat where rownum=1)
    and (a.name like '%ga %'
         or a.name like '%direct temp%')
    union all
    select 'total: ' || a.name, sum(b.value)
    from v$statname a, v$sesstat b, v$session c
    where a.statistic# = b.statistic#
    and (a.name like '%ga %'
         or a.name like '%direct temp%')
    and b.sid = c.sid
    and c.username is not null
    group by 'total: ' || a.name
    );
    l_first_time := false;
end if;
end loop;
end;
/
commit;

```

Этот сценарий сортирует большую таблицу T, используя автоматическое управление памятью PGA. Затем для текущего сеанса он захватывает все настройки памяти PGA/UGA, а также активность, связанную с сортировкой на диске. Кроме того, конструкция UNION ALL добавляет те же показатели (общая память PGA, общая память UGA и т.д.) для уровня системы. Сценарий был запущен в базе данных со следующими параметрами инициализации:

```

*.compatible='12.1.0.1'
*.control_files='/u01/dbfile/ORAI2CR1/control01.ctl','/u02/dbfile/ORAI2CR1/control02.ctl'
*.db_block_size=8192
*.db_name='ORAI2CR1'
*.pga_aggregate_target=256m
*.sga_target=256m
*.open_cursors=300
*.processes=600
*.remote_login_passwordfile='EXCLUSIVE'
*.resource_limit=TRUE
*.undo_tablespace='UNDOTBS1'

```

Эти параметры указывают на применение автоматического управления памятью PGA со значением PGA_AGGREGATE_TARGET, равным 256 Мбайт, т.е. Oracle отведет до 256 Мбайт памяти PGA для сортировки.

Я подготовил еще один сценарий для запуска в других сеансах с целью генерации крупной нагрузки, связанной с сортировкой. В этом сценарии организован цикл и используется встроенный пакет DBMS_ALERT для проверки необходимости продолжения обработки. Если обработка должна продолжаться, в сценарии запускается тот же самый большой запрос, сортирующий все содержимое таблицы T. По завершении эмуляции сеанс сигнализирует всем процессам сортировки, т.е. генераторам нагрузки, о необходимости “останова” и завершения. Сценарий (сохраненный в файле gen_load.sql), применяемый для выполнения сортировки, выглядит следующим образом:

```
declare
    l_msg long;
    l_status number;
begin
    dbms_alert.register( 'WAITING' );
    for i in 1 .. 999999 loop
        dbms_application_info.set_client_info( i );
        dbms_alert.waitone( 'WAITING', l_msg, l_status, 0 );
        exit when l_status = 0;
        for x in ( select * from t order by 1, 2, 3, 4 )
            loop
                null;
            end loop;
        end loop;
    end;
/
exit
```

А вот сценарий (сохраненный в файле stop.sql), который останавливает выполнение процессов:

```
begin
    dbms_alert.signal( 'WAITING', '' );
    commit;
end;
/
```

Чтобы понаблюдать за отличиями в объемах ОЗУ, выделенных измеряемому сеансу, вначале я выполнил запрос SELECT в изоляции — как единственный сеанс. Я захватил ту же статистическую информацию и сохранил ее в таблице SESS_STATS вместе со счетчиком активных сеансов. Затем я добавил в систему 25 сеансов (т.е. запустил приведенный ранее эталонный сценарий (gen_load.sql) с циклом for i in 1 .. 999999 loop в 25 новых сеансах). Выждав некоторое время (одну минуту, чтобы система смогла приспособиться к этой новой нагрузке), я создал новый сеанс и запустил запрос для выполнения сортировки, захватывая статистические сведения, как и ранее. Я проделывал это многократно, постепенно увеличивая количество одновременно работающих пользователей до 500.

Совет. Сценарии, задействованные в этом эксперименте, доступны для загрузки на веб-сайте издательства. Сценарий `run.sql` в каталоге `ch04` автоматизирует тест, описанный в настоящем разделе.

Следует отметить, что здесь я потребовал у экземпляра базы данных сделать невозможное. Как упоминалось ранее, первый запуск сценария `watch_stat.sql` свидетельствует о том, что каждое подключение к базе данных Oracle — еще до выполнения первой операции сортировки — потребляет чуть более 0,5 Мбайт ОЗУ. При наличии 500 пользователей только их вход в базу данных, не говоря уже о действительном выполнении какой-то работы, привел бы к потреблению памяти с объемом, очень близким к пределу, который определен в параметре `PGA_AGGREGATE_TARGET`! Это служит дополнительным свидетельством того, что параметр `PGA_AGGREGATE_TARGET` является целью, а не директивой. Мы можем, да и будем, превышать его значение по разным причинам. Теперь все готово для формирования отчета о результатах; из-за ограниченного пространства мы прекратили вывод на 275 пользователях, т.к. данные стали повторяться:

```

EODA@ORA12CR1> column active format 999
EODA@ORA12CR1> column pga format 999.9
EODA@ORA12CR1> column "tot PGA" format 999.9
EODA@ORA12CR1> column pga_diff format 999.99
EODA@ORA12CR1> column "temp write" format 9,999
EODA@ORA12CR1> column "tot writes temp" format 99,999,999
EODA@ORA12CR1> column writes_diff format 9,999,999
EODA@ORA12CR1> select active,
2      pga,
3      "tot PGA",
4      "tot PGA"-lag( "tot PGA" ) over (order by active) pga_diff,
5      "temp write",
6      "tot writes temp",
7      "tot writes temp"-lag( "tot writes temp" )
      ↳over (order by active) writes_diff
8  from (
9  select *
10  from (
11  select active,
12         name,
13         case when name like '%ga mem%' then round(value/1024/1024,1)
      ↳else value end val
14  from sess_stats
15  where active < 275
16  )
17  pivot ( max(val) for name in (
18         'session pga memory' as "PGA",
19         'total: session pga memory' as "tot PGA",
20         'physical writes direct temporary tablespace' as "temp write",
21         'total: physical writes direct temporary tablespace'
      ↳as "tot writes temp"
22         ) )
23  )
24  order by active
25  /

```

| ACTIVE | PGA | tot PGA | PGA_DIFF | temp write | tot writestemp | WRITES_DIFF |
|--------|------|---------|----------|------------|----------------|-------------|
| 0 | 3.5 | 7.6 | 0 | 0 | | |
| 1 | 15.2 | 19.5 | 11.90 | 0 | 0 | 0 |
| 26 | 15.2 | 195.6 | 176.10 | 0 | 243,387 | 243,387 |
| 51 | 7.7 | 292.7 | 97.10 | 1,045 | 518,246 | 274,859 |
| 76 | 5.2 | 188.7 | -104.00 | 3,066 | 941,324 | 423,078 |
| 101 | 5.2 | 232.6 | 43.90 | 6,323 | 1,834,035 | 892,711 |
| 126 | 5.2 | 291.8 | 59.20 | 6,351 | 3,021,485 | 1,187,450 |
| 151 | 5.1 | 345.0 | 53.20 | 6,326 | 4,783,879 | 1,762,394 |
| 177 | 5.0 | 403.3 | 58.30 | 6,321 | 8,603,295 | 3,819,416 |
| 201 | 5.2 | 453.2 | 49.90 | 6,327 | 12,848,568 | 4,245,273 |
| 226 | 4.8 | 507.5 | 54.30 | 6,333 | 15,225,399 | 2,376,831 |
| 251 | 5.1 | 562.2 | 54.70 | 6,315 | 17,579,502 | 2,354,103 |

12 rows selected.

12 строк выбрано.

Прежде чем анализировать результаты, давайте посмотрим на запрос, использованный для формирования отчета. В нем применяется средство, доступное начиная с версии Oracle 11g Release 1, которое называется *вращением* (pivot) результирующего набора. Ниже приведен альтернативный способ записи строк с 11 по 24 этого SQL-запроса, который будет работать в Oracle 10g Release 2 и предшествующих версиях:

```

11 select active,
12    max( decode(name, 'session pga memory', val) ) pga,
13    max( decode(name, 'total: session pga memory', val) ) as "tot PGA",
14    max( decode(name, 'physical writes direct temporary tablespace',
15               val) ) as "temp write",
15    max( decode(name, 'total: physical writes direct temporary tablespace',
16               val) ) as "tot writes temp"
16 from (
17 select active,
18    name,
19    case when name like '%ga mem%' then round(value/1024/1024,1)
20    else value end val
21 from sess_stats
22 where active < 275
23 )
24 group by active

```

Эта часть запроса извлекает записи из таблицы показателей, когда есть менее 275 активных сеансов, преобразовывает показатели памяти (UGA/PGA) из байтов в мегабайты и затем вращает четыре интересующих нас показатели, делая из строк столбцы. После получения этих четырех показателей в одной записи с помощью аналитической функции LAG() к каждой строке были добавлены значения общего наблюдаемого объема PGA и общего наблюдаемого количества операций ввода-вывода во временную область для предыдущих строк. Это позволяет легко видеть последовательные различия в этих значениях. Вернемся к данным — как видите, при наличии небольшого числа активных сеансов все операции сортировки происходили полностью в памяти. При количестве активных сеансов от 1 до 50 сортировка происходила полностью в памяти. Однако к моменту, когда в систему вошло 50 пользователей, активно выполняющих сортировку, базе данных стало не хватать объема памяти, который было

разрешено использовать одновременно. Для того чтобы объем памяти PGA уместился в допустимые пределы (запрошенные 256 Мбайт), потребовалась пара минут, но в конечном итоге это происходило на таком низком уровне одновременно работающих пользователей. Объем памяти PGA, выделенной сеансу, снизился с 15,2 Мбайт до 7,7 Мбайт и остановился на показателе примерно 5,2 Мбайт (вспомните, что части памяти PGA распределяются не для рабочей области (сортировки), а для других операций; один лишь вход пользователя приводит к выделению 0,5 Мбайт памяти PGA). Общий объем памяти PGA, задействованной системой, оставался в приемлемых пределах до тех пор, пока в системе не появилось приблизительно 126 пользователей. Затем объем памяти начал регулярно превышать значение `PGA_AGGREGATE_TARGET` и это сохранилось вплоть до конца теста. В этом случае перед экземпляром базы была поставлена неразрешимая задача — выделенных 256 Мбайт ОЗУ просто не хватало для обеспечения одновременной работы 126 пользователей, большинство из которых выполняли код PL/SQL и операции сортировки. Это просто невозможно было сделать. Таким образом, каждый сеанс применял минимально возможный объем памяти, но вынужден был выделять столько памяти, сколько ему требовалось. К моменту завершения теста активные сеансы использовали всего около 560 Мбайт памяти PGA — настолько мало, насколько это было возможно.

Тем не менее, мы должны обдумать, как этот вывод выглядел бы при ручном управлении памятью. Предположим, что параметр `SORT_AREA_SIZE` установлен в 5 Мбайт. Математическая оценка очень проста: каждый сеанс смог бы выполнять сортировку в ОЗУ (или в виртуальной памяти при нехватке реального ОЗУ) и, следовательно, каждый сеанс потреблял бы от 6 до 7 Мбайт ОЗУ (объем, задействованный без выполнения сортировки на диске в предыдущем случае с единственным пользователем). Использование памяти выглядело примерно так:

```
EODA@ORA12CR1> column total_pga format 9,999
EODA@ORA12CR1> with data(users)
 2 as
 3 (select 1 users from dual
 4  union all
 5  select users+25 from data where users+25 <= 275)
 6 select users, 7 my_pga, 7*users total_pga
 7   from data
 8   order by users
 9 /
```

| USERS | MY_PGA | TOTAL_PGA |
|-------|--------|-----------|
| ----- | ----- | ----- |
| 1 | 7 | 7 |
| 26 | 7 | 182 |
| 51 | 7 | 357 |
| 76 | 7 | 532 |
| 101 | 7 | 707 |
| 126 | 7 | 882 |
| 151 | 7 | 1,057 |
| 176 | 7 | 1,232 |
| 201 | 7 | 1,407 |
| 226 | 7 | 1,582 |
| 251 | 7 | 1,757 |

11 rows selected.

На заметку! В этом запросе применяется прием рекурсивного выноса подзапроса (recursive subquery factoring), который доступен только в Oracle 11g Release 2 и последующих версиях. В более ранних выпусках приведенный запрос работать не будет.

Если бы я запустил этот тест (на сервере установлено 2 Гбайт реальной памяти, а объем области SGA составляет 256 Мбайт), то к моменту появления в системе 250 пользователей подкачка и страничный обмен происходили бы настолько часто, что продолжение работы оказалось бы невозможным; при 500 пользователях пришлось бы выделить около 3 514 Мбайт ОЗУ! Следовательно, в такой системе администратор базы данных, скорее всего, установил бы значение SORT_AREA_SIZE не в 5 Мбайт, а в 0,5 Мбайт, пытаясь удержать максимальный применяемый объем PGA на приемлемом уровне. При 500 пользователях выделенный объем PGA составил бы около 500 Мбайт, что примерно соответствовало бы ситуации с автоматическим управлением памятью, но при меньшем количестве пользователей вместо сортировки в памяти все равно выполнялась бы запись во временное пространство на диске.

Ручное управление памятью демонстрирует вполне предсказуемое — но далеко не оптимальное — использование памяти при увеличении или уменьшении рабочей нагрузки. Автоматическое управление PGA было разработано именно для того, чтобы позволить небольшому сообществу пользователей задействовать максимально возможный объем ОЗУ, когда он доступен, уменьшая потребление памяти при возрастании рабочей нагрузки и увеличивая объем памяти, выделенной для отдельных операций, при снижении нагрузки.

Использование параметра `PGA_AGGREGATE_TARGET` для управления выделением памяти

Ранее было указано, что “теоретически” параметр `PGA_AGGREGATE_TARGET` можно применять для управления общим объемом памяти PGA, используемой экземпляром. Однако последний пример показал, что этот предел не является жестким. Экземпляр попытается остаться в границах, определенных `PGA_AGGREGATE_TARGET`, но если ему это не удастся, он не прекратит обработку, а просто превысит данный порог.

Еще одна причина того, что этот предел является “теоретическим”, связана с тем, что хотя рабочие области составляют крупную часть памяти PGA, они не являются единственными ее фрагментами. Область PGA состоит из множества частей, и только рабочие области находятся под контролем экземпляра базы данных. Если создать и выполнить блок кода PL/SQL, который заполняет данными большой массив в режиме выделенного сервера, когда область UGA размещена в области PGA, то Oracle не останется ничего другого, как разрешить это.

Рассмотрим следующий короткий пример. Мы создадим пакет, который может хранить некоторые постоянные (глобальные) данные на сервере:

```
EODA@ORA12CR1> create or replace package demo_pkg
2 as
3     type array is table of char(2000) index by binary_integer;
4     g_data array;
5 end;
6 /
Package created.
Пакет создан.
```

Теперь измерим объем памяти, который в текущий момент наш сеанс использует в PGA/UGA (в этом примере применяется выделенный сервер, поэтому UGA является подобластью памяти PGA):

```
EODA@ORA12CR1> select a.name, to_char(b.value, '999,999,999') bytes,
2      to_char(round(b.value/1024/1024,1), '99,999.9' ) mbytes
3  from v$statname a, v$mystat b
4  where a.statistic# = b.statistic#
5    and a.name like '%ga memory%';
```

| NAME | BYTES | MBYTES |
|------------------------|-----------|--------|
| ----- | ----- | ----- |
| session uga memory | 1,526,568 | 1.5 |
| session uga memory max | 1,526,568 | 1.5 |
| session pga memory | 2,208,088 | 2.1 |
| session pga memory max | 2,208,088 | 2.1 |

Первоначально наш сеанс использует около 2,1 Мбайт памяти PGA (как результат компиляции пакета PL/SQL, выполнения запроса и т.д.). Теперь снова запустим этот запрос применительно к таблице T с применением того же значения PGA_AGGREGATE_TARGET, равного 256 Мбайт (это было сделано в экземпляре, который в противном случае бы простаивал; мы затребуем выделение памяти для сеанса прямо сейчас):

```
EODA@ORA12CR1> set autotrace traceonly statistics;
EODA@ORA12CR1> select * from t order by 1,2,3,4;

72616 rows selected.
72616 строк выбрано.
```

Statistics

| | |
|---------|--|
| ----- | |
| 105 | recursive calls |
| 0 | db block gets |
| 1103 | consistent gets |
| 993 | physical reads |
| 0 | redo size |
| 3665844 | bytes sent via SQL*Net to client |
| 53795 | bytes received via SQL*Net from client |
| 4843 | SQL*Net roundtrips to/from client |
| 1 | sorts (memory) |
| 0 | sorts (disk) |
| 72616 | rows processed |

```
EODA@ORA12CR1> set autotrace off
```

Как видите, сортировка была выполнена полностью в памяти, и фактически взглянув на использование памяти PGA/UGA со стороны сеанса, можно выяснить объем потребляемой памяти:

```
EODA@ORA12CR1> select a.name, to_char(b.value, '999,999,999') bytes,
2      to_char(round(b.value/1024/1024,1), '99,999.9' ) mbytes
3  from v$statname a, v$mystat b
4  where a.statistic# = b.statistic#
5    and a.name like '%ga memory%';
```

| NAME | BYTES | MBYTES |
|-------------------------------|-------------------|-------------|
| ----- | ----- | ----- |
| session uga memory | 1,854,008 | 1.8 |
| session uga memory max | 11,213,280 | 10.7 |
| session pga memory | 2,470,232 | 2.4 |
| session pga memory max | 12,104,024 | 11.5 |

Легко заметить, что используются 11,5 Мбайт ОЗУ с примерным диапазоном в 15 Мбайт, который мы наблюдали ранее в предыдущем тесте сортировки. Теперь заполним массив CHAR, имеющийся в пакете (тип данных CHAR дополняется пробелами, поэтому каждый элемент массива имеет длину точно 2000 символов):

```

EODA@ORA12CR1> begin
2       for i in 1 .. 200000
3         loop
4           demo_pkg.g_data(i) := 'x';
5         end loop;
6 end;
7 /

```

PL/SQL procedure successfully completed.

Процедура PL/SQL успешно завершена.

Измерение текущей утилизации памяти PGA нашего сеанса дает следующие результаты:

```

EODA@ORA12CR1> select a.name, to_char(b.value, '999,999,999') bytes,
2       to_char(round(b.value/1024/1024,1), '99,999.9' ) mbytes
3   from v$statname a, v$mystat b
4  where a.statistic# = b.statistic#
5    and a.name like '%ga memory%';

```

| NAME | BYTES | MBYTES |
|------------------------|-------------|--------|
| ----- | ----- | ----- |
| session uga memory | 469,569,304 | 447.8 |
| session uga memory max | 469,569,304 | 447.8 |
| session pga memory | 470,921,560 | 449.1 |
| session pga memory max | 470,921,560 | 449.1 |

Теперь область PGA содержит выделенные сегменты, которыми сам экземпляр управлять не может. Мы уже превысили предел, установленный в PGA_AGGREGATE_TARGET для всего экземпляра *в этом отдельном сеансе*, и база данных ничего не может с этим поделать. Если бы она действовала иначе, то ей пришлось бы отклонить наш запрос, и она сделала бы это только в случае, если бы операционная система сообщила об отсутствии памяти для выделения (ORA-04030). При желании можно было бы зарезервировать больше памяти для этого массива и поместить в него больше данных — и экземпляр всего лишь бы выполнил такое действие.

Тем не менее, экземпляр осведомлен о том, что мы сделали. Он вовсе не игнорирует память, которой не может управлять; он просто идентифицирует, что память используется, и соответствующим образом уменьшает размер памяти, выделенной для рабочих областей. Поэтому, если снова выполнить тот же запрос сортировки, то выяснится, что на этот раз сортировка осуществляется на диске — экземпляр не предоставит около 12 Мбайт ОЗУ, которые необходимы для выполнения этой операции в памяти, поскольку значение PGA_AGGREGATE_TARGET уже превышено:


```

EODA@ORA12CR1> set autotrace traceonly statistics;
EODA@ORA12CR1> select * from t order by 1,2,3,4;

72616 rows selected.
72616 строк выбрано.

```

Statistics

```

-----
          9  recursive calls
          8  db block gets
         986  consistent gets
        2025  physical reads
           0  redo size
    3665844  bytes sent via SQL*Net to client
         53795  bytes received via SQL*Net from client
         4843  SQL*Net roundtrips to/from client
           0  sorts (memory)
           1  sorts (disk)
        72616  rows processed

```

```
EODA@ORA12CR1> set autotrace off
```

Итак, из-за того, что часть памяти PGA находится вне контроля Oracle, легко превысить значение `PGA_AGGREGATE_TARGET`, просто выделяя память для множества действительно крупных структур данных в коде PL/SQL. Я ни в коем случае не рекомендую поступать так, а просто подчеркиваю, что значение `PGA_AGGREGATE_TARGET` — скорее просьба, а не жесткое ограничение.

Выбор между ручным и автоматическим управлением памятью

Так какой же метод следует применять: ручной или автоматический? Я однозначно предпочитаю по умолчанию использовать автоматическое управление памятью PGA.

Внимание! В этой книге я время от времени повторяю: не вносите никаких изменений в производственную — реально действующую — систему без предварительного тестирования на предмет любых побочных эффектов. Например, предположим, что вы не читали эту главу, проверили свою систему и обнаружили, что в ней применяется ручное управление памятью, а затем просто активизировали автоматическое управление. В таком случае планы запросов могут измениться, и это может повлиять на производительность. Возможен один из трех сценариев:

- производительность системы останется в точности прежней;
- производительность увеличится;
- производительность значительно уменьшится.

Будьте внимательны при внесении изменений; в первую очередь проверяйте предложенное изменение.

Одной из наиболее трудных задач для администратора базы данных может быть установка индивидуальных параметров, особенно таких, как `SORT_HASH_AREA_SIZE` и т.д. Мне часто приходилось сталкиваться с системами, работающими с невероятно низкими значениями этих параметров — настолько низкими, что это приводи-

ло к очень серьезному снижению производительности. Возможно, это объясняется тем, что значения таких параметров, устанавливаемые по умолчанию, очень малы: 64 Кбайт для сортировки и 128 Кбайт для хеширования. Существует множество противоречивых мнений по поводу их оптимальных значений. Кроме того, оптимальные значения могут изменяться в течение рабочего дня. В 8 часов утра при наличии двух пользователей размер области сортировки, равный 50 Мбайт, выделенный для одного зарегистрированного пользователя, может быть вполне допустимым. Но в полдень, когда пользователей становится 500, значение 50 Мбайт может оказаться совершенно неприемлемым. Именно в этой ситуации целесообразно применять параметры `WORKAREA_SIZE_POLICY=AUTO` и `PGA_AGGREGATE_TARGET`. Концептуально установка значения `PGA_AGGREGATE_TARGET` — объема памяти, который СУБД Oracle должна быть вправе использовать для выполнения операций сортировки и хеширования — проще, чем попытка подбора идеальных значений параметров `SORT_HASH_AREA_SIZE`, тем более что таких идеальных значений попросту не существует. Оптимальные значения зависят от конкретной рабочей нагрузки.

По исторически сложившейся традиции администраторы баз данных конфигурировали объем памяти, применяемый Oracle, путем установки размера области SGA (кеш буферов, журнальный буфер, разделяемый пул, большой пул и пул Java). Оставшаяся память компьютера должна была использоваться выделенным и разделяемым серверами в области PGA. Администратор базы данных обладал очень незначительным контролем над тем, сколько этой памяти должно было быть задействовано или не задействовано. Он мог установить параметр `SORT_AREA_SIZE`, но при одновременном выполнении 10 операций сортировки СУБД Oracle могла бы применять до $10 * \text{SORT_AREA_SIZE}$ байт ОЗУ. При параллельном выполнении 100 операций сортировки Oracle использовала бы $100 * \text{SORT_AREA_SIZE}$ байт, при 1000 операциях — $1000 * \text{SORT_AREA_SIZE}$ и т.д. Учитывая, что в PGA помещаются и другие объекты, администратор лишен реальной возможности управлять максимальным объемом памяти, который задействуется из области PGA системы.

Желательно, чтобы использование этой памяти регулировалось реальными потребностями в системе. Чем больше пользователей, тем меньший объем ОЗУ каждый из них должен занимать, а чем их меньше, тем больший объем памяти должен быть доступен для применения каждым из них. Установка `WORKAREA_SIZE_POLICY=AUTO` как раз и является способом достижения этой цели. Теперь администратор базы данных указывает единственный размер — `PGA_AGGREGATE_TARGET`, представляющий собой максимальный объем памяти PGA, который база данных должна стремиться использовать. СУБД Oracle будет распределять эту память между активными сеансами так, как сочтет нужным. Более того, Oracle9i Release 2 и последующие версии предлагают консультативную справку по размеру PGA (часть Statspack и AWR, доступная через динамическое представление информации производительности `V$` и отображаемая в окне диспетчера предприятия), похожую на консультативную справку по настройке кеша буферов. Со временем эта справка подскажет, какое значение `PGA_AGGREGATE_TARGET` будет оптимальным для минимизации количества физических операций ввода-вывода во временных табличных пространствах данного экземпляра. Эту информацию можно применять либо для динамического интерактивного изменения размера PGA (при наличии достаточного объема ОЗУ), либо для принятия решения о том, нуждается ли сервер в дополнительном объеме ОЗУ с целью достижения оптимальной производительности.

Однако бывают ли ситуации, в которых пользоваться автоматическим управлением памятью нежелательно? Безусловно, но к счастью они являются скорее исключением, чем правилом. Этот метод был разработан для многопользовательской среды. В ожидании присоединения новых пользователей к системе автоматическое управление памятью будет ограничивать объем выделяемой памяти в виде некоторого процента от значения `PGA_AGGREGATE_TARGET`. Но как быть, если требуется задействовать весь доступный объем памяти? Что ж, в таком случае придется с помощью команды `ALTER SESSION` отключить автоматическое управление памятью в своем сеансе (оставив его включенным во всех остальных) и вручную установить нужные значения параметров `SORT|HASH_AREA_SIZE`. Например, каким образом поступить, когда есть крупный пакетный процесс, который запускается в 2:00 и выполняет большие хеш-соединения, строит индексы и тому подобное? Ему должно быть разрешено эксплуатировать все ресурсы, доступные экземпляру. Этому процессу не нужно “скромничать” в плане использования памяти — ему необходима вся доступная память, поскольку в данный момент он является единственным выполняющимся в базе данных. Совершенно очевидно, что это пакетное задание может выдать команду `ALTER SESSION` и задействовать все доступные ресурсы.

Короче говоря, я предпочитаю применять автоматическое управление памятью для сеансов конечных пользователей — для приложений, которые ежедневно выполняются в базе данных. Ручное управление памятью целесообразно использовать при запуске больших пакетных заданий, функционирующих в те периоды времени, когда они являются единственными активными процессами в экземпляре.

Заключительные соображения по поводу использования областей PGA и UGA

К настоящему времени мы рассмотрели две структуры памяти: PGA и UGA. Теперь вы должны понимать, что область PGA является закрытой для процесса. Она представляет собой набор переменных, которые выделенный или разделяемый сервер Oracle должен хранить независимо от сеанса. Область PGA — это “куча” памяти, в которой могут распределяться другие структуры. Область UGA также является “кучей” памяти, в которой могут быть определены различные структуры, характерные для сеанса. Область UGA распределяется внутри области PGA при подключении к базе данных Oracle посредством выделенного сервера и в области SGA при подключении через разделяемый сервер. Из этого следует, что при использовании разделяемого сервера необходимо определять размер большого пула области SGA так, чтобы обеспечить достаточный объем для обслуживания любого возможного пользователя, который когда-либо подключится к базе данных параллельно. Таким образом, область SGA базы данных, которая поддерживает подключения посредством разделяемого сервера, в общем случае намного больше области SGA аналогично сконфигурированной базы данных, поддерживающей только режим выделенного сервера. Область SGA подробно рассматривается далее в главе.

Системная глобальная область

Каждый экземпляр Oracle имеет одну крупную структуру памяти, которая называется *системной глобальной областью* (System Global Area — SGA). Это большая

структура памяти совместного использования, к которой в тот или иной момент времени будет обращаться каждый процесс Oracle. Ее размер варьируется от десятков мегабайт в небольших тестовых системах до нескольких гигабайт в средних и больших системах и сотен гигабайт в действительно крупных системах.

В среде UNIX/Linux область SGA — это физическая сущность, которую можно “видеть” в командной строке операционной системы. Физически она реализована в виде сегмента разделяемой памяти — автономного фрагмента памяти, к которому могут присоединяться процессы. Допускается иметь в системе область SGA без единого процесса Oracle — память оказывается предоставленной самой себе. Однако следует отметить, что наличие области SGA безо всяких процессов Oracle свидетельствует об аварийном отказе базы данных по какой-либо причине. Эта ситуация необычна, но может встречаться. Вот как область SGA “выглядит” в системе Oracle Linux:

```
$ ipcs -m | grep ora
0x27ba944c 887324675 oracle 640 14680064 82
0x00000000 887357444 oracle 640 1061158912 41
0x749a2e08 887947269 oracle 640 14680064 72
0x00000000 887980038 oracle 640 511705088 36
0x00000000 888537095 oracle 640 8388608 16
0x00000000 888569864 oracle 640 260046848 16
0xc6e51dc4 888602633 oracle 640 2097152 16
```

На заметку! На моей тестовой/демонстрационной машине функционирует несколько экземпляров. Мне понадобилось множество экземпляров для тестирования разнообразных концепций, представленных в этой книге, в разных выпусках. Единственно разумным и корректным количеством экземпляров на рабочей машине является *один*. В реальных условиях никогда не устанавливайте более одного экземпляра на каждый отдельный рабочий сервер. Если вам нужно более одного экземпляра на физическом сервере, вы должны применять виртуализацию, чтобы разделить этот один сервер на несколько виртуальных серверов — каждый с собственным экземпляром Oracle.

Здесь представлены три области SGA, и в отчете показана учетная запись операционной системы, которая владеет SGA (oracle для всех областей в этом примере), а также размер SGA — 1 Гбайт (вторая строка) для первого примера. В среде Windows вы не сможете видеть SGA как отдельную сущность, что возможно в UNIX/Linux. Из-за того, что на платформе Windows база данных Oracle выполняется как одиночный процесс с единственным адресным пространством, область SGA выделяется как закрытая память процесса `oracle.exe`. С помощью диспетчера задач Windows или какого-то другого инструмента мониторинга производительности можно посмотреть, сколько памяти распределено для `oracle.exe`, но видеть SGA отдельно от остальной распределенной памяти не получится.

На заметку! Если только вы не используете мои настройки параметров и не работаете в точности с такой же версией Oracle и в точно той же операционной системе, вы почти наверняка увидите другие показатели. Размеры SGA сильно зависят от версии, операционной системы и параметров.

Внутри самой СУБД Oracle информацию об области SGA можно просматривать независимо от платформы, применяя еще одно “магическое” представление V\$ по имени V\$SGASTAT. Информация может выглядеть примерно так:

```

EODA@ORA12CR1> compute sum of bytes on pool
EODA@ORA12CR1> break on pool skip 1
EODA@ORA12CR1>
EODA@ORA12CR1> select pool, name, bytes
  2   from v$sgastat
  3   order by pool, name;

```

| POOL | NAME | BYTES |
|---------------------|---------------------------|-----------|
| java pool | free memory | 4194304 |
| ***** | | ----- |
| sum | | 4194304 |
| large pool | PX msg pool | 491520 |
| | free memory | 3702784 |
| ***** | | ----- |
| sum | | 4194304 |
| shared pool | 1063.kgght | 36784 |
| | 11G QMN so | 4144 |
| | 177.kggfa | 39840 |
| ... | | |
| | zlllab Group Tree Heap De | 160 |
| ***** | | ----- |
| sum | | 314572800 |
| | buffer_cache | 184549376 |
| | fixed_sga | 2290264 |
| | log_buffer | 7938048 |
| | shared_io_pool | 4194304 |
| ***** | | ----- |
| sum | | 198971992 |
| 1064 rows selected. | | |
| 1064 строк выбрано. | | |

Область SGA разделена на разнообразные пулы. Основные пулы описаны ниже.

- **Пул Java (Java pool).** Это фиксированный объем памяти, выделенной для виртуальной машины Java (Virtual Java Machine — JVM), которая действует в базе данных. В Oracle10g и последующих версиях пул Java может оперативно изменяться во время работы базы данных.
- **Большой пул (large pool).** Этот пул используется подключениями посредством разделяемого сервера для памяти сеанса, средствами параллельного выполнения для буферов сообщений и резервным копированием RMAN для буферов дискового ввода-вывода. Большой пул допускает оперативное изменение размера.
- **Разделяемый пул (shared pool).** Этот пул содержит разделяемые курсоры, хранимые процедуры, объекты состояния, кеша словаря данных и многие десятки других элементов данных. Размеры этого пула могут оперативно изменяться, начиная с версии Oracle9i.
- **Пул Streams (Streams pool).** Этот пул памяти применяется инструментами совместного использования данных, такими как Oracle GoldenGate, Oracle Streams и т.д. Пул Streams доступен в Oracle 10g и последующих версиях и допускает

оперативное изменение своего размера. Если он не сконфигурирован, но функциональность Streams задействована, то Oracle выделит под память потоков до 10% объема разделяемого пула.

- **Неопределенный пул (“Null” pool).** В действительности этот пул не имеет названия. Он представляет собой память, предназначенную для хранения буферов блоков (кешированных блоков базы данных), буфера журнала повторения и “фиксированной области SGA”.

Типичная область SGA может выглядеть так, как показано на рис. 4.1.

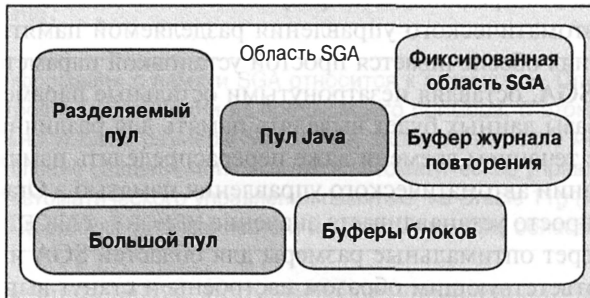


Рис. 4.1. Типичная область SGA

Наибольшее влияние на общий размер области SGA оказывают следующие параметры.

- **JAVA_POOL_SIZE.** Управляет размером пула Java.
- **SHARED_POOL_SIZE.** Управляет размером разделяемого пула (до некоторой степени).
- **LARGE_POOL_SIZE.** Управляет размером большого пула.
- **STREAMS_POOL_SIZE.** Управляет размером пула Streams.
- **DB_*_CACHE_SIZE.** Восемь параметров **CACHE_SIZE** управляют размерами различных доступных кешей буферов.
- **LOG_BUFFER.** Управляет размером буфера повторения (до некоторой степени).
- **SGA_TARGET.** Применяется с автоматическим управлением памятью SGA в Oracle 10g и последующих версиях; может изменяться оперативным образом.
- **SGA_MAX_SIZE.** Используется для управления размером области SGA.
- **MEMORY_TARGET.** Применяется с автоматическим управлением памятью (как для PGA, так и для SGA) в Oracle 11g и последующих версиях.
- **MEMORY_MAX_SIZE.** Используется для управления максимальным объемом памяти, которого СУБД Oracle должна стараться придерживаться в отношении размеров PGA и SGA при автоматическом управлении памятью в Oracle 11g и последующих версиях. В действительности это лишь цель; область PGA может превосходить оптимальный размер, если количество пользователей вырастает сверх некоторого уровня, либо когда сеанс (сеансы) распределяет большие неуправляемые фрагменты памяти, как демонстрировалось ранее.

В Oracle9i различные компоненты SGA требуют настройки своих размеров администратором базы данных вручную. Однако в Oracle 10g и последующих версиях появилась возможность автоматического управления разделяемой памятью, посредством которого экземпляр базы данных будет выделять и перераспределять память для разнообразных компонентов SGA во время выполнения в соответствии с изменяющейся рабочей нагрузкой. Более того, в Oracle 11g доступен еще один вариант: автоматическое управление памятью, при помощи которого экземпляр базы данных не только автоматически управляет памятью SGA и PGA, но также подбирает оптимальные размеры для областей SGA и PGA, автоматически перераспределяя эти части памяти, когда это представляется разумным.

Применение автоматического управления разделяемой памятью в Oracle 10g и последующих версиях обеспечивается простой установкой параметра SGA_TARGET в желаемый размер SGA, оставляя незатронутыми остальные параметры, связанные с SGA. Экземпляр базы данных будет выделять память для различных пулов по мере необходимости и с течением времени даже перераспределять память между пулами.

При использовании автоматического управления памятью в Oracle 11g и последующих версиях вы просто устанавливаете значение MEMORY_TARGET. Экземпляр базы данных затем выберет оптимальные размеры для областей SGA и PGA, и эти компоненты будут соответствующим образом настроены и станут выполнять собственное управление памятью в отведенных им пределах. Более того, база данных может и будет изменять размеры выделенных областей SGA и PGA по мере дальнейшего изменения рабочей нагрузки.

Независимо от применяемого управления памятью — автоматического или ручного — вы обнаружите, что память различным пулам в SGA выделяется блоками, которые называются *гранулами*. Одна гранула (granule) — это область памяти размером 4, 8 или 16 Мбайт. Гранула является наименьшей единицей выделения памяти, поэтому если запросить для пула Java память 5 Мбайт при размере гранулы 4 Мбайт, то в действительности Oracle выделит для пула Java объем 8 Мбайт (8 — наименьшее целое число, большее или равное 5 и кратное размеру гранулы, равному 4). Размер гранулы определяется размером области SGA (в определенной степени это утверждение выглядит рекурсивным, т.к. размер SGA зависит от размера гранулы). Размеры гранул, используемых для каждого пула, можно выяснить, запросив представление V\$SGA_DYNAMIC_COMPONENTS. Фактически это представление можно применять для выяснения влияния общего размера SGA на размер гранул:

```
EODA@ORA12CR1> show parameter sga_target
```

| NAME | TYPE | VALUE |
|------------|-------------|-------|
| sga_target | big integer | 256M |

```
EODA@ORA12CR1> select component, granule_size from v$sga_dynamic_
components;
```

| COMPONENT | GRANULE_SIZE |
|----------------------|--------------|
| shared pool | 4194304 |
| large pool | 4194304 |
| java pool | 4194304 |
| streams pool | 4194304 |
| DEFAULT buffer cache | 4194304 |

```

KEEP buffer cache                                4194304
RECYCLE buffer cache                             4194304
DEFAULT 2K buffer cache                         4194304
DEFAULT 4K buffer cache                         4194304
DEFAULT 8K buffer cache                         4194304
DEFAULT 16K buffer cache                       4194304
DEFAULT 32K buffer cache                       4194304
Shared IO Pool                                  4194304
Data Transfer Cache                             4194304
ASM Buffer Cache                                4194304
15 rows selected.
15 строк выбрано.

```

На заметку! Эта информация о памяти SGA относится к экземпляру Oracle, запущенному с файлом параметров инициализации из предыдущего примера. В этом файле параметров мы самостоятельно указали размеры SGA и PGA. Вследствие этого мы используем автоматическое управление разделяемой памятью и автоматическое управление памятью PGA, но не средство “автоматического управления памятью” из Oracle 11g (и последующих версий), которое само должно устанавливать и изменять размеры областей PGA/SGA.

В этом примере применялось автоматическое управление разделяемой памятью, а размер SGA устанавливался через единственный параметр SGA_TARGET. Когда размер SGA не превышает приблизительно 1 Гбайт, размер гранулы составляет 4 Мбайт. Когда размер SGA увеличивается до некоторого порога, превышающего 1 Гбайт (конкретное значение варьируется в зависимости от операционной системы и даже от выпуска), размер гранулы возрастает. Сначала мы вернемся к использованию хранимого файла параметров, чтобы облегчить изменение SGA_TARGET:

```

SYS@ORA12CR1> create spfile from pfile;
File created.
Файл создан.

SYS@ORA12CR1> startup force;
ORACLE instance started.
Экземпляр ORACLE запущен.

Total System Global Area                267227136 bytes
Fixed Size                              2287336 bytes
Variable Size                          180357400 bytes
Database Buffers                        79691776 bytes
Redo Buffers                            4890624 bytes
Database mounted.
Database opened.
Общий размер системной глобальной области      267227136 байтов
Размер фиксированной области                  2287336 байтов
Размер переменной области                    180357400 байтов
Буферы базы данных                          79691776 байтов
Буферы повторения                           4890624 байтов
База данных смонтирована.
База данных открыта.

```

На заметку! Если ваш экземпляр в текущий момент функционирует, то команда STARTUP FORCE завершит его работу (режим прекращения) и затем перезапустит.

Далее мы модифицируем параметр SGA_TARGET:

```
SYS@ORA12CR1> alter system set sga_target = 1512m scope=spfile;
System altered.
```

Система изменена.

```
SYS@ORA12CR1> startup force
```

```
ORACLE instance started.
```

Экземпляр ORACLE запущен.

```
Total System Global Area          1586708480 bytes
Fixed Size                          2288824 bytes
Variable Size                      402654024 bytes
Database Buffers                   1174405120 bytes
Redo Buffers                        7360512 bytes
```

```
Database mounted.
```

```
Database opened.
```

```
SYS@ORA12CR1> show parameter sga_target
```

| NAME | TYPE | VALUE |
|------------|-------------|-------|
| sga_target | big integer | 1520M |

Теперь посмотрим на компоненты SGA:

```
SYS@ORA12CR1> select component, granule_size from v$sga_dynamic_components;
```

| COMPONENT | GRANULE_SIZE |
|--------------------------|--------------|
| shared pool | 16777216 |
| large pool | 16777216 |
| java pool | 16777216 |
| streams pool | 16777216 |
| DEFAULT buffer cache | 16777216 |
| KEEP buffer cache | 16777216 |
| RECYCLE buffer cache | 16777216 |
| DEFAULT 2K buffer cache | 16777216 |
| DEFAULT 4K buffer cache | 16777216 |
| DEFAULT 8K buffer cache | 16777216 |
| DEFAULT 16K buffer cache | 16777216 |
| DEFAULT 32K buffer cache | 16777216 |
| Shared IO Pool | 16777216 |
| Data Transfer Cache | 16777216 |
| ASM Buffer Cache | 16777216 |

```
15 rows selected.
```

15 строк выбрано.

Несложно заметить, что при размере SGA в 1,5 Гбайт память под пулы будет выделяться с гранулами 16 Мбайт, поэтому размер любого пула будет кратным 16 Мбайт. Имея это в виду, давайте по очереди рассмотрим основные компоненты SGA.

Фиксированная область SGA

Фиксированная область SGA (fixed SGA) — это компонент SGA, размер которого зависит от платформы и выпуска. Эта область “компилируется” в сам двоичный модуль Oracle во время установки (отсюда и название “фиксированная”). Фиксированная область SGA содержит набор переменных, которые указывают на другие компонен-

ты SGA, и переменных, содержащих значения различных параметров. Размер фиксированной области SGA не доступен для управления со стороны администратора базы данных или программиста и, как правило, он очень мал. Эту область можно считать разделом “начальной загрузки” SGA — чем-то таким, что Oracle использует внутренне для нахождения других элементов и фрагментов области SGA.

Буфер повторения

Буфер повторения (redo buffer) — это место, где данные, которые должны быть записаны в оперативные журналы повторения (online redo log), будут временно кешироваться перед их записью на диск. Поскольку передача из памяти в память выполняется значительно быстрее передачи из памяти на диск, применение буфера журнала повторения может ускорить работу базы данных. Данные будут храниться в буфере повторения не очень долго. В действительности процесс LGWR инициирует сброс содержимого этой области на диск в одном из перечисленных ниже сценариев:

- каждые три секунды;
- когда выдается COMMIT или ROLLBACK;
- когда процесс LGWR получает запрос на переключение журнальных файлов;
- когда буфер повторения заполняется на одну треть или объем содержащихся в нем кешированных данных журнала повторения составляет 1 Мбайт.

По этим причинам лишь очень редкая система выиграет от буфера повторения с размером в пару десятков мегабайт. Крупная система с множеством параллельных транзакций может извлечь некоторое преимущество из большого буфера журнала повторения, поскольку пока процесс LCWR (отвечающий за сброс на диск буфера журнала повторения) записывает порцию журнального буфера, другие сеансы смогут его заполнять. В общем случае длительно выполняющаяся транзакция, которая генерирует много данных повторения, получит максимальный выигрыш от использования большого, а не обычного журнального буфера, т.к. она будет постоянно заполнять часть буфера журнала повторения пока процесс LCWR занят записью на диск другой части буфера (мы детально раскроем феномен записывания незафиксированных данных в главе 9). Чем крупнее и дольше длится транзакция, тем больше преимуществ она может извлечь из большого журнального буфера.

Стандартный размер буфера повторения, управляемый параметром LOG_BUFFER, варьируется в широких пределах в зависимости от операционной системы, версии базы данных и настроек других параметров. Вместо попытки объяснить, каким должен быть наиболее распространенный стандартный размер (его просто нет), я отсылаю вас к документации по применяемой версии Oracle (руководство *Oracle Database Reference*). Мое стандартное значение LOG_BUFFER, учитывая ранее запущенный экземпляр с SGA размером 1,5 Гбайт, отображается с помощью следующего запроса:

```

ODA@ORA12CR1> select value, isdefault
2  from v$parameter
3  where name = 'log_buffer'
4  /

```

| VALUE | ISDEFAULT |
|---------|-----------|
| 7036928 | TRUE |

Размер составляет около 7 Мбайт. Минимальный размер стандартного журнального буфера зависит от операционной системы. Если вы хотите выяснить, каков он, просто установите LOG_BUFFER в 1 байт и перезапустите базу данных. Например, в своем экземпляре Oracle Linux я получаю такой вывод:

```
EODA@ORA12CR1> alter system set log_buffer=1 scope=spfile;
System altered.
Система изменена.
```

```
EODA@ORA12CR1> connect / as sysdba;
Connected.
Подключено.
```

```
SYS@ORA12CR1> startup force;
ORACLE instance started.
Экземпляр ORACLE запущен.
```

```
Total System Global Area          1586708480 bytes
Fixed Size                        2288824 bytes
Variable Size                     402654024 bytes
Database Buffers                  1174405120 bytes
Redo Buffers                      7360512 bytes
```

Database mounted.

Database opened.

```
SYS@ORA12CR1> show parameter log_buffer
```

| NAME | TYPE | VALUE |
|------------|---------|---------|
| log_buffer | integer | 1703936 |

```
SYS@ORA12CR1> select 1703936/1024/1024 from dual;
```

```
1703936/1024/1024
-----
1.625
```

Наименьшим размером буфера, который я действительно могу иметь, независимо от установок в этой системе будет 1,625 Мбайт.

Кеш буферов блоков

До сих пор мы рассматривали сравнительно небольшие компоненты области SGA. Теперь нам предстоит взглянуть на компонент с потенциально очень большим размером. В кеше буферов блоков Oracle хранит блоки базы данных перед их записью на диск и после их чтения с диска. Для нас это критически важная область SGA. Если сделать ее слишком малой, то запросы будут выполняться бесконечно долго. Если же сделать ее слишком большой, то это “посадит на голодный паек” остальные процессы (т.е. выделенному серверу не останется достаточного пространства памяти для создания его области PGA, и мы даже не сможем запустить базу данных).

В ранних выпусках Oracle существовал только один кеш буферов блоков, и все блоки из любых сегментов помещались в эту единственную область. Начиная с Oracle 8.0, для хранения кешированных блоков из различных сегментов в SGA можно использовать три местоположения.

- **Стандартный пул (default pool).** Место, где обычно кешируются блоки всех сегментов. Это исходный — и ранее единственный — буферный пул.

- **Удерживающий пул (keep pool).** Запасной буферный пул, в котором по соглашению будут храниться часто используемые сегменты, перемещенные из стандартного буферного пула в связи с необходимостью освобождения места для других сегментов.
- **Рециклирующий пул (recycle pool).** Запасной буферный пул, где по соглашению будут храниться большие, очень редко используемые сегменты, которые приводят к длительному сбросу данных на диск, но не предоставляют никаких особых преимуществ, поскольку к тому моменту, когда блок потребуется снова, он уже будет удален из кеша по причине устаревания. Эти сегменты можно выделять из сегментов стандартного и удерживающего пулов для предотвращения их полного удаления из кеша по истечении заданного времени хранения.

Обратите внимание, что при описании удерживающего и рециклирующего пулов применялось выражение “по соглашению”. В действительности никто не принуждает использовать эти пулы только так, как описано. Фактически все три пула управляют блоками почти идентичным образом; *применяемые ими алгоритмы удаления вследствие устаревания или кеширования блоков не имеют существенных отличий.* Цель использования этих пулов — предоставление администратору базы данных возможности разделения сегментов на горячие, теплые и “не заслуживающие кеширования” области. Теоретически объекты в стандартном пуле должны были быть достаточно горячими (т.е. применяться довольно часто), чтобы гарантировано оставаться в кеше при любых обстоятельствах. Кеш будет удерживать такие объекты в памяти из-за их очень высокой популярности. Некоторые сегменты могут быть достаточно популярными, но не по-настоящему горячими, тогда эти блоки считаются теплыми. Они могут сбрасываться из кеша на диск с целью освобождения места для блоков, которые применяются не настолько часто (“не заслуживающие кеширования” блоки). Для сохранения таких теплых блоков в кеше можно выполнить одно из следующих действий.

- Назначить эти сегменты удерживающему пулу, чтобы позволить теплым блокам оставаться в буферном кеше подольше.
- Назначить “не заслуживающие кеширования” сегменты рециклирующему пулу, сохраняя размер этого пула достаточно маленьким, чтобы операции помещения блоков в кеш и их удаления из него выполнялись чаще (это позволит снизить накладные расходы, связанные с управлением всеми блоками).

Любое из этих двух действий увеличивает объем работы, которую должен выполнять администратор базы данных, поскольку ему придется думать о трех кешах, их размерах и назначенных им объектах. Следует также помнить, что между этими областями отсутствует механизм обмена данными, поэтому, если удерживающий пул содержит большой объем неиспользуемого пространства, то он не сможет передать его перегруженному стандартному или рециклирующему пулу. В общем, эти пулы представляют собой очень эффективные средства низкоуровневой настройки, к которым, однако, следует прибегать после применения всех других средств (если запрос можно переписать так, чтобы он выполнял только одну десятую часть операций ввода-вывода, я предпочту эту возможность вместо настройки множества буферных пулов).

Начиная с Oracle9i, в дополнение к стандартному, удерживающему и рециклирующему пулам администраторы баз данных получили в свое распоряжение еще четыре необязательных кеша, DB_nK_CACHE_SIZE. Эти кеши были добавлены для поддержки в базе данных нескольких размеров блоков. До выхода версии Oracle9i база данных должна была иметь единственный размер блока (обычно 2, 4, 8, 16 или 32 Кбайт). Начиная с Oracle9i, база данных может располагать стандартным размером блока, который задает размер блоков, хранящихся в стандартном, удерживающем и рециклирующем пулах, а также до четырех нестандартных размеров блоков, как объяснялось в главе 3.

Управление блоками в этих буферных кешах производится тем же способом, что и управление блоками в исходном стандартном пуле — используемые алгоритмы не претерпели никаких изменений. Давайте теперь посмотрим, как осуществляется управление блоками в этих пулах.

Управление блоками в кеше буферов

Для простоты предположим, что в системе существует только стандартный пул. Поскольку управление остальными пулами выполняется точно так же, достаточно обсудить только один из них.

По существу управление блокам в кеше буферов производится из единственной области с помощью двух различных указывающих на них списков:

- список *грязных* (dirty) блоков, которые должны записываться процессом записи блоков базы данных (DBWn; мы рассмотрим этот процесс несколько позже);
- список *чистых* (nondirty) блоков.

В Oracle 8.0 и предшествующих версиях список чистых блоков назывался списком недавно использовавшихся блоков (Least Recently Used — LRU). В нем блоки перечислялись в порядке их использования. В Oracle8i и последующих версиях алгоритм управления этими блоками был несколько модифицирован. Вместо того чтобы поддерживать список блоков, упорядоченный в определенном физическом порядке, Oracle задействует алгоритм со счетчиком обращений, который инкрементирует значение счетчика, связанного с блоком, при обнаружении данного блока в кеше. Значение счетчика увеличивается не при каждом попадании блока, а приблизительно каждые три секунды, если блок встречается постоянно. Увидеть этот алгоритм в работе можно на примере одного из действительно загадочных наборов таблиц — таблиц X\$. Эти таблицы абсолютно не документированы Oracle, но время от времени информация о них просачивается во внешний мир.

На заметку! В последующих примерах я имею дело с пользователем, подключенным с привилегиями SYSDBA, потому что таблицы X\$ по умолчанию видимы только этой учетной записи. На практике вы не должны применять учетную запись с привилегиями SYSDBA для выполнения запросов. Необходимость в запрашивании информации о блоках в кеше буферов — редкое исключение из этого правила.

Таблица X\$BH отображает информацию о блоках в кеше буферов блоков (она предоставляет больше информации, чем документированное представление V\$BH). Здесь можно наблюдать инкрементирование значения счетчика обращений при обнаружениях блоков.

Для отыскания пяти “самых горячих в настоящий момент блоков” и соединения этой информации с представлением DBA_OBJECTS для выяснения сегментов, которым они принадлежат, можно выполнить приведенный ниже запрос к упомянутому представлению. Запрос упорядочивает строки в таблице X\$BH по столбцу TCH (touch count — счетчик обращений) и сохраняет первые пять строк. Затем информация X\$BH соединяется с представлением DBA_OBJECTS по соответствию столбца X\$BH.OBJ с DBA_OBJECTS.DATA_OBJECT_ID.

```
SYS@ORA12CR1> select tch, file#, dbablk,
2      case when obj = 4294967295
3      then 'rbs/compat segment'
4      else (select max( '||object_type||' ) ||
5      owner || '.' || object_name ) ||
6      decode( count(*), 1, ' ', ' maybe!' )
7      from dba_objects
8      where data_object_id = X.OBJ )
9      end what
10 from (
11 select tch, file#, dbablk, obj
12 from x$bh
13 where state <> 0
14 order by tch desc
15 ) x
16 where rownum <= 5
17 /
```

| TCH | FILE# | DBABLK | WHAT |
|-----|-------|--------|---------------------------|
| 98 | 1 | 2825 | (INDEX) SYS.I_JOB_NEXT |
| 13 | 1 | 337 | (INDEX) SYS.I_OBJ1 |
| 13 | 1 | 62117 | (INDEX) SYS.I_OBJ1 |
| 11 | 1 | 4377 | (INDEX) SYS.SYS_C00819 |
| 11 | 1 | 209 | (TABLE) SYS.USER\$ maybe! |

На заметку! Выражение $(2^{32} - 1)$, или 4 294 967 295, в операторе case — это “магическое” число, используемое для пометки “специальных” блоков. Если необходимо понять, с чем связан лежащий в основе блок, воспользуйтесь запросом `select * from dba_extents where file_id = <FILE#> and block_id <= <DBABLK> and block_id+blocks-1 >= <DBABLK>`.

У вас может возникнуть вопрос о том, что означает ' maybe! ' и зачем применяется функция MAX() в предыдущем скалярном подзапросе. Причина в том, что, как видно в следующем запросе, столбец DATA_OBJECT_ID не является “первичным ключом” в представлении DBA_OBJECTS:

```
SYS@ORA12CR1> select data_object_id, count(*)
2 from dba_objects
3 where data_object_id is not null
4 group by data_object_id
5 having count(*) > 1;
```

| DATA_OBJECT_ID | COUNT (*) |
|----------------|-----------|
| 337 | 2 |
| 6 | 3 |
| 29 | 3 |
| 620 | 7 |
| 2 | 18 |
| 781 | 3 |
| 8 | 3 |
| 750 | 3 |
| 64 | 2 |
| 10 | 3 |

10 rows selected.

10 строк выбрано.

Это обусловлено наличием кластеров (обсуждаются в главе 10), которые могут содержать несколько таблиц. Поэтому при выполнении соединения таблицы X\$BH с представлением DBA_OBJECTS для вывода имени сегмента формально необходимо перечислить имена всех объектов в кластере, т.к. блок базы данных не принадлежит постоянно какой-то одной таблице. Можно даже наблюдать, каким образом Oracle инкрементирует значения счетчика обращений к блоку, который запрашивается многократно. В этом примере мы будем использовать “магическую” таблицу DUAL — нам известно, что она содержит одну строку и один столбец.

На заметку! До выхода Oracle 10g запрос таблицы DUAL привел бы к полному сканированию реальной таблицы по имени DUAL, хранящейся в словаре данных. При включенной функции автоматической трассировки и выполнении запроса `SELECT DUMMY FROM DUAL` во всех выпусках Oracle вы наблюдали бы наличие ряда операций ввода-вывода (согласованных операций чтения). В Oracle 9i и предшествующих версиях в случае выдачи запроса `SELECT SYSDATE FROM DUAL` или переменная `:= SYSDATE` в среде PL/SQL также будут происходить операции реального ввода-вывода. Тем не менее, в Oracle 10g этот запрос `SELECT SYSDATE` распознается как не нуждающийся в действительном обращении к таблице DUAL (из-за того, что вы не запрашиваете из DUAL столбец или идентификатор строки) и выполняется подобно вызову функции. Следовательно, таблица DUAL не подвергается полному сканированию — приложению возвращается только значение SYSDATE. Это небольшое изменение может радикально уменьшить количество согласованных чтений в системе, которая интенсивно работает с таблицей DUAL.

Таким образом, при каждом выполнении следующего запроса мы должны обращаться к реальной таблице DUAL (поскольку явно ссылаемся на столбец DUMMY):

```

SYS@ORA12CR1> select tch, file#, dbablk, DUMMY
  2   from x$bh, (select dummy from dual)
  3   where obj = (select data_object_id
  4                 from dba_objects
  5                 where object_name = 'DUAL'
  6                 and data_object_id is not null)
  7 /

```

| TCH | FILE# | DBABLK | D |
|-----|-------|--------|---|
| 1 | 1 | 929 | X |
| 2 | 1 | 928 | X |

```
SYS@ORA12CR1> exec dbms_lock.sleep(3.2);
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно выполнена.
SYS@ORA12CR1> /
```

| TCH | FILE# | DBABLK | D |
|-----|-------|--------|---|
| 2 | 1 | 1416 | X |
| 2 | 1 | 1417 | X |

```
SYS@ORA12CR1> exec dbms_lock.sleep(3.2);
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно выполнена.
SYS@ORA12CR1> /
```

| TCH | FILE# | DBABLK | D |
|-----|-------|--------|---|
| 4 | 1 | 1416 | X |
| 4 | 1 | 1417 | X |

```
SYS@ORA12CR1> exec dbms_lock.sleep(3.2);
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно выполнена.
SYS@ORA12CR1> /
```

| TCH | FILE# | DBABLK | D |
|-----|-------|--------|---|
| 5 | 1 | 1416 | X |
| 5 | 1 | 1417 | X |

Я ожидаю, что вывод будет варьироваться в зависимости от выпуска Oracle; вы вполне можете наблюдать возвращение более двух строк. Может оказаться, что значение счетчика TCH инкрементируется не каждый раз. В многопользовательской системе результаты будут еще более непредсказуемыми. СУБД Oracle будет пытаться инкрементировать значение TCH каждые три секунды (время последнего обновления столбца TCH отображается в столбце TIM), но точное значение не столь уж важно, поскольку оно достаточно близко к действительному. Кроме того, через определенные интервалы времени Oracle будет намеренно “охлаждать” блоки и декрементировать значение счетчика TCH. Поэтому, если вы запускаете этот запрос в своей системе, то будьте готовы к получению отличающихся результатов.

Итак, в Oracle8i и последующих версиях буфер блоков больше не перемещается в голову списка, как было ранее; вместо этого он остается в списке на своем месте и имеет увеличивающийся счетчик обращений. Однако блоки будут иметь тенденцию к “перемещению” по списку с течением времени. Слово “перемещение” приведено в кавычках, потому что физически блок не перемещается. В действительности поддерживаются несколько списков, которые указывают на блоки, а блок будет “перемещаться” из списка в список. Например, измененные блоки задаются грязным списком (блоков, которые должны быть записаны на диск процессом DBWn). Кроме того, по мере их повторного использования с течением времени, когда буферный кеш оказывается практически заполненным и какой-то блок с низким значением счетчика обращений освобождается, такой блок будет “помещен” приблизительно в середину списка с новым блоком данных.

Полный алгоритм, применяемый для управления этими списками, достаточно сложен и слегка изменяется от выпуска к выпуску Oracle по мере внесения усовершенствований. Подробности реализации не имеют особого значения для разработчиков; достаточно помнить о том, что интенсивно используемые блоки будут кешироваться, а блоки, которые применяются редко, не будут кешироваться на длительный период.

На заметку! Если вы следите за изложением, запуская все примеры в своей базе данных, воспользуйтесь возможностью выйти от имени учетной записи SYSDBA и затем войти под своей учетной записью!

Использование нескольких размеров блоков

Начиная с Oracle9i, в одной базе данных можно иметь несколько размеров блоков. До этого все блоки в одной базе данных обладали одним и тем же размером, и для изменения размера блоков нужно было перестраивать всю базу данных. Теперь в базе данных можно одновременно применять “стандартный” размер блоков (размер блоков, указанный при первоначальном создании базы данных; он используется для табличных пространств SYSTEM и TEMPORARY) и до четырех других размеров блоков. Каждому уникальному размеру блоков должна быть назначена собственная область буферного кеша. Стандартный, удерживающий и рециклирующий пулы будут кешировать только блоки стандартного размера. Чтобы применять в базе данных нестандартные блоки, необходимо соответствующим образом сконфигурировать буферный пул.

В следующем примере размер стандартных блоков составляет 8 Кбайт. Я попытаюсь создать табличное пространство с размером блоков 16 Кбайт:

```

EODA@ORA12CR1> create tablespace ts_16k
  2 datafile '/tmp/ts_16k.dbf'
  3 size 5m
  4 blocksize 16k;
create tablespace ts_16k
*
ERROR at line 1:
ORA-29339: tablespace block size 16384 does not match configured block sizes
ОШИБКА в строке 1:
ORA-29339: размер блоков табличного пространства, равный 16384,
не соответствует сконфигурированным размерам блоков
EODA@ORA12CR1> show parameter 16k

```

| NAME | TYPE | VALUE |
|-------------------|-------------|-------|
| db_16k_cache_size | big integer | 0 |

В настоящий момент из-за того, что не был сконфигурирован кеш 16 Кбайт, такое табличное пространство создать нельзя. Для исправления ситуации можно было бы предпринять одно из следующих действий. Можно установить параметр DB_16K_CACHE_SIZE и перезапустить базу данных. Можно сжать один из других компонентов SGA, чтобы освободить место в существующей области SGA под кеш 16 Кбайт. Или же, если бы значение параметра SGA_MAX_SIZE было больше текущего размера SGA, можно было бы просто выделить кеш 16 Кбайт.

На заметку! Начиная с версии Oracle9i, размеры разнообразных компонентов SGA можно изменять во время работы базы данных. Если нужна возможность “увеличения” размера SGA свыше начального распределения, понадобится установить для SGA_MAX_SIZE значение, которое больше выделенного объема SGA. Например, если после запуска размер SGA составлял 800 Мбайт и требуется добавить к буферному кешу дополнительные 200 Мбайт, то придется установить SGA_MAX_SIZE в 1 Гбайт или большее значение, чтобы разрешить подобное разрастание.

В этом примере я устанавливаю DB_16K_CACHE_SIZE и перезапускаю экземпляр, поскольку использую автоматическое управление разделяемой памятью и не хочу вручную настраивать другие кешы:

```

EODA@ORA12CR1> alter system set sga_target=300m scope=spfile;
System altered.
Система изменена.

EODA@ORA12CR1> alter system set db_16k_cache_size = 16m scope=spfile;
System altered.
Система изменена.

EODA@ORA12CR1> connect / as sysdba
Connected.
Подключено.

SYS@ORA12CR1> startup force
ORACLE instance started.

ORACLE instance started.
Экземпляр ORACLE запущен.

Total System Global Area          313159680 bytes
Fixed Size                        2287864 bytes
Variable Size                     180356872 bytes
Database Buffers                  125829120 bytes
Redo Buffers                      4685824 bytes
Database mounted.
Database opened.
SYS@ORA12CR1> show parameter 16k

```

| NAME | TYPE | VALUE |
|-------------------|-------------|-------|
| db_16k_cache_size | big integer | 16M |

Итак, я определил еще один буферный кеш, предназначенный для кеширования блоков размером 16 Кбайт. Стандартный пул будет потреблять остальное пространство буферного кеша, что можно увидеть, запросив V\$SGASTAT. Эти два буферных кеша являются взаимно исключаящими: если один из них “заполняется”, он не может использовать пространство в другом. Это предоставляет администратору базы данных очень высокий уровень контроля над использованием памяти, но за счет увеличения сложности и трудоемкости управления. Функциональность множества размеров блоков не задумывалась как средство повышения производительности или настройки (если вам нужно несколько кешей, то у вас уже есть стандартный, удерживающий и рециклирующий пулы). Вместо этого она направлена на поддержку переносимых табличных пространств — возможности изъятия файлов форматированных данных из одной базы данных и их перемещение либо присоединение

к другой базе данных. Средство было реализовано для того, чтобы можно было извлечь файлы данных из транзакционной системы, в которой применяются блоки размером 8 Кбайт, и перенести эту информацию в хранилище данных с блоками размером 16 или 32 Кбайт.

Тем не менее, поддержка нескольких размеров блоков очень удобна при проверке теоретических предположений. Если нужно посмотреть, как база данных будет работать с другим размером блока (например, какой объем памяти будет занимать определенная таблица, если вместо блоков 8 Кбайт использовать блоки 4 Кбайт), то это легко проверить, не создавая экземпляр полностью новой базы данных.

Несколько размеров блоков можно применять также в качестве инструмента очень точной настройки определенного набора сегментов, предоставляя им собственные закрытые буферные пулы. Или же в смешанной системе пользователи, выполняющие транзакции, могли бы применять один набор данных, а пользователи учетных/складских данных — другой. Транзакционные данные получают преимущество от блоков меньшего размера из-за меньшего соперничества за блоки (меньше данных/строк на блок означает, что в общем случае меньше людей будут одновременно обращаться к одному и тому же блоку), а также из-за лучшей утилизации буферного кеша (пользователи считывают в кеш только интересующие их данные — одиночную строку или же небольшой набор строк). Учетные/складские данные, которые могут быть основаны на транзакционных данных, получают выигрыш от блоков большего размера частично из-за меньших накладных расходов, связанных с блоками (в целом для их хранения будет требоваться меньше пространства), и, возможно, из-за больших размеров логических блоков ввода-вывода. И поскольку в случае учетных/складских данных соперничество за обновление не возникает, то наличие большего количества строк в блоке — не проблема, а преимущество. Более того, пользователи, выполняющие транзакции, в действительности получают в свое распоряжение собственный буферный кеш; им не придется беспокоиться о том, что запросы на генерацию отчетов переполняют их кеш.

Но в общем случае стандартного, удерживающего и рециклирующего пулов должно быть вполне достаточно для оптимальной настройки кеша буфера блоков, а множество размеров блоков следует применять главным образом для переноса данных из одной базы данных в другую и, возможно, в смешанных системах с генерацией отчетов и выполнением транзакций.

Разделяемый пул

Разделяемый пул (shared pool) — одна из наиболее важных частей области SGA, особенно с точки зрения производительности и масштабируемости. Слишком маленький разделяемый пул может привести к такому снижению производительности, что система будет казаться “зависшей”. Слишком большой разделяемый пул может дать аналогичный эффект. Некорректное использование этого пула также может вызвать аварийную ситуацию.

Что же собой представляет разделяемый пул? В этой области Oracle кеширует множество “программных” данных. При выполнении разбора запроса здесь же кешируется проанализированное представление. Прежде чем приступить к разбору всего запроса, Oracle осуществляет поиск в разделяемом пуле, чтобы выяснить, не была ли уже сделана эта работа. Запущенный на выполнение код PL/SQL также ке-

шируется в разделяемом пуле, чтобы при следующем запуске Oracle не пришлось снова считывать его с диска. При наличии 1000 сеансов, выполняющих один и тот же код, только одна копия кода загружается и совместно используется всеми сеансами. В разделяемом пуле Oracle сохраняет системные параметры. Здесь же находится кеш словаря (кешированная информация об объектах базы данных). Короче говоря, в разделяемом пуле хранится буквально все кроме разве что кухонной утвари.

Разделяемый пул характеризуется множеством небольших (обычно 4 Кбайт или меньше) сегментов памяти. Имейте в виду, что 4 Кбайт не является жестким пределом. В ряде случаев размеры сегментов могут быть больше, но основная цель заключается в применении небольших сегментов памяти для предотвращения фрагментации, которая возникала бы, если бы память выделялась сегментами, размеры которых сильно отличались друг от друга — от очень маленьких до очень больших. Управление памятью в разделяемом пуле осуществляется на основе алгоритма LRU. В этом отношении разделяемый пул похож на буферный кеш — если он не используется, то теряется. Для изменения такого поведения можно применять дополнительный пакет DBMS_SHARED_POOL, который позволяет принудительно закреплять объекты в разделяемом пуле. Эту процедуру можно задействовать для загрузки часто используемых процедур и пакетов во время начального запуска базы данных, обеспечив при этом невозможность их устаревания. Однако в обычных условиях, если какое-то время часть содержимого памяти разделяемого пула не используется повторно, то это содержимое устареет. Даже код PL/SQL, который может быть довольно большим, управляется механизмом подкачки страниц, так что при выполнении кода очень большого пакета в разделяемый пул загружается лишь действительно необходимый код, причем небольшими порциями. Если с ним не работать в течение продолжительного периода времени, то этот код устареет и будет удален из разделяемого пула в случае его заполнения и необходимости освобождения места для других объектов.

Простейший способ разрушения разделяемого пула Oracle — отказ от применения переменных привязки. Как было показано в главе 1, отказ от переменных привязки может “поставить систему на колени” по двум причинам:

- система тратит непомерно большую часть времени процессора на разбор запросов;
- система использует слишком много ресурсов для управления объектами в разделяемом пуле как результат того, что запросы никогда не используются повторно.

Если бы каждый отправленный базе данных Oracle запрос был уникальным (из-за того, что в нем жестко закодированы уникальные значения), то концепция разделяемого пула была бы большей частью ликвидирована. Разделяемый пул разрабатывался с расчетом на многократное использование планов запросов. Если каждый запрос является совершенно новым, никогда ранее не встречавшимся запросом, то кеширование ведет только к увеличению накладных расходов. Разделяемый пул становится компонентом, который подавляет производительность. Часто, но совершенно неоправданно, для решения этой проблемы многие пытаются увеличивать размер разделяемого пула, что обычно только ухудшает положение. Когда разделяемый пул неизбежно снова заполняется, он становится еще большей обузой, чем пул

меньшего размера, по той простой причине, что управление большим заполненным разделяемым пулом будет более трудоемким, нежели управление маленьким полным пулом.

Единственно правильным решением этой проблемы является применение разделяемого кода SQL для повторного использования запросов. Ранее в главе 1 мы бегло рассмотрели параметр `CURSOR_SHARING`, который может служить кратковременным средством решения данной проблемы. Однако единственный реальный способ ее решения — применение SQL-кода многократного использования. Даже в самых больших системах применяется максимум 10 000–20 000 уникальных SQL-операторов. В большинстве систем выполняются только несколько сотен уникальных запросов.

Следующий реальный пример демонстрирует, к каким негативным последствиям может привести некорректное использование разделяемого пула. Меня попросили поработать над системой, в которой стандартной рабочей процедурой была остановка базы данных на ночь, очистка области SGA и повторный ее запуск в чистом состоянии. Причина выполнения этих действий заключалась в том, что в течение рабочего дня в системе возникали проблемы, связанные с перегрузкой процессора, и если база данных продолжала работать более одного дня, ее производительность действительно начинала деградировать. В системе применялся разделяемый пул в 1 Гбайт внутри области SGA размером 1,1 Гбайт. Это было действительно так: 0,1 Гбайт было выделено для буферного кеша блоков и других элементов, а 1 Гбайт был предназначен для кеширования уникальных запросов, которые никогда более не выдавались снова. Причина холодного старта обуславливалась тем, что в случае продолжения работы системы дольше одного дня объем свободной памяти в разделяемом пуле полностью исчерпывался. К этому моменту накладные расходы, связанные с наличием устаревших структур (особенно таких крупных), становились настолько большими, что система переставала с ними справляться, и производительность стремительно снижалась (хотя она не была высокой в любом случае, поскольку системе приходилось управлять разделяемым пулом в 1 Гбайт). Кроме того, в связи с тем, что полный разбор SQL-запросов требовал большой загрузки процессора, персонал, работающий с этой системой, постоянно стремился наращивать количество процессоров компьютера. Исправление приложения и разрешение ему применять переменные привязки не только позволило снизить предъявляемые к компьютеру аппаратные требования (компьютер и без того обладал значительно более мощным процессором, чем действительно было необходимо), но и высвободить память, выделенную для различных пулов. В результате вместо 1 Гбайт разделяемому пулу были выделены 100 Мбайт, которые никогда не заполнялись полностью в течение многих недель непрерывной работы базы данных.

Приведу последний комментарий по поводу разделяемого пула и параметра `SHARED_POOL_SIZE`. В Oracle9i и предшествующих версиях не существовало никакой прямой взаимосвязи между результатом выполнения следующего запроса:

```
ops$tkyte@ORA9IR2> select sum(bytes) from v$sgastat where pool = 'shared pool';
SUM(BYTES)
-----
100663296
```

и параметром `SHARED_POOL_SIZE`:

```
ops$tkyte@ORA9IR2> show parameter shared_pool_size
```

| NAME | TYPE | VALUE |
|------------------|-------------|----------|
| shared_pool_size | big integer | 83886080 |

кроме того факта, что значение `SUM(BYTES) FROM V$SGASTAT` всегда было больше значения `SHARED_POOL_SIZE`. Разделяемый пул содержит множество других структур, которые выходят за рамки, определяемые этим параметром. Как правило, `SHARED_POOL_SIZE` вносит наибольший вклад в размер разделяемого пула, сообщаемый столбцом `SUM(BYTES)`, но он не является единственной составляющей. Например, параметр `CONTROL_FILES` занимает 264 байта для каждого файла в разделе “miscellaneous” (“прочие”) разделяемого пула. Выбор имени “shared pool” в представлении `V$SGASTAT` и имени параметра `SHARED_POOL_SIZE` не особенно удачен, поскольку этот параметр вносит наибольший вклад в размер разделяемого пула, но он не единственный, кто вносит свой вклад.

Однако в Oracle 10g и последующих версиях наблюдается полное соответствие при условии использования ручного управления памятью SGA (т.е. при установке значения параметра `SHARED_POOL_SIZE` администратором базы данных):

```
ops$tkyte@ORA10G> select sum(bytes)/1024/1024 mbytes from v$sgastat where
pool = 'shared pool';
```

```
      MBYTES
-----
      128
```

```
ops$tkyte@ORA10G> show parameter shared_pool_size;
```

| NAME | TYPE | VALUE |
|------------------|-------------|-------|
| shared_pool_size | big integer | 128M |

На заметку! В этом примере использовалось ручное управление разделяемой памятью!

Это достаточно важное изменение по сравнению с Oracle9i и версиями до Oracle 10g. В Oracle 10g параметр `SHARED_POOL_SIZE` управляет размером разделяемого пула, в то время как в Oracle9i и предшествующих версиях он лишь вносил наибольший вклад в размер разделяемого пула. Рекомендую выяснить действительный размер разделяемого пула в Oracle9i и предшествующих версиях (на основе `VSSGASTAT`) и воспользоваться этой информацией для установки значения параметра `SHARED_POOL_SIZE` в Oracle 10g и последующих версиях. Теперь предполагается, что память для различных дополнительных компонентов, вносящих свой вклад в размер разделяемого пула, выделяется администратором базы данных.

Большой пул

Большой пул (large pool) назван так не потому, что является “большой” структурой (хотя его размер вполне может быть достаточно большим), а из-за того, что он применяется для выделения крупных сегментов памяти, превышающих по размеру те, для обработки которых предназначен разделяемый пул.

До появления большого пула в Oracle 8.0 вся память выделялась в разделяемом пуле. Это было нежелательно при использовании функций, которые требовали “крупные” сегменты памяти, такие как сегменты памяти UGA разделяемого сервера. Проблема еще больше усложнялась тем, что обработка, которая обычно требовала очень большого объема выделенной памяти, должна была применять память иначе, чем разделяемый пул управлял ею. Разделяемый пул управляет памятью по алгоритму LRU, который прекрасно подходит для кеширования и повторного использования данных. Однако выделение больших сегментов обычно предполагает получение сегмента памяти, его применение и прекращение работы с ним — необходимость в кешировании содержимого этой памяти отсутствует.

СУБД Oracle нуждалась в механизме, подобном реализации рециклирующего и удерживающего пулов кеша буферов блоков. Именно такими являются большой и разделяемый пулы в настоящее время. Большой пул представляет собой область памяти, аналогичную рециклирующему пулу, а разделяемый пул больше похож на удерживающий буферный пул — если оказывается, что какой-то объект используется часто, он сохраняется в кеше.

Управление памятью, выделенной в большом пуле, осуществляется в куче во многом подобно тому, как язык C управляет памятью посредством процедур `malloc()` и `free()`. Как только вы “освобождаете” сегмент памяти, он может быть задействован другими процессами. Разделяемый пул не имел никаких механизмов освобождения сегментов памяти. Он должен был выделять ее, использовать, а затем прекратить потребление. По истечении некоторого времени, если возникала необходимость повторного использования этой области памяти, СУБД Oracle должна была удалить содержимое устаревшего сегмента. Проблема с применением одного только разделяемого пула состоит в том, что единственный размер подходит не для всех ситуаций.

Большой пул используется:

- *подключениями посредством разделяемого сервера* для выделения области UGA внутри области SGA;
- *параллельным выполнением операторов*, чтобы сделать возможным выделение буферов сообщений для взаимодействия между процессами, которые применяются для координации серверов параллельных запросов;
- *резервным копированием* для выделения буферов дискового ввода-вывода RMAN в некоторых случаях.

Как видите, ни одно из этих выделений памяти не должно управляться в буферном пуле типа LRU, предназначенном для управления небольшими сегментами памяти. Например, память, используемая подключением посредством разделяемого сервера, никогда не применяется повторно после выхода из сеанса, поэтому она должна немедленно возвращаться в пул. Кроме того, область памяти разделяемого сервера имеет тенденцию быть “крупной”. Если вы просмотрите приведенные ранее примеры работы с параметром `SORT_AREA_RETAINED_SIZE` или `PGA_AGGREGATE_TARGET`, то вспомните, что область UGA может становиться очень большой и определенно превышать порции в 4 Кбайт. Помещение памяти многопоточного сервера в разделяемый пул ведет к ее разбиению на фрагменты произвольных размеров. Более того, окажется, что большие сегменты памяти, которые никогда не будут применяться

повторно, приведут к устареванию и тех сегментов, которые могли бы использоваться многократно. В результате база данных будет вынуждена позже выполнять дополнительную работу по перестройке этой структуры памяти.

Все сказанное относится и к буферам сообщений параллельных запросов, поскольку они не пригодны для управления по алгоритму LRU. Они выделяются и не могут быть освобождены до тех пор, пока их применение не будет прекращено. Как только они доставили свои сообщения, они больше не требуются и должны быть незамедлительно освобождены. Это в еще большей степени относится к буферам резервного копирования — их размер очень велик, и как только Oracle прекращает их использование, они должны немедленно “исчезнуть”.

В случае применения подключений через разделяемый сервер наличие большого пула не обязательно, но настоятельно рекомендуется. Если большой пул отсутствует и производится подключение посредством разделяемого сервера, то память выделяется в разделяемом пуле, как это всегда происходило в Oracle 7.3 и предшествующих версиях. Это неизбежно приведет к снижению производительности по истечении некоторого времени, и такой ситуации следует избегать. Большой пул получит определенный стандартный размер, если параметр `DBWR_IO_SLAVES` или `PARALLEL_MAX_SERVERS` установлен в какое-то положительное значение. Когда применяется функциональное средство, которое задействует большой пул, его размер должен быть установлен вручную. Стандартное значение обычно не будет подходить в такой ситуации.

Пул Java

Пул Java появился в версии Oracle 8.1.5 для поддержки выполнения кода Java в базе данных. Если хранимая процедура написана на языке Java, то Oracle будет использовать этот сегмент памяти при обработке такого кода.

Параметр `JAVA_POOL_SIZE` служит для фиксации объема памяти, которая выделяется пулу Java для хранения всего кода и данных Java, специфичных для сеанса.

Способ применения пула Java зависит от режима функционирования сервера Oracle. В режиме выделенного сервера пул Java включает в себя разделяемые части каждого из классов Java, которые действительно используются в сеансе. В основном это части, предназначенные только для чтения (векторы выполнения, методы и т.п.), и их размер составляет около 4–8 Кбайт на класс. Таким образом, в режиме выделенного сервера общий объем памяти, требуемый для пула Java, является достаточно умеренным и может быть определен, исходя из количества применяемых классов Java. Следует отметить, что в режиме выделенного сервера ни один элемент данных, специфичный для сеанса, не хранится в области SGA, т.к. эта информация содержится в области UGA и, как вы помните, в режиме выделенного сервера область UGA входит в состав PGA. При подключении к Oracle через разделяемый сервер пул Java содержит оба описанных ниже компонента.

- Разделяемая часть каждого класса Java.
- Часть области UGA, используемая для хранения информации о состоянии каждого сеанса и выделяемая из сегмента `JAVA_POOL` внутри области SGA. Остальная часть области UGA будет как обычно размещаться в разделяемом или в большом пуле, если он сконфигурирован.

Поскольку в Oracle9i и предшествующих версиях общий размер пула Java фиксирован, разработчикам приложений придется оценить общий объем памяти, требуемый их приложениям, и умножить это значение на количество одновременно действующих сеансов, которые нужно поддерживать. Полученное значение определит общий размер пула Java. Каждый сегмент Java области UGA будет увеличиваться и уменьшаться по мере необходимости, но следует иметь в виду, что размер этого пула должен быть определен так, чтобы в нем могли одновременно уместиться все сегменты UGA. В Oracle 10g и последующих версиях этот параметр может изменяться, и пул Java может со временем увеличиваться и уменьшаться без перезапуска базы данных.

Пул Streams

Пул Streams (Потоки) — это структура SGA, появившаяся в Oracle 10g. В число продуктов Oracle, которые работают с пулом Streams, входят Oracle GoldenGate, XStream, Oracle Streams, Oracle Advanced Queuing и Oracle Data Pump.

Размер пула Streams регулируется установкой параметра `STREAMS_POOL_SIZE`. Если параметр `SGA_TARGET` имеет ненулевое значение, то автоматическое управление памятью SGA будет применять параметр `STREAMS_POOL_SIZE` в качестве минимального значения для размера пула Streams. Если параметры `SGA_TARGET` и `STREAMS_POOL_SIZE` оба установлены в 0, то пул Streams будет использовать до 10% объема разделяемого пула.

Продукты, работающие с пулом Streams, будут буферизировать сообщения в очередях. Вместо применения постоянных, основанных на диске очередей, с сопровождающими их накладными расходами, эти средства используют очереди внутри памяти. По мере заполнения очереди будут сбрасываться на диск. Если по какой-то причине (из-за программной ошибки, сбоя электропитания и т.п.) в экземпляре Oracle с очередями в памяти происходит отказ, то эти очереди в памяти воссоздаются из журналов повторения транзакций.

Пул Streams важен только в системах, использующих средства (такие как GoldenGate, Streams и т.д.), которым необходимо пространство в этой области памяти. В этих средах должен быть установлен параметр `STREAMS_POOL_SIZE`, чтобы избежать расхода 10% объема разделяемого пула.

Управление памятью SGA

Параметры, связанные с памятью SGA, относятся к одной из двух областей.

- **Автоматически настраиваемые параметры SGA.** В настоящее время это `DB_CACHE_SIZE`, `SHARED_POOL_SIZE`, `LARGE_POOL_SIZE`, `JAVA_POOL_SIZE` и `STREAMS_POOL_SIZE`.
- **Настраиваемые вручную параметры SGA.** К ним относятся `LOG_BUFFER`, `DB_NK_CACHE_SIZE`, `DB_KEEP_CACHE_SIZE` и `DB_RECYCLE_CACHE_SIZE`.

В Oracle 10g и последующих версиях в любой момент можно запросить представление `V$SGAINFO` и посмотреть, какие компоненты SGA поддерживают возможность изменения размера:

```
EOA@ORA12CR1> select * from V$SGAINFO;
```

| NAME | BYTES | RES | CON_ID |
|---------------------------------|-----------|-----|--------|
| Fixed SGA Size | 2287336 | No | 0 |
| Redo Buffers | 4890624 | No | 0 |
| Buffer Cache Size | 67108864 | Yes | 0 |
| Shared Pool Size | 184549376 | Yes | 0 |
| Large Pool Size | 4194304 | Yes | 0 |
| Java Pool Size | 4194304 | Yes | 0 |
| Streams Pool Size | 0 | Yes | 0 |
| Shared IO Pool Size | 4194304 | Yes | 0 |
| Data Transfer Cache Size | 0 | Yes | 0 |
| Granule Size | 4194304 | No | 0 |
| Maximum SGA Size | 267227136 | No | 0 |
| Startup overhead in Shared Pool | 169940696 | No | 0 |
| Free SGA Memory Available | 0 | | 0 |

```
13 rows selected.
```

```
13 строк выбрано.
```

Для компонентов памяти SGA, которые могут автоматически настраиваться, предусмотрены три способа управления ими.

- Ручное управление разделяемой памятью. Установка всех необходимых параметров пула и кеша.
- Автоматическое управление разделяемой памятью (или памятью SGA), доступное в Oracle 10g и последующих версиях. Установка параметра SGA_TARGET. За счет установки параметра SGA_TARGET вы позволяете экземпляру задавать и изменять размеры разнообразных компонентов SGA.
- Автоматическое управление памятью, доступное в Oracle 11g и последующих версиях. Установка параметра MEMORY_TARGET. За счет установки параметра MEMORY_TARGET вы позволяете экземпляру задавать и изменять размеры областей памяти SGA и PGA.

Ниже мы обсудим эти способы по очереди.

Ручное управление разделяемой памятью

Если требуется определенный уровень контроля над автоматически настраиваемыми областями памяти SGA, установите параметры MEMORY_TARGET и SGA_TARGET в нулевое значение. Когда в ноль установлен параметр MEMORY_TARGET, отключается автоматическое управление памятью, а когда в ноль установлен параметр SGA_TARGET, отключается автоматическое управление разделяемой памятью.

На заметку! В Oracle9i и предшествующих версиях было доступно только ручное управление разделяемой памятью — параметр SGA_TARGET не существовал, а параметр SGA_MAX_SIZE указывал максимальный размер области SGA.

Когда автоматическое управление памятью отключено, вы можете вручную устанавливать размер области SGA, указывая значения для следующих параметров памяти: DB_CACHE_SIZE, SHARED_POOL_SIZE, LARGE_POOL_SIZE, JAVA_POOL_SIZE и STREAMS_POOL_SIZE. Каждый из перечисленных параметров имеет стандартное значение, которое Oracle будет применять в случае, если параметр не установлен

явным образом. Например, параметр `DB_CACHE_SIZE` будет установлен в значение 48 Мбайт или 4 Мбайт, умноженное на количество процессоров, в зависимости от того, какая величина окажется больше в конкретной системе.

Ниже приведен пример содержимого файла инициализации, в котором включается ручное управление разделяемой памятью:

```
*.compatible='12.1.0.1'
*.control_files='/u01/dbfile/ORAl2CR1/control01.ctl',
                '/u02/dbfile/ORAl2CR1/control02.ctl'
*.db_block_size=8192
*.db_name='ORAl2CR1'
*.memory_target=0
*.sga_target=0
*.db_cache_size=1G
*.shared_pool_size=256M
*.pga_aggregate_target=256m
*.open_cursors=300
*.processes=600
*.remote_login_passwordfile='EXCLUSIVE'
*.resource_limit=TRUE
*.undo_tablespace='UNDOTBS1'
```

В Oracle 11g Release 2 и последующих версиях с ручным управлением разделяемой памятью связан один аспект, о котором следует знать: даже когда вы явно отключаете все автоматическое управление памятью (устанавливая в ноль параметры `MEMORY_TARGET` и `SGA_TARGET`), Oracle может по-прежнему делать ряд автоматических перераспределений памяти из кеша буферов базы данных в разделяемый пул. В случае исчерпания пространства в разделяемом пуле Oracle будет автоматически добавлять к нему пространство, чтобы избежать появления ошибки `ORA-04031 "Unable to allocate %s bytes of shared memory"` ("Не удастся выделить %s байтов разделяемой памяти").

Понаблюдать за автоматическим изменением размера памяти SGA можно, запросив представление `V$MEMORY_RESIZE_OPS`. Столбец `OPER_MODE` будет содержать значение `DEFERRED` или `IMMEDIATE` для любых операций автоматического изменения размера SGA. Когда используется ручное управление памятью SGA, автоматическое изменение размера SGA отключается для запросов в режиме `DEFERRED`, но разрешено для запросов в режиме `IMMEDIATE`. Следовательно, в случае применения ручного управления разделяемой памятью вы можете обнаружить наличие операций `GROW` и `SHRINK` (в столбце `OPER_TYPE`) для запросов автоматической настройки `IMMEDIATE` в отношении размеров областей памяти `DB_CACHE_SIZE` и `SHARED_POOL`.

Продemonстрируем утверждения из предыдущих абзацев на простом примере. Для начала создадим пакет `DBMS_SHARED_POOL`, чтобы можно было закрепить объекты в разделяемом пуле:

```
SYS@ORA12CR1> @?/rdbms/admin/dbmspool
Session altered.
Сеанс изменен.
Package created.
Пакет создан.
Grant succeeded.
Выдано успешно.
```

Session altered.

Сеанс изменен.

```
SYS@ORA12CR1> grant execute on dbms_shared_pool to eoda;
```

Grant succeeded.

Выдано успешно.

Далее запустим код, который быстро начнет заполнять разделяемый пул (закрепив множество процедур PL/SQL в разделяемом пуле):

```
SYS@ORA12CR1> conn eoda/foo
```

Connected.

Подключено.

```
EODA@ORA12CR1> declare
```

```
2   k varchar2(30);
```

```
3   ss varchar2(2000);
```

```
4   begin
```

```
5     for i in 1 .. 100000 loop
```

```
6       ss := 'create or replace procedure SP' || i || ' is
```

```
7         a number;
```

```
8         begin
```

```
9           a := 123456789012345678901234567890;
```

```
10          a := 123456789012345678901234567890;
```

```
11          a := 123456789012345678901234567890;
```

```
12          a := 123456789012345678901234567890;
```

```
13          a := 123456789012345678901234567890;
```

```
14          a := 123456789012345678901234567890;
```

```
15          a := 123456789012345678901234567890;
```

```
16          a := 123456789012345678901234567890;
```

```
17          a := 123456789012345678901234567890;
```

```
18          a := 123456789012345678901234567890;
```

```
19          a := 123456789012345678901234567890;
```

```
20          a := 123456789012345678901234567890;
```

```
21          a := 123456789012345678901234567890;
```

```
22          a := 123456789012345678901234567890;
```

```
23          a := 123456789012345678901234567890;
```

```
24        end;';
```

```
25      execute immediate ss;
```

```
26      k := 'SP' || i;
```

```
27      sys.dbms_shared_pool.keep(k);
```

```
28    end loop;
```

```
29  end;
```

```
30  /
```

А теперь из другого сеанса запросим словарь данных, чтобы увидеть операции изменения размера памяти, связанные с тем, что Oracle перемещает память в разделяемый пул:

```
EODA@ORA12CR1> select component, parameter, oper_type,
```

```
...           oper_mode from v$memory_resize_ops;
```

| | | | |
|----------------------|---------------|--------|-----------|
| DEFAULT buffer cache | db_cache_size | SHRINK | IMMEDIATE |
|----------------------|---------------|--------|-----------|

| | | | |
|-------------|------------------|------|-----------|
| shared pool | shared_pool_size | GROW | IMMEDIATE |
|-------------|------------------|------|-----------|

| | | | |
|----------------------|---------------|--------|-----------|
| DEFAULT buffer cache | db_cache_size | SHRINK | IMMEDIATE |
|----------------------|---------------|--------|-----------|

| | | | |
|----------------------|---------------|--------|-----------|
| DEFAULT buffer cache | db_cache_size | SHRINK | IMMEDIATE |
|----------------------|---------------|--------|-----------|

| | | | |
|----------------------|------------------|--------|-----------|
| shared pool | shared_pool_size | GROW | IMMEDIATE |
| DEFAULT buffer cache | db_cache_size | SHRINK | IMMEDIATE |
| shared pool | shared_pool_size | GROW | IMMEDIATE |

Итак, вы должны знать о возможности выполнения в Oracle 11g Release 2 и последующих версиях определенного автоматического перераспределения памяти SGA, даже когда данное средство было отключено. Как отмечалось ранее, это делается Oracle для того, чтобы процессы не исчерпали пространство в разделяемом пуле.

Совет. Дополнительные сведения об автоматическом изменении размера области SGA во время использования ручного управления разделяемой памятью можно найти в примечании MOS под номером 1269139.1.

Автоматическое управление разделяемой памятью

При автоматическом управлении разделяемой памятью основным параметром указания полного размера автоматически настраиваемых компонентов является SGA_TARGET, который может динамически изменяться во время работы базы данных вплоть до значения, установленного в SGA_MAX_SIZE. Его стандартное значение равно величине SGA_TARGET, поэтому если вы планируете увеличение SGA_TARGET, то перед запуском экземпляра базы данных должны установить SGA_MAX_SIZE в более высокое значение. База данных будет работать с памятью с объемом, равным значению SGA_TARGET минус значения размеров любых других компонентов с вручную заданными размерами, подобных DB_KEEP_CACHE_SIZE, DB_RECYCLE_CACHE_SIZE и т.д. Этот объем памяти будет использоваться для установки размеров стандартного, разделяемого, большого пулов и пула Java. Динамически во время выполнения экземпляр будет выделять и перераспределять память между этими четырьмя областями по мере необходимости. Вместо возвращения пользователю сообщения об ошибке ORA-04031 “Unable to allocate %s bytes of shared memory” (“Не удастся выделить %s байтов разделяемой памяти”) в случае нехватки памяти в разделяемом пуле экземпляр может уменьшить буферный кеш на определенное количество мегабайт (размер гранулы) и увеличить размер разделяемого пула на это число.

На заметку! Чтобы можно было применять автоматическое управление разделяемой памятью, параметр STATISTICS_LEVEL должен быть установлен в TYPICAL или ALL. Если сбор статистики не включен, база данных не будет располагать хронологической информацией, требуемой для принятия необходимых решений относительно установки размеров.

Со временем, когда потребности экземпляра в памяти прояснятся, размеры разнообразных компонентов SGA станут более или менее фиксированными. База данных запоминает также размеры этих четырех компонентов между запуском и остановом, поэтому выяснять их заново при каждом запуске экземпляра не придется. Это делается посредством четырех параметров, имена которых начинаются с двух символов подчеркивания: __db_cache_size, __java_pool_size, __large_pool_size и __shared_pool_size. Во время нормального или немедленного завершения работы база данных будет записывать эти параметры в *файл хранимых параметров*, а затем использовать их при запуске для установки стандартных размеров каждой области.

На заметку! Последнее средство сохранения рекомендованных значений для пулов работает, только если применяются файлы хранимых параметров (также известные как *файлы spfile*).

Кроме того, если известно определенное минимальное значение размера для одной из указанных пяти областей, этот параметр можно установить в дополнение к `SGA_TARGET`. Экземпляр будет использовать это значение в качестве нижней границы, т.е. наименьшего размера, который может иметь данная конкретная область.

Автоматическое управление памятью

В Oracle 11g Release 1 и последующих версиях база данных также предлагает *автоматическое управление памятью* — своего рода источник значений для всех настроек памяти. С появлением Oracle 10g и автоматического управления памятью SGA администратору базы данных оставалось устанавливать только две основные настройки памяти: `PGA_AGGREGATE_TARGET` и `SGA_TARGET`. База данных автоматически выделяла и перераспределяла порции памяти внутри каждой области, как было описано ранее. В Oracle 11g администратору базы данных необходимо установить всего один параметр памяти — `MEMORY_TARGET`, представляющий суммарный объем памяти, в пределах которого должны находиться объединенные области SGA и PGA (помните, что память PGA может быть в некоторой степени неконтролируемой). База данных динамически определит подходящие размеры областей SGA и PGA на основе хронологии рабочей нагрузки. По мере изменения рабочей нагрузки с течением времени размеры выделенных областей SGA и PGA также будут изменяться. Например, если в течение дня у вас производится интенсивная оперативная обработка транзакций (OLTP) и интенсивная пакетная обработка по ночам, вы можете обнаружить, что дневной размер SGA намного больше, чем PGA, а ночью все наоборот. Это отражает разные потребности в памяти указанных двух типов приложений.

На заметку! Перед внедрением автоматического управления памятью предусмотрите изучение любых примечаний MOS, специфичных для имеющейся операционной системы, таких как 749851.1 (для Linux) и 1399274.1 (для Solaris).

Как и при автоматическом управлении памятью SGA, администратор базы данных может определять нижние границы для размеров всех областей памяти, устанавливая параметры `SGA_TARGET` и `PGA_AGGREGATE_TARGET`, или нижнюю границу каждого пула в SGA, устанавливая в нее значения их параметров. База данных запомнит оптимальные настройки для пулов и областей SGA и PGA в файле хранимых параметров, если он применяется.

Например, в одной из своих тестовых систем я установил следующие параметры:

- `memory_target = 756m`
- `sga_target = 256m`
- `pga_aggregate_target = 256m`

Файл хранимых параметров для данной базы данных в настоящий момент выглядит так:

```

SYS@ORA12CR1> create pfile='/tmp/pfile' from spfile;
File created.
Файл создан.

SYS@ORA12CR1> !cat /tmp/pfile;
ORA12CR1.__data_transfer_cache_size=0
ORA12CR1.__db_cache_size=67108864
ORA12CR1.__java_pool_size=4194304
ORA12CR1.__large_pool_size=4194304
ORA12CR1.__oracle_base='/orahome/app/oracle'#ORACLE_BASE set from environment
ORA12CR1.__pga_aggregate_target=520093696
ORA12CR1.__sga_target=272629760
ORA12CR1.__shared_io_pool_size=0
ORA12CR1.__shared_pool_size=184549376
ORA12CR1.__streams_pool_size=0
*.compatible='12.1.0.1'
*.control_files='/u01/dbfile/ORA12CR1/control01.ctl',
                '/u02/dbfile/ORA12CR1/control02.ctl'
*.db_block_size=8192
*.db_name='ORA12CR1'
*.memory_target=792723456
*.open_cursors=300
*.pga_aggregate_target=268435456
*.processes=600
*.remote_login_passwordfile='EXCLUSIVE'
*.resource_limit=TRUE
*.sga_target=268435456
*.undo_tablespace='UNDOTBS1'

```

Как видите, параметры с именами, начинающимися с двух символов подчеркивания, которые выделены полужирным, теперь включают настройки **__sga_target** и **__pga_aggregate_target**, а также различные пулы. Эти значения выводятся на основе последних трех параметров, выделенных полужирным, и наблюдаемой рабочей нагрузки сервера. Таким образом, Oracle будет запоминать ваши последние оптимальные настройки SGA/PGA и использовать их при следующем перезапуске.

Резюме

В этой главе мы взглянули на структуры памяти Oracle. Мы начали с уровня процессов и сеансов, исследовав области PGA и UGA и взаимосвязь между ними. Мы выяснили, какое влияние на организацию памяти оказывает режим подключения к базе данных Oracle. Подключение через выделенный сервер предусматривает более интенсивное использование памяти процессом сервера, чем подключение посредством разделяемого сервера, но подключение через разделяемый сервер требует выделения значительно большего объема для области SGA. Затем мы обсудили основные структуры самой области SGA. Мы обнаружили отличия между разделяемым и большим пулами, и установили причины, по которым может требоваться большой пул для “экономии” объема разделяемого пула. Мы раскрыли особенности пула Java и способы его применения в разнообразных ситуациях, а также рассмотрели буферный кеш блоков и возможность его разделения на меньшие, более специализированные пулы. Теперь можно переходить к исследованию физических процессов, которые формируют остальную часть экземпляра Oracle.

Процессы Oracle

Наконец, мы добрались до последнего фрагмента архитектурной головоломки. Мы уже исследовали базу данных и набор образующих ее физических файлов. Рассмотрев память, используемую Oracle, мы раскрыли первую половину экземпляра. Осталось ознакомиться с еще одним вопросом, касающимся архитектуры — набором *процессов*, которые образуют вторую половину экземпляра.

Каждый процесс в Oracle будет выполнять отдельную задачу или набор задач, и каждый из них будет иметь внутреннюю память (память PGA), выделенную им для выполнения своего задания. С экземпляром Oracle связаны три обширных класса процессов.

- **Серверные процессы.** Эти процессы выполняют свою работу на основе клиентского запроса. Мы уже до определенной степени знакомы с выделенным и разделяемым серверами. Они являются серверными процессами.
- **Фоновые процессы.** Эти процессы запускаются вместе с базой данных и выполняют разнообразные задачи обслуживания, такие как запись блоков на диск, ведение оперативных журналов повторения, очистка после прерывания процессов, поддержка автоматического репозитория рабочей нагрузки (Automatic Workload Repository — AWR) и т.д.
- **Подчиненные процессы.** Эти процессы подобны фоновым, но выполняют дополнительную работу по поручению либо фонового, либо серверного процесса.

Мы уже упоминали некоторые из этих процессов, например, процесс записи блоков базы данных (DBWn) и процесс записи в журнал (LGWR), но здесь мы более подробно рассмотрим функционирование каждого процесса, выполняемые ими действия и причины такой их реализации.

На заметку! Термин *процесс* в этой главе должен трактоваться как синоним термина *поток* в операционных системах, в средах которых СУБД Oracle реализована посредством потоков (например, Windows). В контексте главы термин *процесс* охватывает как процессы, так и потоки. Если вы используете реализацию Oracle с множеством процессов, такую как в среде UNIX/Linux, то термин *процесс* полностью адекватен. Если же вы применяете реализацию Oracle с единственным процессом, характерную для среды Windows, то термин *процесс* в действительности будет означать *поток внутри процесса Oracle*. Таким образом, например, когда речь идет о процессе DBWn, его эквивалентом в среде Windows является поток DBWn внутри процесса Oracle.

Серверные процессы

Серверные процессы — это процессы, которые выполняют работу от имени сеанса клиента. Такие процессы в итоге принимают и действуют согласно SQL-операторам, которые наши приложения отправляют базе данных.

В главе 2 мы кратко посмотрели на два основных типа подключений к Oracle.

- **Подключение посредством выделенного сервера.** На сервере для подключения имеется выделенный процесс. При этом между подключением к базе данных и серверным процессом или потоком существует однозначное соответствие.
- **Подключение посредством разделяемого сервера.** Множество сеансов совместно используют пул серверных процессов, порожденных и управляемых экземпляром Oracle. Подключение производится к диспетчеру базы данных, а не к процессу выделенного сервера, который создан специально для подключения.

На заметку! Важно понимать разницу между подключением и сеансом в контексте терминологии Oracle. *Подключение* (connection) — это просто физический путь передачи информации между процессом клиента и экземпляром Oracle (например, сетевое соединение между пользователем и экземпляром). С другой стороны, *сеанс* (session) — это логическая сущность в базе данных, в которой процесс клиента может выполнять SQL-запросы и т.д. С единственным подключением может быть связано несколько независимых сеансов, и эти сеансы могут даже существовать независимо от подключения. Вскоре мы обсудим это более подробно.

Процессы выделенного и разделяемого серверов решают одну и ту же задачу: они обрабатывают все предоставляемые им SQL-запросы. Когда вы отправляете базе данных запрос `SELECT * FROM EMP`, процесс выделенного/разделяемого сервера производит синтаксический разбор запроса и помещает его в разделяемый пул (или находит его в разделяемом пуле, если он там уже присутствует). При необходимости процесс создает план запроса и выполняет его, возможно отыскивая нужные данные в буферном кеше или считывая данные с диска в буферный кеш.

Серверные процессы делают всю рутинную работу. Часто вы будете обнаруживать, что эти процессы являются основными потребителями ресурсов процессора в системе, поскольку они выполняют сортировку, суммирование и соединение — практически все необходимые действия.

Подключения посредством выделенного сервера

В режиме выделенного сервера между подключением клиента и серверным процессом (или потоком) устанавливается однозначное соответствие. При наличии на компьютере UNIX/Linux ста подключений через выделенный сервер от их имени будут выполняться сто процессов. Эта ситуация графически представлена на рис. 5.1.

Клиентское приложение будет содержать скомпонованные с ним библиотеки Oracle. Эти библиотеки предоставляют API-интерфейсы, которые требуются для обмена информацией с базой данных. Таким API-интерфейсам известно, каким образом отправлять запрос базе данных и обрабатывать возвращаемый курсор. Они также умеют объединять запросы в сетевые вызовы, которые выделенный сервер будет распаковывать.

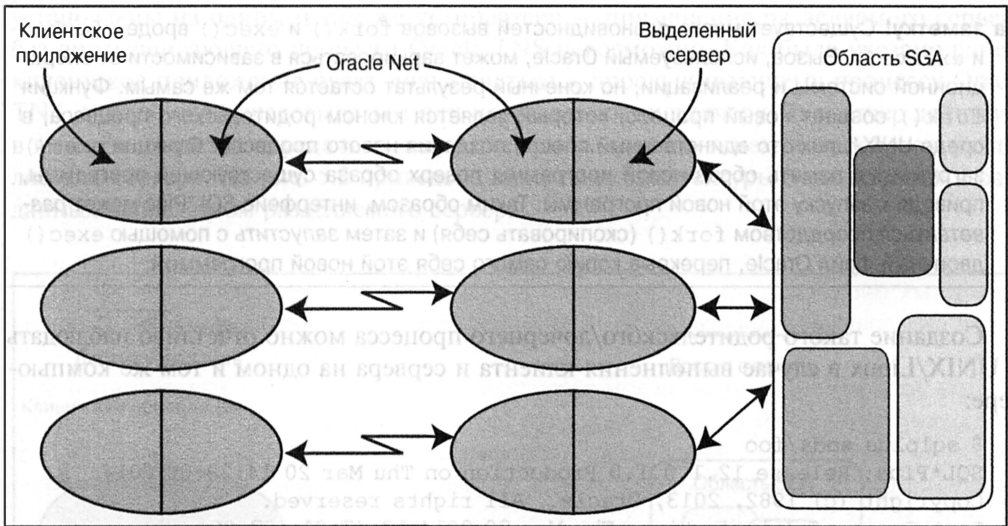


Рис. 5.1. Типичное подключение посредством выделенного сервера

Эта часть программного обеспечения называется *Oracle Net* (Сеть Oracle), хотя в предшествующих версиях ее называли *SQL*Net* или *Net8*. Она представляет собой сетевое программное обеспечение/протокол, который Oracle применяет для реализации клиент-серверной обработки (в многозвенной архитектуре также присутствует клиент-серверное программное обеспечение). СУБД Oracle использует эту же архитектуру, даже если Oracle Net формально не принимает участия. То есть такая двух-процессная архитектура (которую также называют *двухзадачной*) задействована даже в ситуации, когда клиент и сервер находятся на одном и том же компьютере. Эта архитектура предоставляет два преимущества.

- **Возможность удаленного выполнения.** Ситуация, когда клиентское приложение выполняется на компьютере, отличном от того, где функционирует база данных, является вполне естественной.
- **Обеспечение изоляции адресного пространства.** Серверный процесс имеет доступ по чтению и записи к области SGA. Если бы клиентский и серверный процессы были физически связаны друг с другом, то ошибочный указатель в клиентском процессе легко мог бы повредить структуры данных в SGA.

В главе 2 было показано, как эти выделенные серверы *порождаются*, или создаются, прослушивающим процессом Oracle. Я не буду описывать это снова; взамен мы кратко рассмотрим, что происходит в тех случаях, когда прослушивающий процесс не задействован. Механизм во многом похож на тот, который применяется вместе с прослушивающим процессом, но теперь выделенный сервер создается не прослушивающим процессом посредством вызова `fork()`/`exec()` в UNIX/Linux или вызова `IPC` (interprocess communication — взаимодействие между процессами) в Windows, а самим клиентским процессом.

На заметку! Существует много разновидностей вызовов `fork()` и `exec()` вроде `vfork()` и `execve()`. Вызов, используемый Oracle, может варьироваться в зависимости от операционной системы и реализации, но конечный результат остается тем же самым. Функция `fork()` создает новый процесс, который является клоном родительского процесса; в среде UNIX/Linux это единственный способ создания нового процесса. Функция `exec()` загружает в память образ новой программы поверх образа существующей программы, приводя к запуску этой новой программы. Таким образом, интерфейс SQL*Plus может *разветвиться* посредством `fork()` (скопировать себя) и затем *запустить* с помощью `exec()` двоичный файл Oracle, перекрыв копию самого себя этой новой программой.

Создание такого родительского/дочернего процесса можно отчетливо наблюдать в UNIX/Linux в случае выполнения клиента и сервера на одном и том же компьютере:

```
$ sqlplus eoda/foo
SQL*Plus: Release 12.1.0.1.0 Production on Thu Mar 20 14:29:00 2014
Copyright (c) 1982, 2013, Oracle. All rights reserved.
Last Successful login time: Thu Mar 20 2014 13:47:01 -07:00
Connected to:
Oracle Database 12c Enterprise Edition Release 12.1.0.1.0 -
64bit Production
With the Partitioning, OLAP, Advanced Analytics and Real Application
Testing options
```

```
EODA@ORA12CR1> select a.spid dedicated_server, b.process clientpid
2 from v$process a, v$session b
3 where a.addr = b.paddr
4 and b.sid = sys_context('userenv','sid');
```

| DEDICATED_SERVER | CLIENTPID |
|------------------|-----------|
| ----- | ----- |
| 18571 | 18570 |

```
EODA@ORA12CR1> !/bin/ps -fp 18571 18570
UID      PID  PPID  C STIME TTY      STAT   TIME CMD
oracle  18570 11782  0 15:17 pts/4    S+      0:00 sqlplus
oracle  18571 18570  0 15:17 ?        Ss      0:00 oracleORA12CR1
                                     (DESCRIPTION=(LOCAL=...
```

Здесь с помощью запроса выясняется идентификатор процесса (process ID — PID), ассоциированный с выделенным сервером (значение SPID из `V$PROCESS` — это идентификатор PID операционной системы для процесса, который применялся во время выполнения этого запроса). Вывод `/bin/ps -fp` включает идентификатор родительского процесса (parent process ID — PPID) и показывает процесс выделенного сервера, 18571, как дочерний процесс процесса SQL*Plus: идентификатором процесса является 18570.

Подключения посредством разделяемого сервера

Давайте теперь более подробно рассмотрим процесс разделяемого сервера. Этот тип подключения требует использования Oracle Net, даже если клиент и сервер

установлены на одном и том же компьютере — применение разделяемого сервера без прослушивающего процесса Oracle TNS невозможно. Как было указано ранее, клиентское приложение будет подключаться к прослушивающему процессу Oracle TNS, а затем переадресовываться или передаваться диспетчеру. Диспетчер действует в качестве средства передачи между клиентским приложением и процессом разделяемого сервера. На рис. 5.2 приведена диаграмма архитектуры подключения к базе данных посредством разделяемого сервера.

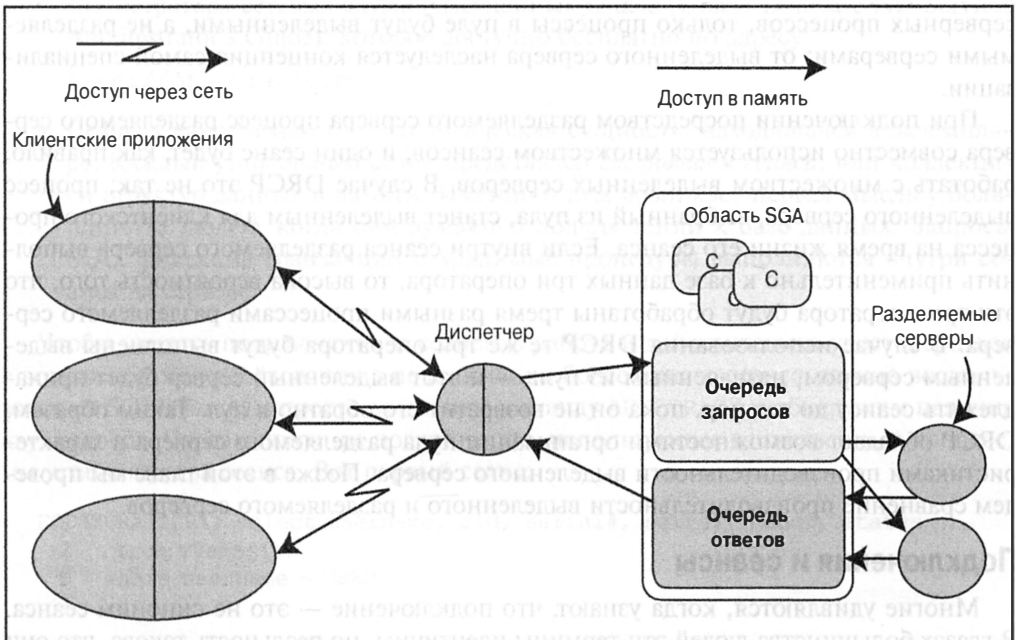


Рис. 5.2. Типичное подключение посредством разделяемого сервера

Как видите, клиентские приложения со скомпонованными в них библиотеками Oracle будут физически подключены к процессу диспетчера. Для любого заданного экземпляра может быть сконфигурировано множество диспетчеров, но не так уж редко единственный диспетчер используется для сотен и даже тысяч пользователей. Диспетчер отвечает лишь за прием входящих запросов от клиентских приложений и помещение их в очередь запросов внутри области SGA. Первый доступный процесс разделяемого сервера из пула заранее созданных процессов выделенного сервера извлечет запрос из очереди и присоединится к области UGA связанного сеанса (на рис. 5.2 это прямоугольники, помеченные “С”). Разделяемый сервер обработает запрос и поместит полученный из него вывод в очередь ответов. Диспетчер постоянно отслеживает очередь ответов на предмет появления в ней результатов и передает их обратно клиентскому приложению. С точки зрения клиента невозможно сказать, подключен он посредством выделенного или разделяемого сервера — оба типа подключения для него выглядят одинаково. Различие проявляется только на уровне базы данных.

Резидентный пул соединений с базой данных

Резидентный пул соединений с базой данных (Database Resident Connection Pooling — DRCP) представляет собой дополнительный метод подключения к базе данных и установления сеанса. Он спроектирован как более эффективный метод организации пула соединений для интерфейсов приложений, которые не имеют практической встроенной поддержки пула соединений, подобных PHP — универсальному языку веб-сценариев. Метод DRCP — это смесь концепций выделенного и разделяемого сервера. Он наследует от разделяемого сервера концепцию организации пула серверных процессов, только процессы в пуле будут выделенными, а не разделяемыми серверами; от выделенного сервера наследуется концепция самой специализации.

При подключении посредством разделяемого сервера процесс разделяемого сервера совместно используется множеством сеансов, и один сеанс будет, как правило, работать с множеством выделенных серверов. В случае DRCP это не так; процесс выделенного сервера, выбранный из пула, станет выделенным для клиентского процесса на время жизни его сеанса. Если внутри сеанса разделяемого сервера выполнить применительно к базе данных три оператора, то высока вероятность того, что эти три оператора будут обработаны тремя разными процессами разделяемого сервера. В случае использования DRCP те же три оператора будут выполнены выделенным сервером, назначенным из пула — и этот выделенный сервер будет принадлежать сеансу до тех пор, пока он не возвратит его обратно в пул. Таким образом, DRCP обладает возможностями организации пула разделяемого сервера и характеристиками производительности выделенного сервера. Позже в этой главе мы проведем сравнение производительности выделенного и разделяемого серверов.

Подключения и сеансы

Многие удивляются, когда узнают, что подключение — это не синоним сеанса. В глазах большинства людей эти термины идентичны, но реальность такова, что они не должны быть одинаковыми. Подключение может иметь ноль, один или более установленных в нем сеансов. Каждый сеанс является отдельным и независимым, даже если все они совместно используют одно и то же физическое подключение к базе данных. Фиксация в одном сеансе не оказывает влияния на любые другие сеансы этого подключения. В действительности каждый сеанс, применяющий одно подключение, мог бы использовать разные идентификационные данные пользователя!

В среде Oracle подключение представляет собой просто физическую линию связи между клиентским процессом и экземпляром базы данных — чаще всего сетевое подключение. Подключение может производиться к процессу выделенного сервера или к диспетчеру. Как утверждалось ранее, подключение может иметь ноль или большее число сеансов, т.е. допустимо наличие подключения без соответствующих ему сеансов. Вдобавок сеанс *может иметь*, а *может и не иметь* подключение. С помощью усовершенствованных функциональных средств Oracle Net, таких как организация пула подключений, клиент может закрывать физическое подключение, оставляя сеанс незатронутым (но бездействующим). Когда клиенту потребуется выполнить определенную операцию в этом *сеансе*, ему придется заново установить физическое *подключение*. Давайте определим эти термины более точно.

- **Подключение.** Подключение (connection) — это физическая линия связи между клиентом и экземпляром Oracle. Подключение устанавливается либо по сети, либо посредством механизма IPC. Как правило, подключение создается между клиентским процессом и выделенным сервером либо диспетчером. Однако с применением диспетчера подключений Oracle (Connection Manager — CMAN) подключение можно устанавливать между клиентом и CMAN и между CMAN и базой данных. Подробное описание CMAN выходит за рамки этой книги, но определенные сведения об этом инструменте можно найти в руководстве администратора сетевых служб базы данных Oracle (Oracle Database Net Services Administrator's Guide), которое доступно бесплатно по адресу:

<http://otn.oracle.com>

- **Сеанс.** Сеанс (session) — это логическая сущность, находящаяся в экземпляре. Каждый уникальный сеанс представлен *состоянием сеанса*, или коллекцией структур данных в памяти. Именно о нем возникает первая мысль у большинства людей, когда они думают о подключении к базе данных. Запросы SQL, фиксация транзакций и хранимые процедуры выполняются внутри сеанса на сервере.

Чтобы увидеть подключения и сеансы в работе и действительно удостовериться в возможности существования в подключении более одного сеанса, можно воспользоваться SQL*Plus. Мы просто применим команду AUTOTRACE и обнаружим наличие двух сеансов. Через одиночное подключение, использующее единственный процесс, мы установим два сеанса. Вот первый сеанс:

```

EODA@ORA12CR1> select username, sid, serial#, server, paddr, status
2   from v$session
3   where username = USER
4   /

```

| USERNAME | SID | SERIAL# | SERVER | PADDR | STATUS |
|----------|-----|---------|-----------|------------------|--------|
| EODA | 10 | 15 | DEDICATED | 00000000727FE9B0 | ACTIVE |

В текущий момент есть один сеанс — сеанс, подключенный через выделенный сервер. В столбце PADDR находится адрес единственного процесса выделенного сервера. На этот раз мы включим функцию AUTOTRACE для получения статистических сведений об операторах, выполняемых в среде SQL*Plus:

```

EODA@ORA12CR1> set autotrace on statistics
EODA@ORA12CR1> select username, sid, serial#, server, paddr, status
2   from v$session
3   where username = USER
4   /

```

| USERNAME | SID | SERIAL# | SERVER | PADDR | STATUS |
|----------|-----|---------|-----------|------------------|----------|
| EODA | 10 | 15 | DEDICATED | 00000000727FE9B0 | ACTIVE |
| EODA | 21 | 721 | DEDICATED | 00000000727FE9B0 | INACTIVE |

Statistics

```

-----
      8 recursive calls
      0 db block gets
      2 consistent gets
      0 physical reads
      0 redo size
    1004 bytes sent via SQL*Net to client
     543 bytes received via SQL*Net from client
        2 SQL*Net roundtrips to/from client
        0 sorts (memory)
        0 sorts (disk)
        2 rows processed

```

```
EODA@ORA12CR1> set autotrace off
```

Теперь в системе присутствуют два сеанса, но они *оба* применяют единственный процесс выделенного сервера, о чем свидетельствует одно и то же значение PADDR. В операционной системе не было создано ни одного нового процесса, и для обоих сеансов используется единственный процесс — одно подключение. Обратите внимание, что один из сеансов (исходный) имеет состояние ACTIVE (активный). Это вполне объяснимо: он выполняет запрос для отображения приведенной информации, так что, естественно, он активен. Но что собой представляет неактивный (INACTIVE) сеанс? Это сеанс AUTOTRACE. Его задачей является слежение за реальным сеансом и сообщение о том, что он делает.

В случае включения функции AUTOTRACE в SQL*Plus при выполнении операций DML (INSERT, UPDATE, DELETE, SELECT и MERGE) среда SQL*Plus будет выполнять перечисленные ниже действия.

1. Она создаст новый сеанс, использующий текущее подключение, если второй сеанс пока еще не существует.
2. Она потребует от нового сеанса запросить представление V\$SESSTAT для запоминания начальных статистических значений сеанса, в котором мы будем выполнять операцию DML. Это очень похоже на функцию, реализованную сценарием `watch_stat.sql` в главе 4.
3. Она выполнит операцию DML в исходном сеансе.
4. После завершения оператора DML среда SQL*Plus снова потребует от нового сеанса выполнения запросить представление V\$SESSTAT и сгенерировать ранее показанный отчет, отображающий отличия в статистических сведениях о сеансе, который выполнил операцию DML.

Если вы отключите функцию AUTOTRACE, то среда SQL*Plus завершит работу этого дополнительного сеанса, и вы больше не увидите его в представлении V\$SESSTAT. Почему SQL*Plus предпринимает такой трюк? Ответ довольно прост. Среда SQL*Plus делает это по той же причине, по которой в главе 4 мы использовали второй сеанс SQL*Plus для мониторинга потребления памяти и временного пространства: если бы мы применяли для отслеживания использования памяти единственный сеанс, то определенный объем памяти был бы потрачен на сам мониторинг. Наблюдая за статистическими сведениями в одном сеансе, мы неизбежно изменили бы их. Если бы среда SQL*Plus применяла единственный сеанс для сообщения

о количестве выполненных операций ввода-вывода, объеме данных, переданных по сети, и числе выполненных операций сортировки, то запросы, используемые для получения этих сведений, также внесли бы свой вклад в статистическую информацию. Они могли бы выполнять сортировку, операции ввода-вывода, передачу данных по сети и т.п. Следовательно, для получения корректных данных измерений необходимо применять другой сеанс.

До сих пор мы видели подключение с одним или двумя сеансами. Теперь было бы желательно воспользоваться SQL*Plus, чтобы посмотреть на подключение, не имеющее сеансов. Это достаточно легко. В том же окне SQL*Plus, которое было открыто в предыдущем примере, введите команду с обманчивым именем DISCONNECT:

```
EODA@ORA12CR1> disconnect
Disconnected from Oracle Database 12c Enterprise Edition Release
12.1.0.1.0 - 64bit Production
With the Partitioning, OLAP, Advanced Analytics and Real Application
Testing options
```

Формально эта команда должна была бы называться DESTROY_ALL_SESSIONS (уничтожить все сеансы), а не DISCONNECT (отключиться), т.к. в действительности она не выполняет физического отключения от базы данных.

На заметку! Настоящее отключение в SQL*Plus обеспечивает команда exit, поскольку для полного разрыва соединения вы должны произвести выход.

Тем не менее, мы закрыли все сеансы. Откроем еще один сеанс с применением учетной записи какого-то другого пользователя и выдадим следующий запрос к базе данных (разумеется, замените EODA именем своей учетной записи):

```
$ sqlplus / as sysdba
SQL*Plus: Release 12.1.0.1.0 Production on Thu Mar 20 15:57:56 2014
Copyright (c) 1982, 2013, Oracle. All rights reserved.
Connected to:
Oracle Database 12c Enterprise Edition Release 12.1.0.1.0 - 64bit Production
With the Partitioning, OLAP, Advanced Analytics and Real Application
Testing options
```

```
SYS@ORA12CR1> select * from v$session where username = 'EODA';
no rows selected
строки не выбраны
```

Легко заметить, что сеансы отсутствуют, но по-прежнему имеется процесс, физическое подключение (используется предыдущее значение PADDR):

```
SYS@ORA12CR1> select username, program
2   from v$process
3  where addr = hextoraw( '00000000727FE9B0' );
```

| USERNAME | PROGRAM |
|----------|----------------------------|
| oracle | oracle@cs-xvm2 (TNS V1-V3) |

Итак, здесь есть подключение без ассоциированных с ним сеансов. Для создания нового сеанса в существующем процессе можно применять также не очень удач-

но названную команду CONNECT из SQL*Plus (более подходящим именем было бы CREATE_SESSION (создать сеанс)). Используя экземпляр SQL*Plus, в котором было произведено отключение, мы выполняем следующие команды:

```
EODA@ORA12CR1> connect eoda/foo
Connected.
Подключено.
EODA@ORA12CR1> select username, sid, serial#, server, paddr, status
2      from v$session
3      where username = USER;
```

| USERNAME | SID | SERIAL# | SERVER | PADDR | STATUS |
|----------|-----|---------|-----------|------------------|--------|
| EODA | 10 | 25 | DEDICATED | 00000000727FE9B0 | ACTIVE |

Обратите внимание, что значение PADDR осталось таким же, как ранее, следовательно, мы имеем дело с тем же самым физическим подключением. Но значение SID потенциально может быть другим. *Потенциально* означает, что сеансу может быть назначен и тот же самый идентификатор SID — это просто зависит от того, вошли ли другие пользователи в систему, пока мы были от нее отключены, и от доступности первоначального SID.

На заметку! В Windows или других операционных системах, основанных на потоках, вы можете получить отличающиеся результаты — адрес процесса может изменяться, поскольку вы подключаетесь к многопоточному процессу, а не к специализированному процессу, как это происходит в UNIX/Linux.

До сих пор все тесты выполнялись с применением подключения посредством выделенного сервера, поэтому значение PADDR представляло адрес процесса выделенного сервера. А что происходит, когда используется разделяемый сервер?

На заметку! Для подключения через разделяемый сервер экземпляр базы данных должен быть запущен с необходимой настройкой. Тема конфигурирования разделяемого сервера выходит за рамки этой книги, но подробно рассматривается в руководстве администратора сетевых служб базы данных Oracle (*Oracle Database Net Services Administrator's Guide*).

Давайте войдем в систему с применением разделяемого сервера и выполним запрос в этом сеансе:

```
EODA@ORA12CR1> select a.username, a.sid, a.serial#, a.server,
2      a.paddr, a.status, b.program
3      from v$session a left join v$process b
4      on (a.paddr = b.addr)
5      where a.username = 'EODA'
6      /
```

| USERNAME | SID | SERIAL# | SERVER | PADDR | STATUS | PROGRAM |
|----------|-----|---------|--------|------------------|--------|-----------------------|
| EODA | 485 | 1105 | SHARED | 0000000071C4BE68 | ACTIVE | oracle@cs-xvm2 (S004) |

Наше подключение посредством разделяемого сервера ассоциировано с процессом — есть значение PADDR и с помощью соединения с V\$PROCESS можно извлечь

имя этого процесса. В этом случае мы видим, что процесс является разделяемым сервером, о чем свидетельствует строка S000.

Однако если для запроса той же информации воспользоваться другим окном SQL*Plus, оставив наш сеанс разделяемого сервера бездействующим, мы получим примерно такой результат:

```
$ sqlplus / as sysdba
```

```
SYS@ORA12CR1> select a.username, a.sid, a.serial#, a.server,
2      a.paddr, a.status, b.program
3      from v$session a left join v$process b
4      on (a.paddr = b.addr)
5      where a.username = 'EODA'
6      /
```

| USERNAME | SID | SERIAL# | SERVER | PADDR | STATUS | PROGRAM |
|----------|-----|---------|--------|------------------|----------|-----------------------|
| EODA | 485 | 1105 | NONE | 0000000071C46788 | INACTIVE | oracle@cs-xvm2 (D000) |

Обратите внимание, что значение PADDR отличается, и имя процесса, с которым мы связаны, также изменилось. Бездействующее соединение с разделяемым сервером теперь ассоциировано с диспетчером, D000. Следовательно, имеется еще один метод наблюдения за множеством сеансов, указывающих на один процесс. Диспетчер может иметь сотни или даже тысячи указывающих на него сеансов.

Интересная особенность соединений с разделяемым сервером заключается в том, что применяемый нами процесс разделяемого сервера может меняться от вызова к вызову. Если я являюсь единственным пользователем системы (и это так для проводимых тестов), то многократный запуск этого запроса от имени EODA должен производить одно и то же значение PADDR — 0000000071C4BE68. Однако если открыть больше подключений посредством разделяемого сервера и начать их использовать в других сеансах, то можно заметить, что применяемый разделяемый сервер варьируется.

Рассмотрим пример. Я запрошу информацию о текущем сеансе, в том числе имя используемого разделяемого сервера. Затем в другом сеансе разделяемого сервера я выполню длительную операцию (т.е. монополизирую этот разделяемый сервер). При повторном запросе базы данных о применяемом разделяемом сервере, скорее всего, отобразится другой сервер (как если бы исходный сервер обслуживал другой сеанс). В приведенном ниже примере выделенный полужирным код представляет второй сеанс SQL*Plus, подключенный через разделяемый сервер:

```
EODA@ORA12CR1> select a.username, a.sid, a.serial#, a.server,
2      a.paddr, a.status, b.program
3      from v$session a left join v$process b
4      on (a.paddr = b.addr)
5      where a.username = 'EODA'
6      /
```

| USERNAME | SID | SERIAL# | SERVER | PADDR | STATUS | PROGRAM |
|----------|-----|---------|--------|------------------|--------|-----------------------|
| EODA | 485 | 1105 | SHARED | 0000000071C4BE68 | ACTIVE | oracle@cs-xvm2 (S004) |

В другом терминальном сеансе подключимся к базе данных от имени пользователя SCOTT (текст shared в следующей строке подключения отображается на запись в файле tnsnames.ora, которая инструктирует SQL*Plus использовать подключение посредством разделяемого сервера):

```
$ sqlplus scott/tiger@shared
```

```
SCOTT@ORA12CR1> exec dbms_lock.sleep(20);
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.
```

В исходном подключении как EODA снова запустим запрос для просмотра информации в столбце PROGRAM:

```
EODA@ORA12CR1> select a.username, a.sid, a.serial#, a.server,
2      a.paddr, a.status, b.program
3      from v$session a left join v$process b
4      on (a.paddr = b.addr)
5      where a.username = 'EODA'
6      /
```

| USERNAME | SID | SERIAL# | SERVER | PADDR | STATUS | PROGRAM |
|----------|-----|---------|--------|------------------|--------|-----------------------|
| EODA | 485 | 1105 | SHARED | 0000000071C478E8 | ACTIVE | oracle@cs-xvm2 (S000) |

На заметку! Вам необходимо использовать учетную запись, которая имеет привилегии выполнения для пакета DBMS_LOCK. Я выдал демонстрационной учетной записи SCOTT привилегии выполнения для DBMS_LOCK следующим образом: SYS@ORA12CR1> grant execute on dbms_lock to scott;

Обратите внимание, что при первой выдаче запроса в качестве разделяемого сервера применялся S004. Затем во втором сеансе (от имени SCOTT) был запущен длительно выполняющийся оператор, который монополизировал разделяемый сервер; на этот раз им оказался S004. Выполнение работы назначается первому незанятому разделяемому серверу, и поскольку в этом случае никто другой не запрашивал использование разделяемого сервера S004, команда DBMS_LOCK передается именно ему. Когда я снова выдал запрос в первом сеансе SQL*Plus, он был назначен процессу другого разделяемого сервера, S000, т.к. разделяемый сервер S004 был занят.

Интересно отметить, что синтаксический разбор запроса (который пока не возвращает каких-либо строк) мог быть произведен разделяемым сервером S000, извлечение первой строки — сервером S001, извлечение второй строки — сервером S002, а закрытие курсора — сервером S003. То есть отдельный оператор может последовательно обрабатываться множеством разделяемых серверов.

Итак, в этом разделе мы убедились, что подключение — физическая линия связи между клиентом и экземпляром базы данных — может иметь ноль, один и более установленных для него сеансов. Мы видели один из возможных сценариев, когда применяли средство AUTOTRACE интерфейса SQL*Plus. Эту возможность эксплуатируют также и многие другие инструменты. Например, компонент Oracle Forms организует несколько сеансов одного подключения для реализации своих функций отладки. Средство многозвенной прокси-аутентификации Oracle, обеспечивающее

сквозную идентификацию пользователей от браузера до базы данных, интенсивно использует концепцию единственного подключения с множеством сеансов, но в каждом сеансе потенциально может быть задействована учетная запись другого пользователя. Мы видели, что с течением времени сеансы могут применять множество процессов, особенно в среде с разделяемым сервером. Кроме того, в случае использования пула подключений Oracle Net сеанс может быть вообще не связан ни с одним процессом; клиент будет разрывать соединение по истечении установленного интервала бездействия и прозрачным образом восстанавливать его при обнаружении активности.

Короче говоря, между подключениями и сеансами существует отношение типа “многие ко многим”. Однако наиболее распространенным вариантом, с которым большинство из нас сталкивается день ото дня, является отношение “один к одному” между выделенным сервером и единственным сеансом.

Сравнение режимов выделенного сервера, разделяемого сервера и DRCP

Прежде чем продолжить исследования остальных процессов, давайте обсудим причины, по которым существуют три режима подключения и когда один из них может оказываться более подходящим, чем другой.

Когда следует использовать выделенный сервер

Как отмечалось ранее, в режиме выделенного сервера между клиентским подключением и серверным процессом имеется соответствие типа “один к одному”. В настоящее время для всех приложений, основанных на SQL, этот метод подключения к базе данных Oracle является наиболее распространенным. Он очень прост в настройке, предоставляет самый легкий способ установления подключений и требует незначительного конфигурирования или вообще в нем не нуждается.

Поскольку есть однозначное соответствие, можно не беспокоиться о том, что длительно выполняющиеся транзакции заблокируют другие транзакции. Эти другие транзакции просто будут обрабатываться собственными выделенными процессами. Следовательно, такой метод подключения будет единственным, который вы должны применять в среде, отличной от OLTP (On-Line Transaction Processing — оперативная обработка транзакций), где могут существовать длительные транзакции. Конфигурация с выделенным сервером является рекомендуемой для Oracle, и она довольно хорошо масштабируется. До тех пор, пока сервер располагает оборудованием (процессором и ОЗУ), достаточным для обслуживания необходимого числа процессов выделенного сервера, выделенный сервер может обеспечивать поддержку тысяч параллельных подключений.

Определенные операции, такие как запуск или останов базы данных, должны выполняться только в режиме выделенного сервера, поэтому каждая база данных будет содержать либо обе конфигурации, либо только конфигурацию с выделенным сервером.

Когда следует использовать разделяемый сервер

Установка и конфигурирование разделяемого сервера хотя и не сложны, но все же требуют выполнения дополнительных действий по сравнению с настройкой вы-

деленного сервера. Тем не менее, основное отличие между этими конфигурациями связано не с их установкой и конфигурированием, а с режимом работы. В случае выделенного сервера между клиентскими подключениями и серверными процессами имеется однозначное соответствие. В случае разделяемого сервера существует отношение “многие к одному”: с одним разделяемым сервером связано несколько клиентов.

Как следует из названия, разделяемый сервер является совместно используемым ресурсом, тогда как выделенный сервер — нет. При работе с совместно используемым ресурсом необходимо соблюдать осторожность, чтобы не монополизировать его на длительные периоды времени. Как уже было показано ранее, выдача простой команды `DBMS_LOCK.SLEEP(20)` в одном сеансе монополизировала бы процесс разделяемого сервера на 20 секунд. Монополизация ресурсов разделяемого сервера может создавать видимость “зависания” системы.

На рис. 5.2 были изображены два разделяемых сервера. При наличии трех клиентов, которые все пытаются выполнить 45-секундный процесс приблизительно в одно и то же время, два из них получили бы ответ за 45 секунд, а третий — за 90 секунд. Отсюда правило номер один для разделяемого сервера: обеспечьте, чтобы все транзакции были короткими по продолжительности. Они могут выполняться часто, но должны быть краткими (в соответствии с характеристиками систем OLTP). Если они не являются короткими, вы попадете в ситуацию, когда работа всей системы будет выглядеть замедленной из-за того, что разделяемые ресурсы монополизированы несколькими процессами. В экстремальных ситуациях, если все разделяемые серверы заняты, система будет казаться зависшей всем пользователям за исключением нескольких счастливиц, которые монополизировали разделяемые серверы.

Еще одна интересная ситуация, которую можно наблюдать при работе с разделяемым сервером — *искусственная взаимоблокировка*. В случае применения разделяемого сервера несколько серверных процессов совместно используются потенциально большим сообществом пользователей. Представим ситуацию с пятью разделяемыми серверами и сотней установленных сеансов пользователей. При таких условиях в каждый момент времени могут быть активными максимум пять из этих сеансов пользователей. Предположим, что один из сеансов пользователей обновляет строку и не выполняет фиксацию. Пока этот пользователь сидит и размышляет над своей транзакцией, пять других сеансов пользователей пытаются заблокировать эту же строку. Естественно, они окажутся заблокированными и будут вынуждены терпеливо дожидаться момента, когда эта строка станет доступной. Наконец, сеанс пользователя, который установил блокировку данной строки, предпринимает попытку зафиксировать свою транзакцию (снимая тем самым блокировку со строки). Этот сеанс пользователя обнаружит, что все разделяемые серверы монополизированы пятью ожидающими сеансами. Возникает ситуация искусственной взаимоблокировки: владелец блокировки не получит в свое распоряжение разделяемый сервер, чтобы выполнить фиксацию транзакции, если только один из ожидающих сеансов не уступит свой разделяемый сервер. Но если для ожидающих сеансов не определен интервал истечения времени ожидания, то они никогда не уступят свои разделяемые серверы (разумеется, можно потребовать от администратора, чтобы он уничтожил свой сеанс разделяемого сервера для выхода из этой тупиковой ситуации).

По описанным причинам разделяемый сервер подходит только для систем оперативной обработки транзакций (OLTP), которые характеризуются короткими, час-

то выполняющимися транзакциями. В системе OLTP транзакции выполняются в течение нескольких миллисекунд; нет ничего, что занимало бы более доли секунды. Разделяемый сервер совершенно не подходит для информационных хранилищ. В них могут выполняться запросы, которые занимают одну, две, пять и более минут. В режиме разделяемого сервера это было бы смертельно опасным. Если система на 90% является OLTP и на 10% — “не совсем OLTP”, то в одном и том же экземпляре можно совместно применять выделенные серверы и разделяемый сервер. Подобным образом можно радикально уменьшить количество серверных процессов для пользователей OLTP и одновременно предотвратить монополизацию своих разделяемых серверов пользователями, которые не являются “в полной мере OLTP”. Вдобавок администратор базы данных может использовать встроенный диспетчер ресурсов (Resource Manager) для еще большего контроля над утилизацией ресурсов.

Естественно, веской причиной применения разделяемого сервера является отсутствие выбора. Многие усовершенствованные функции подключения требуют использования разделяемого сервера. Если нужны пулы подключений Oracle Net, то применение разделяемого сервера обязательно. Концентрация соединений нескольких баз данных также требует использования разделяемого сервера для таких подключений.

На заметку! Если в приложении вы уже применяете пул подключений (например, пул подключений J2EE), и его размер соответствующим образом установлен, то использование разделяемого сервера приведет только к снижению производительности. Пул подключений уже настроен на поддержку такого количества подключений, которое будет иметь место в любой момент времени; вы хотите, чтобы каждое из этих подключений было прямым подключением к выделенному серверу. В противном случае вы просто получите пул подключений, соединяющийся с еще одним пулом подключений.

Потенциальные преимущества разделяемого сервера

Каковы преимущества применения разделяемого сервера с учетом того, что необходимо соблюдать определенную осторожность в отношении типов транзакций, которым разрешено его использовать? Разделяемый сервер делает три вещи: уменьшает количество процессов/потоков операционной системы, искусственно ограничивает уровень параллелизма и снижает требуемый объем памяти в системе. Давайте обсудим эти моменты более подробно.

Уменьшение количества процессов/потоков операционной системы

В системе с несколькими тысячами пользователей попытка управления тысячами процессов может быстро приводить к перегрузке операционной системы. В типичной системе в любой момент времени одновременно активна только часть из этих тысяч пользователей. Например, недавно мне пришлось работать с системой, которая обслуживала 5000 параллельных пользователей. В каждый конкретный момент времени максимум 50 из них были активными. Эта система могла бы эффективно работать при организации 50 процессов разделяемых серверов, что позволило бы сократить количество процессов, которыми должна управлять операционная система, на два порядка (в 100 раз). При такой конфигурации операционная система может в значительной степени избегать переключения контекста.

Искусственное ограничение уровня параллелизма

Мне, как человеку, который принимал участие во множестве сравнительных тестов, это преимущество совершенно очевидно. При выполнении сравнительных тестов люди часто просят создавать максимальное количество пользователей вплоть до отказа работы системы. Одним из результатов таких сравнительных тестов всегда является график зависимости числа транзакций от количества одновременно работающих пользователей (рис. 5.3).

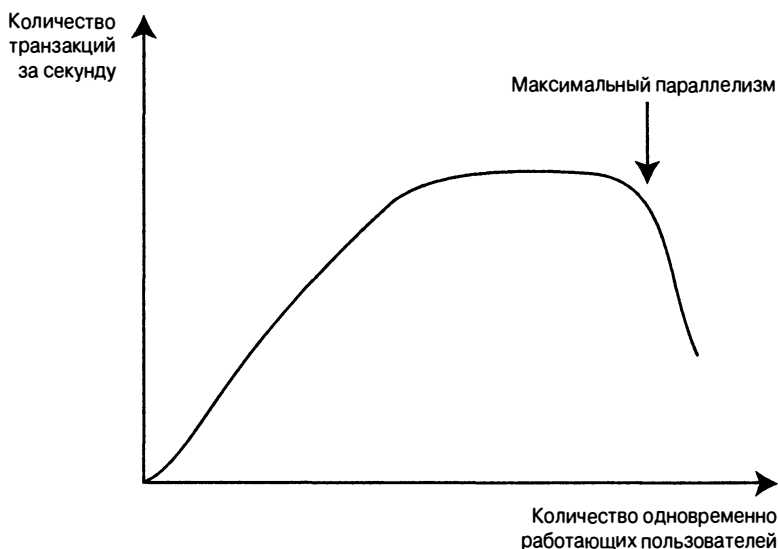


Рис. 5.3. Зависимость числа транзакций в секунду от количества одновременно работающих пользователей

Первоначально по мере добавления одновременно работающих пользователей количество транзакций увеличивается. Однако в определенный момент добавление пользователей не приводит к увеличению числа транзакций, выполняемых в секунду; кривая графика начинает понижаться. Пропускная способность достигла максимума, и теперь начало расти время отклика. Другими словами, вы делаете то же самое количество транзакций в секунду, но конечные пользователи отмечают увеличившееся время отклика. Дальнейшее добавление пользователей ведет к действительному снижению пропускной способности. Число одновременно работающих пользователей, предшествующее этому снижению, соответствует максимальному уровню параллелизма, который допустим для данной системы. Выше этого порога система становится переполненной и для выполнения работы начинают формироваться очереди. Подобно затору на пункте дорожного сбора, система больше не может справляться с нагрузкой. В этот момент не только существенно возрастает время отклика, но также может снижаться пропускная способность системы, поскольку накладные расходы, обусловленные переключением контекста и совместным использованием ресурсов между слишком большим количеством потребителей, сами требуют дополнительных ресурсов. Если ограничить максимальный уровень параллелизма значением, непосредственно предшествующим этому спаду графика,

можно сохранить максимальную пропускную способность и свести к минимуму рост времени отклика для большинства пользователей. Разделяемый сервер позволяет ограничить максимальный уровень параллелизма в системе этим значением.

В качестве аналогии этого процесса можно привести простую дверь. Ширина двери и ширина каждого человека ограничивают количество людей, которые могут пройти через дверь за одну минуту. При низкой нагрузке нет никаких проблем; однако по мере того, как подходит больше людей, некоторые из них вынуждены ждать (квант времени процессора). Если пройти через дверь хотят много людей, возникает эффект обхода — слишком многие говорят “после вас” и топчутся на месте, приводя к падению пропускной способности. Каждый человек проходит с задержкой. Применение очереди означает повышение пропускной способности. Некоторые люди проходят через дверь почти так же быстро, как если бы не было никакой очереди, в то время как другие (попавшие в хвост очереди) ожидают долго и могут беспокоиться о том, что “это была плохой идеей”. Но когда вы измеряете, насколько быстро каждый человек проходил через дверь (включая последнего в очереди), выясняется, что модель с очередью (разделяемый сервер) работает лучше свободного для всех прохода (даже в случае хорошо воспитанных людей; но вообразите себе двери в магазине в день большой распродажи, когда все грубо толкаются, пытаются попасть внутрь).

Снижение требований к памяти системы

Это один из наиболее громко рекламируемых доводов в пользу разделяемого сервера: он снижает объем требуемой памяти. Действительно, снижение есть, но не настолько значительное, как вы могли бы подумать, особенно с учетом автоматического управления памятью PGA, обсуждаемого в главе 4, когда рабочие области выделяются процессу, используются и освобождаются — причем их размеры изменяются на основе параллельной рабочей нагрузки. Таким образом, этот довод был *более весомым* в старых версиях Oracle, но в наши дни он не так важен. Кроме того, помните, что когда вы применяете разделяемый сервер, область UGA располагается в области SGA. Это означает, что при переходе на разделяемый сервер вы должны быть способны точно определить ожидаемые потребности в памяти UGA и выделить соответствующий объем в области SGA через параметр `LARGE_POOL_SIZE`. Требования к памяти SGA для конфигурации с разделяемым сервером обычно очень высоки. Эта память должна, как правило, выделяться заранее и, следовательно, может использоваться только экземпляром базы данных.

На заметку! Это правда, что при наличии области SGA изменяемого размера вы можете увеличивать и уменьшать ее размер с течением времени, но в целом ею будет владеть экземпляр базы данных, и она не доступна для использования другими процессами.

Сравните это с выделенным сервером, где любой желающий может работать с любой памятью, не выделенной в области SGA. Но если размер области SGA намного больше из-за нахождения в ней области UGA, то откуда берется экономия памяти? Она обусловлена тем, что выделяется гораздо меньше областей PGA. Каждый выделенный/разделяемый сервер имеет область PGA. В ней хранится информация о процессе. Здесь же расположены области сортировки, области хеширования и другие структуры, связанные с процессом. Именно эту память вы изымаете из системы

за счет применения разделяемого сервера. Если вы переходите от 5000 выделенных серверов к 100 разделяемым серверам, то экономия памяти, связанная с использованием разделяемых серверов, будет равна суммарному размеру 4900 областей PGA (исключая их области UGA), которые больше не нужны.

DRCP

А что можно сказать о DRCP — средстве, доступном в Oracle 11g и последующих версиях? Оно обладает многими преимуществами разделяемого сервера, такими как сокращение количества процессов (применяется пул), что делает возможной экономию памяти без отрицательных последствий. Здесь нет шансов столкнуться с искусственной взаимоблокировкой; например, сеанс, удерживающий блокировку на ресурсе в рассмотренном ранее примере, будет иметь собственный выделенный сервер, *назначенный* ему из пула, и этот сеанс, в конечном счете, сможет снять эту блокировку. Средство DRCP не поддерживает многопоточность разделяемого сервера; когда клиентский процесс получает выделенный сервер из пула, он *владеет* этим процессом до тех пор, пока клиентский процесс не освободит его. Следовательно, средство DRCP лучше подходит для клиентских приложений, которые часто подключаются, выполняют некоторое относительно короткое действие и отключаются — и так снова и снова; короче говоря, оно предназначено для клиентских процессов, имеющих API-интерфейс, который не располагает собственным эффективным пулом соединений.

Заключительные соображения по поводу выделенного/разделяемого сервера

Если только ваша система не перегружена или вы должны использовать разделяемый сервер для отдельной функции, то выделенный сервер, скорее всего, будет самым подходящим вариантом. Он прост в установке (фактически, вообще ее не требует) и легче в настройке.

На заметку! В случае подключений с помощью разделяемого сервера информация трассировки сеанса (вывод SQL_TRACE=TRUE) может распространяться по множеству отдельных трассировочных файлов; таким образом, воспроизведение того, что сделал сеанс, становится труднее. С появлением пакета DBMS_MONITOR в Oracle 10g и последующих версиях большая часть сложностей устранена, но воспроизведение по-прежнему является непростой задачей. Кроме того, при наличии многочисленных связанных трассировочных файлов, сгенерированных сеансом, можно применять утилиту TRCSESS для объединения всех этих файлов.

Если вы имеете дело с очень большим сообществом пользователей и *знаете*, что развертывание будет осуществляться посредством разделяемого сервера, я настоятельно рекомендую проводить *разработку и тестирование* с использованием разделяемого сервера. Разработка с применением только выделенного сервера без тестирования на разделяемом сервере увеличит вероятность возникновения сбоя. Обеспечьте максимальную нагрузку в системе, проведите сравнительные тесты производительности и удостоверьтесь, что приложение будет нормально работать с разделяемым сервером. То есть убедитесь, что приложение не монополизирует разделяемые серверы на слишком длительный период. Если вы обнаружите, что оно делает

это, во время разработки, то ситуацию гораздо легче исправить на данном этапе, чем на стадии развертывания. Для преобразования длительно выполняющегося процесса в предположительно короткий можно воспользоваться средствами наподобие AQ (Advanced Queuing — усовершенствованная организация очередей), но это должно быть *спроектировано* внутри приложения. Такого рода вещи лучше делать во время разработки. Кроме того, исторически сложилось так, что между наборами функциональности, которые доступны подключениям через разделяемый сервер и подключениям через выделенный сервер, имеются отличия. Например, мы уже обсуждали отсутствие автоматического управления памятью PGA в Oracle9i, но в прошлом для подключений через разделяемый сервер не были доступны даже такие важные функции, как хеш-соединения двух таблиц. (Хеш-соединения доступны в текущем выпуске Oracle9i и последующих версиях с разделяемым сервером!)

Фоновые процессы

Экземпляр Oracle образован двумя компонентами: памятью SGA и фоновыми процессами. Фоновые процессы выполняют обычные задачи по обслуживанию, необходимые для поддержания базы данных в рабочем состоянии. Например, существует процесс, который обслуживает кеш буферов блоков, по мере необходимости записывая блоки в файлы данных. Другой процесс отвечает за копирование файла оперативного журнала повторения при его заполнении в каталог архива. Еще один процесс выполняет очистку за прерванными процессами и т.д. Каждый из этих процессов довольно точно ориентирован на решение своей задачи, но работает во взаимодействии со всеми остальными процессами. Например, когда процесс, отвечающий за запись журнальных файлов, заполнит один журнал и перейдет к следующему, он уведомит процесс, который отвечает за архивирование заполненного журнального файла, о появившейся для него работе.

Доступно представление `V$`, которое можно применять для просмотра всех *возможных* фоновых процессов Oracle и выяснения, какие из них используются в системе в текущий момент:

```
EODA@ORA12CR1> select paddr, name, description
2   from v$bgprocess
3   order by paddr desc
4   /
```

| PADDR | NAME | DESCRIPTION |
|------------------|------|---------------------------------|
| 0000000072FD44C8 | MMON | Manageability Monitor Process |
| 00000000727FD850 | MMNL | Manageability Monitor Process 2 |
| 00000000723FE138 | LREG | Listener Registration |
| 00000000723FCFD8 | SMON | System Monitor Process |
| 00000000723F9BB8 | CKPT | checkpoint |
| 00000000723F8A58 | LGWR | Redo etc. |
| 00000000723F78F8 | DBW0 | db writer process 0 |
| ... | | |
| 00 | VMBO | Volume Membership 0 |
| 00 | ACFS | ACFS CSS |
| 00 | SCRB | ASM Scrubbing Master |
| 00 | XDMG | cell automation manager |
| 00 | XDWK | cell automation worker actions |

401 rows selected.

Строки этого представления со значением PADDR, отличающимся от 00, представляют процессы (потоки), сконфигурированные и выполняющиеся в системе.

Совет. Другой способ просмотра фоновых процессов, выполняющихся в текущий момент, предусматривает запрашивание строк представления V\$PROCESS, в которых значение PNAME не равно null.

Различают два класса фоновых процессов: ориентированные на конкретные задачи (как только что описанный) и предназначенные для выполнения многообразия других задач (т.е. служебные процессы). Например, существует служебный фоновый процесс для обработки внутренних очередей заданий, доступный через пакеты DBMS_JOB/DBMS_SCHEDULER. Этот процесс отслеживает очереди заданий и запускает все, что в них появляется. Во многих отношениях он напоминает процесс выделенного сервера, но без клиентского подключения. Давайте исследуем каждый из этих фоновых процессов, начав с тех, которые ориентированы на конкретные задачи, а затем перейдя к служебным процессам.

Специализированные фоновые процессы

Количество, имена и типы специализированных фоновых процессов варьируются в зависимости от выпуска. На рис. 5.4 показан *типичный набор* фоновых процессов Oracle, имеющих конкретное назначение.

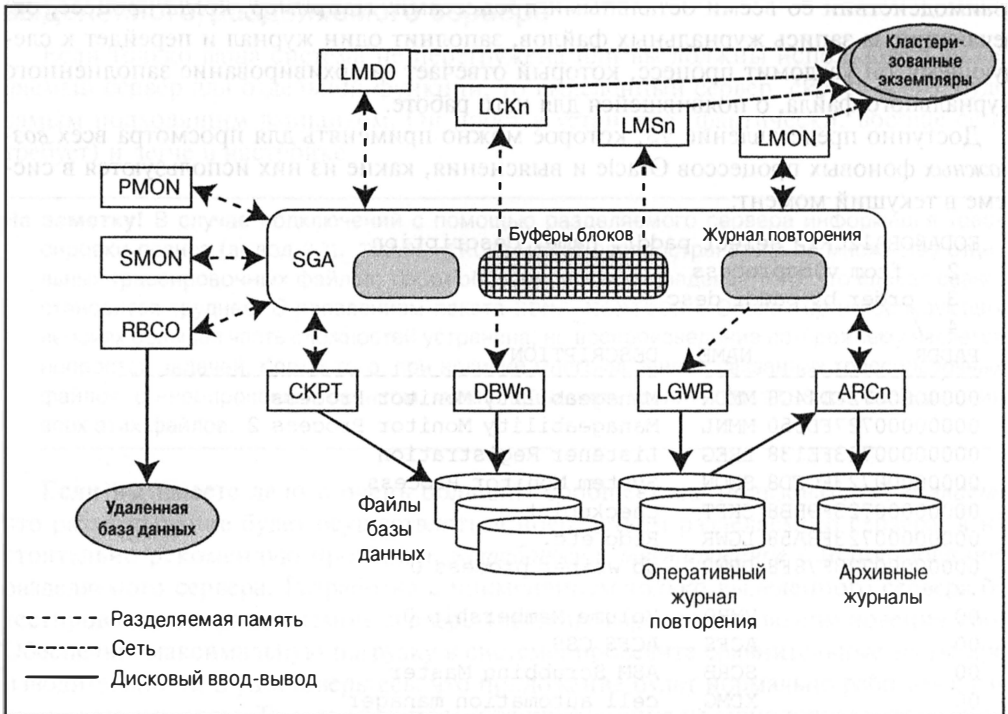


Рис. 5.4. Специализированные фоновые процессы

Совет. Полная диаграмма, отображающая все фоновые процессы и структуры памяти Oracle 12c, доступна в файле по адресу:

www.oracle.com/technetwork/tutorials/posterfiles-1974103.pdf

Например, рассмотрим базу данных, запускаемую с применением минимального числа параметров `init.ora`, в версии Oracle 12c Release 1.

```
SYS@ORA12CR1> create pfile='/tmp/pfile' from spfile;
File created.
Файл создан.

SYS@ORA12CR1> !cat /tmp/pfile
...
*.compatible='12.1.0.1'
*.control_files='/u01/dbfile/ORA12CR1/control01.ctl','/u02/dbfile/
ORA12CR1/control02.ctl'
*.db_block_size=8192
*.db_name='ORA12CR1'
*.memory_target=500M
*.undo_tablespace='UNDOTBS1'
```

В Oracle 12c Release 1 запускается около 22 фоновых процессов:

```
EOA@ORA12CR1> select paddr, name, description
2   from v$bgprocess
3  where paddr <> '00'
4  order by paddr desc
5  /
```

| PADDR | NAME | DESCRIPTION |
|------------------|------|---------------------------------------|
| 000000007F188558 | VKRM | Virtual sKeduler for Resource Manager |
| 000000007F17FA58 | CJQ0 | Job Queue Coordinator |
| 000000007F17E8F8 | AQPC | AQ Process Coord |
| 000000007F17D798 | FBDA | Flashback Data Archiver Process |
| 000000007F17C638 | SMCO | Space Manager Process |
| 000000007F17A378 | TMON | Transport Monitor |
| 000000007F175DF8 | MMNL | Manageability Monitor Process 2 |
| 000000007F174C98 | LREG | Listener Registration |
| 000000007F173B38 | RECO | distributed recovery |
| 000000007F1729D8 | SMON | System Monitor Process |
| 000000007F16F5B8 | CKPT | checkpoint |
| 000000007F16E458 | LGWR | Redo etc. |
| 000000007F16D2F8 | DBW0 | db writer process 0 |
| 000000007F16C198 | DIA0 | diagnosibility process 0 |
| 000000007F16B038 | DBRM | DataBase Resource Manager |
| 000000007F169ED8 | DIAG | diagnosibility process |
| 000000007F168D78 | MMON | Manageability Monitor Process |
| 000000007F167C18 | MMAN | Memory Manager |
| 000000007F166AB8 | GEN0 | generic0 |
| 000000007F165958 | VKTM | Virtual Keeper of TiMe process |
| 000000007F1647F8 | PSP0 | process spawner 0 |
| 000000007F163698 | PMON | process cleanup |

22 rows selected.

22 строк выбрано.

В версии Oracle 11g Release 2 с тем же самым файлом `init.ora` запускается около 17 фоновых процессов:

```
ops$tkyte@ORA11GR2> select paddr, name, description
2   from v$bgprocess
3   where paddr <> '00'
4   order by paddr desc
5   /
```

| PADDR | NAME | DESCRIPTION |
|----------|------|---------------------------------|
| 32AF0E64 | CJQO | Job Queue Coordinator |
| 32AEF8B4 | QMNC | AQ Coordinator |
| 32AEE304 | MMNL | Manageability Monitor Process 2 |
| 32AED82C | MMON | Manageability Monitor Process |
| 32AECD54 | RECO | distributed recovery |
| 32AEC27C | SMON | System Monitor Process |
| 32AEB7A4 | CKPT | checkpoint |
| 32AEACCC | LGWR | Redo etc. |
| 32AEA1F4 | DBW0 | db writer process 0 |
| 32AE971C | MMAN | Memory Manager |
| 32AE8C44 | DIA0 | diagnosibility process 0 |
| 32AE816C | PSP0 | process spawner 0 |
| 32AE7694 | DBRM | DataBase Resource Manager |
| 32AE6BBC | DIAG | diagnosibility process |
| 32AE60E4 | GEN0 | generic0 |
| 32AE560C | VKTM | Virtual Keeper of Time process |
| 32AE4B34 | PMON | process cleanup |

17 rows selected.

Воспользовавшись тем же файлом `init.ora`, но с заменой в нем `MEMORY_TARGET` на `SGA_TARGET` и `PGA_AGGREGATE_TARGET`, в Oracle 10g Release 2 вы увидите всего 12 процессов:

```
ops$tkyte@ORA10GR2> select paddr, name, description
2   from v$bgprocess
3   where paddr <> '00'
4   order by paddr desc
5   /
```

| PADDR | NAME | DESCRIPTION |
|----------|------|---------------------------------|
| 23D27AC4 | CJQO | Job Queue Coordinator |
| 23D27508 | QMNC | AQ Coordinator |
| 23D26990 | MMNL | Manageability Monitor Process 2 |
| 23D263D4 | MMON | Manageability Monitor Process |
| 23D25E18 | RECO | distributed recovery |
| 23D2585C | SMON | System Monitor Process |
| 23D252A0 | CKPT | checkpoint |
| 23D24CE4 | LGWR | Redo etc. |
| 23D24728 | DBW0 | db writer process 0 |
| 23D2416C | MMAN | Memory Manager |
| 23D23BB0 | PSP0 | process spawner 0 |
| 23D235F4 | PMON | process cleanup |

12 rows selected.

Обратите внимание, что при запуске экземпляра вы можете не увидеть все эти процессы, но большинство из них будет присутствовать. Процесс ARcN (архиватор) вы будете наблюдать только в режиме ARCHIVELOG при включенном автоматическом архивировании. Процессы LMD0, LCKn, LMON и LMSn (подробнее они рассматриваются ниже) будут видны только в случае запуска Oracle RAC — конфигурации Oracle, которая позволяет множеству экземпляров на разных компьютерах в кластере монтировать и открывать одну и ту же физическую базу данных.

На заметку! Начиная с выпуска Oracle 12c, на некоторых платформах UNIX/Linux вы можете применять смесь процессов и потоков. Эта возможность включается установкой параметра инициализации `THREADED_EXECUTION` в `TRUE` (стандартным значением является `FALSE`). Фундаментально ничего не меняется — все “процессы” по-прежнему на месте; вы лишь можете не увидеть их с помощью команды `ps` операционной системы, т.к. теперь они выполняются в виде потока в рамках более крупного процесса.

Итак, на рис. 5.4 грубо изображено то, что можно видеть после запуска экземпляра Oracle и последующего монтирования и открытия базы данных. В среде операционной системы, где Oracle реализует архитектуру с множеством процессов, такой как UNIX/Linux, вы можете физически наблюдать эти процессы. После запуска экземпляра я вижу следующую информацию:

```
$ ps -aef | grep ora_... $ORACLE_SID | grep -v grep
oracle  2276      1  0 10:33 ?        00:00:00 ora_pmon_ORA12CR1
oracle  2278      1  0 10:33 ?        00:00:00 ora_psp0_ORA12CR1
oracle  2280      1  0 10:33 ?        00:00:04 ora_vktm_ORA12CR1
oracle  2284      1  0 10:33 ?        00:00:00 ora_gen0_ORA12CR1
oracle  2286      1  0 10:33 ?        00:00:00 ora_mman_ORA12CR1
oracle  2290      1  0 10:33 ?        00:00:00 ora_diag_ORA12CR1
oracle  2292      1  0 10:33 ?        00:00:00 ora_dbrm_ORA12CR1
oracle  2294      1  0 10:33 ?        00:00:00 ora_dia0_ORA12CR1
oracle  2297      1  0 10:33 ?        00:00:00 ora_dbw0_ORA12CR1
oracle  2299      1  0 10:33 ?        00:00:00 ora_lgwr_ORA12CR1
oracle  2301      1  0 10:33 ?        00:00:00 ora_ckpt_ORA12CR1
oracle  2303      1  0 10:33 ?        00:00:00 ora_lg00_ORA12CR1
oracle  2305      1  0 10:33 ?        00:00:00 ora_lg01_ORA12CR1
oracle  2307      1  0 10:33 ?        00:00:00 ora_smon_ORA12CR1
oracle  2309      1  0 10:33 ?        00:00:00 ora_reco_ORA12CR1
oracle  2311      1  0 10:33 ?        00:00:00 ora_lreg_ORA12CR1
oracle  2313      1  0 10:33 ?        00:00:01 ora_mmon_ORA12CR1
oracle  2315      1  0 10:33 ?        00:00:00 ora_mmln_ORA12CR1
oracle  2327      1  0 10:33 ?        00:00:00 ora_p000_ORA12CR1
oracle  2329      1  0 10:33 ?        00:00:00 ora_p001_ORA12CR1
oracle  2331      1  0 10:33 ?        00:00:00 ora_tmon_ORA12CR1
oracle  2333      1  0 10:33 ?        00:00:00 ora_tt00_ORA12CR1
oracle  2335      1  0 10:33 ?        00:00:00 ora_smco_ORA12CR1
oracle  2337      1  0 10:33 ?        00:00:00 ora_fbda_ORA12CR1
oracle  2339      1  0 10:33 ?        00:00:00 ora_aqpc_ORA12CR1
oracle  2343      1  0 10:33 ?        00:00:00 ora_cjq0_ORA12CR1
oracle  2345      1  0 10:33 ?        00:00:00 ora_p002_ORA12CR1
oracle  2347      1  0 10:33 ?        00:00:00 ora_p003_ORA12CR1
oracle  2349      1  0 10:33 ?        00:00:00 ora_p004_ORA12CR1
```

```

oracle  2351    1 0 10:33 ?      00:00:00 ora_p005_ORA12CR1
oracle  2353    1 0 10:33 ?      00:00:00 ora_p006_ORA12CR1
oracle  2355    1 0 10:33 ?      00:00:00 ora_p007_ORA12CR1
oracle  2383    1 0 10:33 ?      00:00:00 ora_w000_ORA12CR1
oracle  2385    1 0 10:33 ?      00:00:00 ora_qm02_ORA12CR1
oracle  2389    1 0 10:33 ?      00:00:00 ora_q002_ORA12CR1
oracle  2391    1 0 10:33 ?      00:00:00 ora_q003_ORA12CR1
oracle  2465    1 0 10:38 ?      00:00:00 ora_w001_ORA12CR1

```

Интересно отметить используемое этими процессами соглашение по именованию. Имя процесса начинается с префикса `ora_`. За ним следуют четыре символа, представляющие действительное имя процесса, а потом — `_ORA12CR1`. Так получилось, что мой `ORACLE_SID` (идентификатор узла) выглядит как `ORA12CR1`. В системе UNIX/Linux это соглашение по именованию позволяет без труда идентифицировать фоновые процессы Oracle и ассоциировать их с конкретным экземпляром (в системе Windows простого способа достичь аналогичной цели не существует, поскольку фоновые задачи являются потоками внутри большого одиночного процесса). Наиболее интересная, но не сразу заметная в приведенном коде особенность, заключается в том, что *все эти процессы в действительности представляют собой одну и ту же двоичную исполняемую программу — не существует отдельного исполняемого файла для каждой “программы”*. Как бы тщательно вы не искали, вы нигде не обнаружите на диске двоичный исполняемый файл с именем `ora_pmon_ORA12CR1`. Вы не найдете и файл `ora_lgwr_ORA12CR1` либо `ora_reco_ORA12CR1`. На самом деле все упомянутые процессы — это `oracle` (таково имя запускаемого двоичного исполняемого файла). Просто во время запуска экземпляра они назначают себе псевдонимы, облегчающие их идентификацию. На платформе UNIX/Linux такой подход обеспечивает эффективное совместное использование значительной части объектного кода. В среде Windows данное обстоятельство не представляет такого интереса, поскольку эти компоненты являются лишь потоками внутри процесса и, следовательно, они входят в один большой двоичный файл.

Давайте теперь взглянем на функцию, выполняемую каждым *крупным* процессом, начиная с основных фоновых процессов Oracle. Полный список возможных фоновых процессов и краткое описание выполняемых ими функций вы найдете в руководстве *Oracle Database Reference*, свободно доступном по адресу:

<http://otn.oracle.com>

PMON: монитор процессов

Этот процесс отвечает за очистку после ненормального закрытия подключений. Например, если выделенный сервер выдает сбой или уничтожается по какой-то причине, то PMON является процессом, который отвечает за исправление ситуации (восстановление или отмену работы) и освобождение ресурсов. Процесс PMON будет инициировать откат незафиксированных транзакций, снятие блокировок и освобождение ресурсов SGA, выделенных отказавшему процессу.

В дополнение к очистке после прерванных подключений процесс PMON отвечает за мониторинг других фоновых процессов Oracle, а также их перезапуск в случае необходимости (если это возможно). Если в разделяемом сервере или диспетчере происходит отказ (авария), процесс PMON вступает в действие и перезапускает другой процесс (после выполнения очистки за отказавшим процессом).

Процесс PMON будет наблюдать за всеми процессами Oracle и либо перезапускать их, либо прекращать работу экземпляра подходящим образом. Например, в случае отказа процесса записи журнала базы данных, LGWR, имеет смысл прервать работу экземпляра. Эта серьезная ошибка, и наиболее безопасным образом действий будет немедленное прекращение работы экземпляра и предоставление программе восстановления возможности исправления данных. (Следует отметить, что подобная ситуация встречается редко, и о ней нужно немедленно сообщать в службу поддержки Oracle Support.)

На заметку! До выхода версии Oracle 12c, процесс PMON обрабатывал задачи регистрации прослушивателем. Начиная с Oracle 12c, фоновый процесс регистрации прослушивателем (LREG) регистрирует экземпляры и службы посредством прослушивателя.

LREG: процесс регистрации прослушивателем

Процесс LREG (доступный, начиная с выпуска Oracle 12c) отвечает за регистрацию экземпляров и служб с помощью прослушивающего процесса Oracle TNS. Когда экземпляр запускается, процесс LREG запрашивает известный адрес порта, если только он не переадресован в другое место, чтобы выяснить, функционирует ли прослушивающий процесс. Известным/стандартным портом, применяемым Oracle, является порт 1521. Что случится, если прослушивающий процесс запущен на каком-то другом порте? В этом случае механизм остается таким же за исключением того, что адрес прослушивающего процесса должен быть явно указан в параметре LOCAL_LISTENER. Вы можете также использовать параметр REMOTE_LISTENER, чтобы проинструктировать LREG о необходимости регистрировать службу посредством удаленного прослушивающего процесса (это распространенная ситуация в средах Oracle RAC).

Если прослушивающий процесс работает во время запуска экземпляра базы данных, то LREG связывается с прослушивающим процессом и передает ему соответствующие параметры, такие как имя службы и показатели нагрузки экземпляра. Если же прослушивающий процесс не был запущен, LREG будет периодически (обычно каждые 60 секунд) пробовать связаться с ним, чтобы зарегистрировать себя. Вы можете также выдать команду ALTER SYSTEM REGISTER, указав LREG о том, что нужно попытаться немедленно зарегистрировать службу с помощью прослушивающего процесса (это полезно в средах с высокой готовностью).

SMON: системный монитор

SMON — это процесс, который делает всю работу на уровне системы. В то время как PMON обслуживает отдельные процессы, SMON занят проблемами системного уровня и является своего рода сборщиком мусора в базе данных. Ниже перечислены некоторые выполняемые им задачи.

- **Очистка временного пространства.** С появлением настоящих временных табличных пространств задача очистки временного пространства стала менее трудоемкой, но не исчезла. Например, при построении индекса экстенды, выделенные для индекса во время создания, помечаются как TEMPORARY (временный). В случае прерывания сеанса CREATE INDEX по какой-либо причине процесс SMON должен очистить эти экстенды. Другие операции также создают временные экстенды, за очистку которых также отвечает PMON.

- **Объединение свободного пространства.** При использовании табличных пространств, управляемых словарем, процесс SMON отвечает за объединение групп свободных, расположенных по соседству друг с другом экстентов табличного пространства в один большой свободный экстенст. Это объединение выполняется только в табличных пространствах, управляемых словарем, со стандартной конструкцией хранения, которая имеет параметр PCTINCREASE, установленный в ненулевое значение.
- **Восстановление активных транзакций с учетом недоступных файлов.** Эта задача аналогична решаемой во время запуска базы данных. В ходе ее выполнения процесс SMON восстанавливает отказавшие транзакции, которые были пропущены во время восстановления экземпляра или состояния после аварии из-за недоступности файла (файлов) процессу восстановления. Например, файл может размещаться на диске, который был недоступен или не был смонтирован. Когда файл станет доступным, процесс SMON восстановит его.
- **Восстановление экземпляра для отказавшего узла в среде RAC.** В конфигурации Oracle RAC при отказе экземпляра базы данных в кластере (например, если отказывает компьютер, на котором выполнялся экземпляр) какой-то другой узел этого кластера откроет файлы журнала повторения отказавшего экземпляра и проведет восстановление всех его данных.
- **Очистка OBJ\$.** Здесь OBJ\$ — это низкоуровневая таблица словаря данных, которая содержит записи для практически всех объектов (таблиц, индексов, триггеров, представлений и т.д.) базы данных. Во многих случаях в ней находятся записи удаленных объектов или “отсутствующих” объектов, используемых механизмом зависимостей Oracle. Процесс SMON удаляет эти строки, которые больше не нужны.
- **Управление сегментами отмены.** Процесс SMON будет выполнять автоматическое присоединение, отсоединение и сегментов отмены.
- **Отсоединение сегментов отката.** Когда применяется ручное управление сегментами отката (не рекомендуется; вы должны использовать автоматическое управление сегментами отката), администратор базы данных может *отсоединять*, или делать недоступным, сегмент отката, который имеет активные транзакции. Активные транзакции могут пользоваться этим отсоединенным сегментом отката. В таком случае сегмент отката не является действительно отсоединенным; он помечается как “ожидающий отсоединения”. Процесс SMON будет в фоновом режиме периодически пробовать отсоединить его, пока это ему не удастся.

Изложенный материал должен дать вам представление о том, что делает процесс SMON. Он выполняет много других действий, таких как сброс на диск статистической информации мониторинга, отображаемой в представлении DBA_TAB_MODIFICATIONS, запись информации преобразования временной метки, хранящейся в таблице SMON_SCN_TIME, и т.д. Со временем процесс SMON может занимать довольно большую долю ресурсов процессора, и это должно считаться нормальным. Периодически SMON пробуждается (самостоятельно или другими фоновыми процессами) для решения своих задач по наведению порядка в системе.

RECO: распределенное восстановление базы данных

Процесс RECO решает очень узкую задачу: он восстанавливает транзакции, которые остались в подготовленном состоянии вследствие аварии или утери подключения во время *двухфазной фиксации* (two-phase commit — 2PC). 2PC — распределенный протокол, который делает возможной автоматическую фиксацию изменений, влияющих на множество физически разрозненных баз данных. Перед фиксацией он предпринимает попытку максимально сузить временное окно распределенного сбоя. В рамках протокола 2PC среди N баз данных одна из них — обычно (но не всегда) база данных, в которую клиент вошел первоначально — будет координатором. Этот узел будет запрашивать остальные $N-1$ узлов об их готовности к фиксации. Каждый из $N-1$ узлов сообщает свое состояние готовности как YES (да) или NO (нет). Если любой из узлов возвращает NO, производится откат всей транзакции. Если все узлы сообщают YES, то узел-координатор распространяет сообщение о необходимости выполнить фиксацию на каждом из $N-1$ узлов.

Если после того, как узел сообщил о своей готовности к фиксации, но до получения от координатора директивы о необходимости действительного выполнения фиксации происходит сбой в сети или какая-то другая ошибка, транзакция становится *сомнительной распределенной транзакцией*. Протокол 2PC пытается ограничить временное окно, в рамках которого это может случиться, но не в состоянии устранить его полностью. Если отказ возник здесь и сейчас, то транзакция переходит под ответственность процесса RECO. Этот процесс попытается связаться с координатором транзакции для выяснения ее результата. До тех пор, пока он не сделает это, транзакция будет оставаться в незафиксированном состоянии. Когда процессу RECO удастся снова связаться с координатором транзакции, он выполнит либо фиксацию, либо откат транзакции.

Следует отметить, что при наличии каких-то невыполненных транзакций и длительном периоде простоя фиксацию/откат этих транзакций можно производить вручную. Это может понадобиться по той причине, что сомнительные распределенные транзакции могут приводить к *блокированию процессов чтения процессами записи* — одна из ситуаций, которая может произойти в среде Oracle. В этом случае администратор базы данных может связаться с администратором другой базы данных и попросить его проверить состояние таких сомнительных транзакций. Затем администратор базы данных может выполнить фиксацию или откат этих транзакций, освобождая процесс RECO от этой задачи.

СКРТ: процесс выполнения контрольных точек

Вопреки своему названию, процесс выполнения контрольных точек не производит запись контрольных точек (контрольные точки обсуждались в главе 3, в разделе, посвященном журналам повторения) — по большей части это задача процесса DBWn. Процесс СКРТ просто облегчает работу процессу записи контрольных точек, обновляя заголовки файлов данных. Раньше СКРТ считался необязательным, но, начиная с версии Oracle 8.0, он запускается всегда. Поэтому если вы введете команду `ps` в UNIX/Linux, то обычно увидите его (здесь указано “обычно” потому, что в Oracle 12c процесс выполнения контрольных точек может быть запущен внутри потока операционной системы и, следовательно, не отображаться в качестве процесса).

Принято считать, что работа по обновлению заголовков файлов данных информацией о контрольных точках принадлежит процессу LGWR. Однако с увеличением количества файлов и ростом размера базы данных с течением времени эта дополнительная задача становится слишком большой нагрузкой для LGWR. Если бы процесс LGWR был вынужден обновлять десятки, сотни или даже тысячи файлов, то вполне вероятно, что сеансам пришлось бы очень долго дожидаться фиксации этих транзакций. Процесс CKPT освобождает LGWR от этой обязанности.

DBWn: процесс записи блоков базы данных

Процесс записи блоков базы данных (DBWn) — фоновый процесс, отвечающий за запись грязных блоков на диск. Этот процесс будет записывать грязные блоки из буферного кеша обычно с целью освобождения в нем места (для подготовки буферов к чтению других данных) или для продвижения контрольной точки (для перехода на следующую позицию в файле оперативного журнала повторения, с которой СУБД Oracle должна будет начинать чтение при восстановлении экземпляра после аварийного отказа). Как обсуждалось в главе 3, при переключении журнальных файлов выдается сигнал о создании контрольной точки. СУБД Oracle необходимо продвинуться к следующей контрольной точке, чтобы больше не нуждаться в только что заполненном файле оперативного журнала повторения. Если это не будет сделано к моменту, когда потребуется снова использовать файл журнала повторения, отобразится сообщение “checkpoint not complete” (“контрольная точка не завершена”) и придется дожидаться завершения операции.

На заметку! Продвижение журнальных файлов — только один из многих способов создания контрольных точек. Существуют инкрементные контрольные точки, управляемые такими параметрами, как FAST_START_MTTR_TARGET, и другие триггеры, которые вызывают сброс грязных блоков на диск.

Как видите, производительность процесса DBWn может быть критически важной. Если он не будет записывать блоки в свободные буферы (буферы, которые могут повторно использоваться для кэширования других блоков) достаточно быстро, то мы заметим увеличение количества и длительности периодов ожидания Free Buffer Waits (Ожидание появления свободного буфера) и Write Complete Waits (Ожидание окончания записи). Мы можем сконфигурировать более одного процесса DBWn; на самом деле в Oracle 11g допускается настраивать до 36 процессов DBWn (DBW0...DBW9, DBWa...DBWz), а в Oracle 12c — до 100, о чем свидетельствует следующий запрос:

```

EODA@ORA12CR1> select name, description from v$bgprocess
2  where description like 'db writer process%';

```

```

NAME DESCRIPTION
-----

```

```

DBW0 db writer process 0
DBW1 db writer process 1
DBW2 db writer process 2
...
BW97 db writer process 97
BW98 db writer process 98
BW99 db writer process 99
100 rows selected.

```

Большинство систем работают с одним процессом записи блоков базы данных, но в крупных многопроцессорных системах могут применяться несколько таких процессов. Обычно это делается для распределения рабочей нагрузки, связанной с поддержанием в чистом виде большого буферного кеша блоков внутри области SGA, путем сбрасывания грязных (модифицированных) блоков на диск.

В оптимальном случае для записи блоков на диск процесс DBWn использует асинхронный ввод-вывод. В этом случае DBWn собирает пакет блоков, подлежащих записи, и передает их операционной системе. Процессу DBWn не приходится ожидать, пока операционная система действительно запишет блоки; вместо этого процесс приступает к сборке следующего пакета, который должен быть записан. Когда операционная система завершает запись, она асинхронно уведомляет об окончании процесс DBWn. Это позволяет процессу DBWn работать намного быстрее, чем при выполнении всех действий последовательным образом. Позже в разделе “Подчиненные процессы” вы узнаете, как применять подчиненные процессы ввода-вывода для эмуляции асинхронного ввода-вывода на платформах или в конфигурациях, которые его не поддерживают.

Я хотел бы сделать финальное замечание относительно процесса DBWn. Можно почти не сомневаться, что процесс DBWn будет записывать блоки вразброс по всему диску — он производит множество операций записи с разнесением данных. При выполнении обновления вы будете модифицировать индексные блоки, хранящиеся в различных областях диска, и блоки данных, которые также распределены по диску случайным образом. С другой стороны, процесс LGWR выполняет множество последовательных операций записи в журнал повторения. Это важное отличие и одна из причин наличия в Oracle журнала повторения и процесса LGWR, а также процесса DBWn. Запись с разнесением данных происходит значительно медленнее последовательных операций записи. Имея в своем распоряжении грязные блоки буфера SGA и процесс LGWR, выполняющий объемные последовательные записи, которые могут воссоздавать эти грязные буферы, можно увеличить производительность системы. То, что процесс DBWn делает свою медленную работу в фоновом режиме, в то время как процесс LGWR запускает свои более быстрые операции в периоды ожидания пользователя, способствует увеличению общей производительности системы. Это действительно так, несмотря на то, что Oracle может формально выполнять больше операций ввода-вывода, чем необходимо (запись в журнал и в файл данных). *Теоретически* запись в оперативный журнал повторения можно было бы пропустить, если бы во время фиксации вместо этого СУБД Oracle физически записывала измененные блоки на диск. На практике дело обстоит иначе: LGWR выполняет запись информации повторения в оперативные журналы повторения для каждой транзакции, а DBWn сбрасывает блоки базы данных на диск в фоновом режиме.

LGWR: процесс записи журналов

Процесс LGWR отвечает за сброс на диск содержимого буфера журнала повторения, находящегося в области SGA. Это делается в одном из следующих сценариев:

- каждые три секунды;
- когда выдается COMMIT или ROLLBACK;
- когда процесс LGWR получает запрос на переключение журнальных файлов;
- когда буфер повторения заполняется на одну треть или объем содержащихся в нем кешированных данных журнала повторения составляет 1 Мбайт.

По этим причинам наличие огромного (объемом в сотни или тысячи мегабайт) буфера журнала повторения нецелесообразно; Oracle никогда не сможет задействовать его целиком, т.к. практически постоянно производит сброс буфера на диск. В отличие от ввода-вывода с разнесением данных, который должен выполняться процессом DBWn, журналы заполняются с помощью последовательных операций записи. Выполнение таких пакетных операций записи значительно эффективнее выполнения множества операций записи с разнесением данных в различные части файла. Это одна из основных причин применения процесса LGWR и журналов повторения. Повышение эффективности за счет использования последовательного ввода-вывода для записи только изменившихся байтов перевешивает недостаток, связанный с увеличением количества операций ввода-вывода. СУБД Oracle могла бы просто записывать блоки базы данных прямо на диск во время фиксации, но это привело бы к многочисленным операциям ввода-вывода с разнесением данных целых блоков, которые происходили бы значительно медленнее последовательной записи изменений процессом LGWR.

На заметку! Начиная с выпуска Oracle 12c, на многопроцессорных машинах будут запускаться дополнительные процессы LG0 (Log Writer Worker — обработчик процесса записи в журнал), направленные на повышение производительности записи в файлы оперативных журналов повторения.

ARCn: архивный процесс

Задачей процесса ARCn является копирование файла оперативного журнала повторения в другое место, когда он заполняется процессом LGWR. Затем эти архивные файлы журналов повторения можно применять для восстановления с носителя. В то время как оперативный журнал повторения используется для исправления файлов данных при сбое электропитания (когда экземпляр прекращает работу), архивные журналы повторения служат для исправления файлов данных в случае аварийного отказа диска. При выходе из строя дискового устройства, содержащего файл данных, /u01/dbfile/ORA12CR1/system01.dbf, можно обратиться к резервным копиям, созданным на прошлой неделе, восстановить эту старую копию файла, и затем потребовать от базы данных применения всех архивных и оперативных журналов повторения, сгенерированных с момента последнего резервного копирования. Это свяжет упомянутый файл с остальными файлами в базе данных, и мы можем продолжать работу без утери какой-либо информации.

Процесс ARCn обычно копирует файлы оперативного журнала повторения, по крайней мере, в два дополнительных местоположения (избыточность — ключевое средство предотвращения утери данных). Этими местоположениями могут быть диски на локальном компьютере или, что более целесообразно, хотя бы один из них должен находиться на совершенно другом компьютере на случай катастрофического отказа. Во многих ситуациях эти архивные файлы журналов повторения копируются другим процессом на какое-то третье устройство хранения, подобное магнитной ленте. Они могут также пересылаться на другой компьютер для применения к резервной базе данных — вариант обхода отказа, предлагаемый Oracle. Мы вскоре обсудим используемые при этом процессы.

DIAG: процесс диагностики

В прежних выпусках процесс DIAG применялся исключительно в среде RAC. Начиная с Oracle 11g, с его средством ADR (Advanced Diagnostic Repository — репозиторий расширенной диагностики), именно процесс DIAG отвечает за мониторинг общей работоспособности экземпляра и фиксирует информацию, необходимую для обработки отказов экземпляра. Это касается как конфигураций с одним экземпляром, так и многоэкземплярных конфигураций RAC.

FBDA: процесс архивации ретроспективных данных

Процесс архивации ретроспективных данных (Flashback Data Archiver — FBDA) доступен в Oracle 11g Release 1 и последующих версиях. Это ключевой компонент средства архивирования ретроспективных данных — возможности запрашивать данные в том виде, в каком они пребывали в определенный момент времени в прошлом (например, запросить данные таблицы, как они выглядели год назад, пять лет назад и т.д.). Такая возможность достигается благодаря ведению хронологии изменений каждой строки в таблице со временем. Эта хронология, в свою очередь, поддерживается процессом FBDA в фоновом режиме. Указанный процесс функционирует, выполняя свою работу вскоре после фиксации транзакции. Процесс FBDA читает данные отката, сгенерированные этой транзакцией, и производит откат сделанных ею изменений. Затем он регистрирует полученные в результате отката строки (исходные значения) в архиве ретроспективных данных.

DBRM: процесс диспетчера ресурсов

Процесс диспетчера ресурсов базы данных (Database Resource Manager — DBRM) реализует планы ресурсов, которые могут быть сконфигурированы для экземпляра базы данных. Он устанавливает эти планы ресурсов и выполняет различные операции, связанные с принудительным применением/внедрением этих планов ресурсов. Диспетчер ресурсов позволяет администраторам базы данных иметь детализированный контроль над ресурсами, используемыми экземпляром базы данных, приложениями, имеющими доступ к базе данных, либо индивидуальными пользователями, обращающимися к базе данных.

GEN0: процесс выполнения общей задачи

Процесс выполнения общей задачи (General Task Execution Process), как следует из его названия, предоставляет поток выполнения общей задачи. Его главная цель — выгрузить потенциально блокирующую обработку (которая может остановить функционирование процесса) из другого процесса и реализовать ее в фоновом режиме. Например, если главному процессу ASM необходимо выполнить некоторую операцию, блокирующую файл, но данная операция может быть безопасно произведена в фоновом режиме (ASM может благополучно продолжать обработку перед завершением операции), то процесс ASM может запросить у процесса GEN0 выполнение этой операции и уведомление о ее завершении. Это похоже по своей природе на подчиненные процессы, которые описаны далее в главе.

Остальные специализированные процессы

В зависимости от применяемых средств Oracle в системе могут наблюдаться и другие специализированные процессы. Некоторые из них снабжаются кратким описанием их функций.

На заметку! Приложение F (“Background Processes” — “Фоновые процессы”) руководства Oracle Database Reference, доступное по адресу <http://otn.oracle.com/>, содержит полный список фоновых процессов и их функций.

Большинство описанных ранее процессов неизменны — они есть в любом выполняющемся экземпляре Oracle. (Процесс ARCn формально не считается обязательным, но, по моему мнению, должен быть неотъемлемой частью любой производственной базы данных!) Описанные далее процессы не являются обязательными, и они будут действовать в системе только при использовании конкретных функциональных средств. Следующие процессы характерны только для экземпляра базы данных, в котором включена функция ASM (Automatic Storage Management — автоматическое управление пространством хранения), описанная в главе 3.

- **Фоновый процесс автоматического управления пространством хранения (Automatic Storage Management Background — ASMB).** Процесс ASMB выполняется в экземпляре базы данных, применяющем средство ASM. Он отвечает за обмен данными с экземпляром ASM, который управляет пространством хранения, снабжение экземпляра ASM обновленными статистическими сведениями и предоставление такта (heartbeat) экземпляру ASM, чтобы знать, действует ли он по-прежнему.
- **Процесс ребалансировки (ReBALance — RBAL).** Процесс RBAL также действует в экземпляре базы данных, использующем функции ASM. Этот процесс отвечает за обработку запроса ребалансировки (запроса перераспределения) при добавлении/удалении дисков из группы дисков ASM.

Перечисленные ниже процессы можно найти в экземпляре Oracle RAC. Как уже упоминалось, RAC — это конфигурация Oracle, где несколько экземпляров, каждый из которых действует в отдельном узле (как правило, на отдельном физическом компьютере) кластера, могут монтировать и открывать одну базу данных. Это позволяет более чем одному экземпляру получать полный доступ по чтению и записи к одиночному набору файлов базы данных. Основных целей RAC две.

- **Высокая доступность.** В случае применения Oracle RAC отказ одного узла/компьютера в кластере из-за программной, аппаратной или человеческой ошибки позволит остальным узлам продолжить функционирование. База данных останется доступной через другие узлы. При этом вычислительная мощность системы может несколько снизиться, но доступ к базе данных не будет утерян.
- **Масштабируемость.** Вместо того чтобы приобретать все более мощные компьютеры для обработки постоянно возрастающей нагрузки (такой подход называют *вертикальным масштабированием*), технология RAC позволяет добавлять ресурсы в форме большего числа компьютеров в кластере (этот подход носит название *горизонтального масштабирования*). Вместо замены 4-процессорного компьютера 8- или 16-процессорным технология RAC позволяет просто добавить еще один сравнительно недорогой 4-процессорный компьютер (или несколько).

Следующие процессы уникальны для среды RAC. В других местах вы их не увидите.

- **Процесс монитора блокировки (lock monitor — LMON).** Процесс LMON отслеживает все экземпляры в кластере на предмет выявления отказа экземпляра. В случае его обнаружения LMON восстанавливает глобальные блокировки, установленные отказавшим экземпляром. Этот процесс отвечает также за реконфигурирование блокировок и других ресурсов в случае удаления экземпляров либо их добавления в кластер (при их отказе и повторном подключении или при добавлении в кластер новых экземпляров в реальном времени).
- **Процесс демона диспетчера блокировки (lock manager daemon — LMD).** Процесс LMD обрабатывает запросы службы диспетчера блокировки, связанные с глобальным кешем (поддерживающие согласованность буферов блоков между экземплярами). Главным образом этот процесс играет роль посредника, отправляя запросы на предоставление ресурсов в очередь, которая обрабатывается процессами LMSn. Процесс LMD обеспечивает выявление/распознавание глобальных взаимоблокировок и мониторинг тайм-аутов, связанных с блокировками, в глобальной среде. Кроме того, начиная с Oracle 12c, для содействия рабочему потоку могут существовать подчиненные процессы LDDn, порождаемые процессом LMD0.
- **Процесс сервера диспетчера блокировки (lock manager server — LMSn).** Как отмечалось ранее, в среде RAC каждый экземпляр Oracle действует на отдельном компьютере в кластере, и все они имеют доступ по чтению и записи в точности к тому же самому набору файлов базы данных. Для достижения этой цели кэши буферов блоков SGA должны быть согласованными друг с другом. Обеспечение такой согласованности — одна из основных задач процесса LMSn. В ранних выпусках параллельного сервера Oracle (Oracle Parallel Server — OPS) это было реализовано посредством *выгрузки из кеша* (ping). То есть, если узел в кластере нуждался в согласованном по чтению представлении блока, который был заблокирован в монопольном режиме другим узлом, то обмен данными осуществлялся посредством сброса на диск (блок был выгружен из кеша). Такая операция была очень затратной даже для чтения данных. Теперь, благодаря процессу LMSn, этот обмен производится через очень быстрый обмен данными между кэшами по высокоскоростному подключению кластеров. Экземпляр может иметь до десяти процессов LMSn.
- **Процесс блокировки (LCK0).** По своему действию этот процесс подобен описанному ранее процессу LMD, но он обрабатывает запросы на предоставление всех других глобальных ресурсов, отличных от буферов блоков базы данных.

Нижe перечислены общие фоновые процессы, которые присутствуют в большинстве одиночных экземпляров или экземпляров RAC.

- **Процесс генератора процессов (Process Spawner — PSP0).** Процесс PSP0 отвечает за порождение (запуск/создание) разнообразных фоновых процессов. Это процесс, который создает новые процессы/потоки для экземпляра Oracle. Большую часть своей работы он выполняет во время запуска экземпляра.
- **Процесс виртуального хранителя времени (Virtual Keeper of Time — VKTM).** Процесс VKTM реализует согласованные и детализированные часы для экземпляра Oracle. Он отвечает за предоставление времени в формате обычных часов

(читабельном для человека), а также высокоточного таймера (который не обязательно использует формат обычных часов, а скорее является секундомером, отслеживающим очень малые единицы времени), применяемого для измерения продолжительностей и интервалов.

- **Процесс виртуального планировщика для диспетчера ресурсов (Virtual Scheduler for Resource Manager — VKRM).** Планировщик для диспетчера ресурсов, который управляет графиком загрузки процессора и управляемыми процессами с активными планами ресурсов.
- **Процесс координатора управления пространством (Space Management Coordinator — SMC0).** Процесс SMC0 является частью инфраструктуры управления. Он координирует работу средств упреждающего управления пространством базы данных, таких как процессы, обнаруживающие пространство, которое может быть освобождено, и процессы, выполняющие это освобождение.

Служебные фоновые процессы

Эти фоновые процессы являются совершенно необязательными и создаются исключительно по необходимости. Они предоставляют функциональные средства, которые не требуются при повседневной работе базы данных, если только вы сами не используете их, например, очереди заданий, или не применяете использующее их средство, такое как диагностические возможности (в Oracle 10g и последующих версиях).

Такие процессы будут видны в среде UNIX/Linux как любые другие фоновые процессы — в выводе команды `ps`. В примере вывода команды `ps`, приведенном в начале раздела “Специализированные фоновые процессы”, легко заметить наличие следующих элементов.

- Сконфигурированные очереди заданий. Процесс `SJQ0` — это координатор очередей заданий.
- Сконфигурированное средство Advanced Queuing (AQ), о чем свидетельствуют процессы `Qnnn` (процесс очереди AQ) и `QMNC` (процесс монитора AQ).
- Включенное автоматическое управление памятью, что подтверждается процессом диспетчера памяти (Memory Manager — `MMAN`).
- Включенные средства управления/диагностики Oracle, о чем свидетельствуют процессы монитора управляемости (Manageability Monitor — `MMON`) и облегченного монитора управляемости (Manageability Monitor Light — `MMNL`).

Давайте рассмотрим процессы, которые могут присутствовать в зависимости от применяемых функциональных средств.

Процессы `SJQ0` и `Jnnn`: очереди заданий

В первом выпуске Oracle 7.0 была предложена возможность репликации в форме объекта базы данных, называемого *моментальной копией* (snapshot). Очереди заданий были внутренним механизмом, с помощью которого эти моментальные копии обновлялись, или делались текущими.

Процесс очереди отслеживал таблицу заданий, содержащую информацию о том, когда необходимо обновлять различные моментальные копии в системе.

В версии Oracle 7.1 эта возможность стала доступной для всех через пакет базы данных по имени DBMS_JOB. Таким образом, процесс, который в версии 7.0 был связан исключительно с обработкой моментальных копий, в версии 7.1 и последующих версиях стал “очередью заданий”. С течением времени имена параметров, управляющих поведением очереди заданий, изменились с SNAPSHOT_REFRESH_INTERVAL и SNAPSHOT_REFRESH_PROCESSES на JOB_QUEUE_INTERVAL и JOB_QUEUE_PROCESSES. В текущем выпуске для настройки пользователем доступен только параметр JOB_QUEUE_PROCESSES.

В системе можно иметь до 1000 процессов очередей заданий с именами J000...J999. Эти процессы интенсивно используются при репликации в качестве составной части процесса обновления материализованного представления. Поточковая репликация (доступная в Oracle9i Release 2 и последующих версиях) применяет для репликации средство AQ и, следовательно, не использует процессы очередей заданий. Разработчики также часто применяют очереди заданий для планирования разовых (фоновых) заданий либо повторяющихся заданий вроде фоновой отправки электронной почты или обработки длительно выполняющегося пакетного процесса, запускаемого в фоновом режиме. Прodelывая часть работы в фоновом режиме, можно создавать видимость значительно более быстрого выполнения длительных задач (нетерпеливому пользователю будет казаться, что задание выполняется быстрее, хотя оно может быть еще не завершено). Это подобно тому, что Oracle делает с помощью процессов LGWR и DBWn; большую часть своей работы они производят в фоновом режиме, поэтому нам не придется дожидаться, пока упомянутые процессы завершат все задачи в реальном времени.

Процессы Jnnn (где nnn — некоторое число) очень похожи на разделяемый сервер, но с рядом аспектов выделенного сервера. Они являются разделяемыми в том смысле, что обрабатывают одно задание за другим, но управляют памятью в манере, больше подобной способу, который используется выделенным сервером (их память UGA размещена в области PGA, а не SGA). Каждый процесс очереди заданий времени будет выполнять только одно задание за раз, одно за другим, вплоть до завершения. Именно поэтому для одновременного выполнения заданий может потребоваться множество процессов. Не существует никакого метода организации многопоточного или приоритетного выполнения заданий. После того как задание запущено, оно будет выполняться до полного завершения (или отказа).

Вы заметите, что с течением времени процессы Jnnn появляются и исчезают, т.е. если система сконфигурирована с поддержкой 1000 таких процессов, то не все 1000 процессов будут запускаться вместе с базой данных. Вместо этого будет запущен единственный процесс координатора очереди заданий (CJQ0), который будет запускать необходимые процессы Jnnn по мере появления заданий, подлежащих выполнению, в таблице очереди заданий. Выполнив свою работу и обнаружив, что новая работа отсутствует, процессы Jnnn начнут завершаться и исчезать. Таким образом, если вы запланируете большинство своих заданий на запуск в 2:00, когда в помещении никого нет, то эти процессы Jnnn можно вообще не заметить.

QMNC и Qnnn: расширенные очереди

Процесс QMNC играет для таблиц AQ ту же роль, которую процесс CJQ0 играет для таблицы заданий. Он отслеживает расширенные очереди и оповещает процессы извлечения из очереди, ожидающие сообщения, о том, что сообщение стало доступ-

ным. Процессы QMNC и Qnnn отвечают также за *распространение* в очереди — т.е. возможность перемещения сообщения, которое было добавлено в очередь в одной базе данных, в очередь в другой базе данных для его последующего извлечения.

Процессы Qnnn играют для процесса QMNC ту же роль, которую процессы Jnnn играют для процесса CJQ0. Процесс QMNC уведомляет их о работе, которая должна быть выполнена, и они занимаются этой работой.

Процессы QMNC и Qnnn представляют собой необязательные фоновые процессы. Параметр AQ_TM_PROCESSES указывает создание до 40 таких процессов с именами Qnnn (где nn — число 0..15 или буква a..z) и одиночного процесса QMNC.

В отличие от процессов Jnnn, применяемых очередями заданий, процессы Qnnn являются постоянными. Если вы установите параметр AQ_TM_PROCESSES в 10, то увидите десять процессов Q0nn и процесс QMNC при запуске базы данных, а также на протяжении всего времени жизни экземпляра.

СУБД Oracle автоматически корректирует количество процессов очередей и по этой причине вам редко придется устанавливать AQ_TM_PROCESSES вручную. Если же все же установите этот параметр, то Oracle по-прежнему будет корректировать количество порождаемых процессов и использовать значение AQ_TM_PROCESSES в качестве минимального числа процессов, подлежащих созданию.

На заметку! Начиная с версии Oracle 12c, имеется процесс координатора процессов расширенных очередей (Advanced Queue Process Coordinator — AQPC). Он предназначен для создания и управления главными процессами расширенной очереди (запуск, останов и т.д.). Статистическую информацию, связанную с этим процессом, можно запросить из представления GV\$AQ_BACKGROUND_COORDINATOR.

EMNn: процессы монитора событий

Процесс EMNn — это часть архитектуры AQ. Он служит для уведомления подписчиков очереди о сообщениях, в которых они могут быть заинтересованы. Такое уведомление производится асинхронно. Существуют функции OCI (Oracle Call Interface — интерфейс уровня вызовов Oracle), позволяющие зарегистрировать обратный вызов для уведомления о сообщении. Обратный вызов — это функция программы OCI, которая будет автоматически вызываться всякий раз, когда интересующее сообщение доступно в очереди. Фоновый процесс EMNn предназначен для уведомления подписчика. Процесс EMNC запускается автоматически при выдаче первого уведомления для экземпляра. Затем приложение может выдать явную команду `message_receive(dequeue)`, чтобы извлечь сообщение.

MMAN: диспетчер памяти

Этот процесс доступен в Oracle 10g и последующих версиях и применяется средством автоматической установки размера области SGA. Процесс MMAN координирует установку и изменение размеров компонентов разделяемой памяти (стандартного буферного пула, разделяемого пула, пула Java и большого пула).

MMON, MMNL и Mnnn: мониторы управляемости

Эти процессы используются для заполнения автоматического хранилища информации о рабочей нагрузке (Automatic Workload Repository — AWR) — функци-

онального средства, которое доступно начиная с версии Oracle 10g. Процесс MMNL сбрасывает статистические сведения из области SGA в таблицы базы данных на основе расписания. Процесс MMON служит для автоматического обнаружения проблем с производительностью базы данных и реализации новых функций самонастройки. Процессы Mnnn аналогичны процессам Jnnn или Qnnn применительно к очередям заданий; процесс MMON будет запрашивать у этих подчиненных процессов выполнения работы от его имени. По своей природе процессы Mnnn являются временными — они будут появляться и исчезать по мере надобности.

CTWR: процесс отслеживания изменений

Это необязательный процесс в Oracle 10g и последующих версиях. Он отвечает за поддержку файла отслеживания изменений, как было описано в главе 3.

RVWR: процесс записи данных восстановления

Этот процесс (доступный в Oracle 10g и последующих версиях) обеспечивает сопровождение *предшествующих* образов блоков в области быстрого восстановления (Fast Recovery Area), описанной в главе 3, которая используется командой FLASHBACK DATABASE.

DMnn/DWnn: главные/рабочие процессы Data Pump

Data Pump — это средство, добавленное в версии Oracle 10g Release 1. Оно было спроектировано как полная замена унаследованных процессов экспорта/импорта. Data Pump выполняется целиком на сервере и API-интерфейсом для взаимодействия с ним служит PL/SQL. Поскольку Data Pump выполняется на сервере, была добавлена поддержка выполнения разнообразных операций Data Pump. Главный процесс Data Pump (DMnn) собирает весь ввод от клиентских процессов (получающих ввод через API-интерфейс) и затем координирует рабочие процессы (DWnn), которые делают фактическую работу — процессы DMnn осуществляют действительную обработку метаданных и данных.

TMON/TT00: процесс монитора перемещения и подчиненный процесс перемещения данных отмены

Начиная с версии Oracle 12c, во время старта экземпляра автоматически запускаются два процесса, связанные с Data Guard: процесс монитора перемещения (transport monitor process — TMON) и подчиненный процесс перемещения данных отмены (redo transport slave — TT00). Процесс TMON запускается и следит за несколькими процессами TT00. Процессы TT00 применяются для информирования процесса LGWR о необходимости генерации такта данных отмены (heartbeat redo). Вы можете видеть эти процессы запущенными, даже если вообще не внедряли Data Guard. Не переживайте по их поводу, просто помните о том, что они существуют и будут использоваться, если вы развернете средство Data Guard.

Если же средство Data Guard внедрено, будет запущено несколько дополнительных процессов, предназначенных для упрощения доставки информации отмены из одной базы данных в другую и для ее применения. За подробными сведениями обращайтесь в руководство *Data Guard Concepts and Administration* (Концепции и администрирование Data Guard), предлагаемое Oracle.

Остальные служебные фоновые процессы

Итак, является ли приведенный список исчерпывающим? Нет, никоим образом. Существует много других процессов подобного рода, зависящих от того, какие функциональные средства реализованы. Например, есть также процессы применения и перехвата потоков Streams, которые присутствуют в случае внедрения таких продуктов, как Oracle GoldenGate, Oracle XStream, Oracle Streams и т.д. Тем не менее, приведенный перечень охватывает большинство наиболее часто встречающихся фоновых процессов.

Подчиненные процессы

Теперь можно приступить к рассмотрению последнего класса процессов Oracle — *подчиненных* процессов. Различают два типа подчиненных процессов Oracle: подчиненные процессы ввода-вывода и подчиненные процессы параллельных запросов.

Подчиненные процессы ввода-вывода

Подчиненные процессы ввода-вывода используются для эмуляции асинхронного ввода-вывода в системах или устройствах, которые его не поддерживают. Например, лентопротяжные устройства (общеизвестно медленные) не поддерживают асинхронный ввод-вывод. За счет применения подчиненных процессов ввода-вывода для лентопротяжных устройств можно имитировать то, что операционная система обычно предлагает для дисководов. Как и при настоящем асинхронном вводе-выводе, процесс, производящий запись на устройстве, помещает большой объем данных в пакет и передает его для записи. После успешного завершения записи данных записывающий процесс (в этом случае подчиненный процесс ввода-вывода, а не операционная система) выдает сигнал процессу, инициировавшему запись, который удаляет этот пакет из своего списка данных, предназначенных для записи. Такой подход позволяет значительно увеличить пропускную способность, поскольку подчиненные процессы ввода-вывода ожидают завершения работы медленного устройства, а в это время вызывающий их процесс выполняет другую важную работу по сбору данных для следующей записи.

Подчиненные процессы ввода-вывода используются в паре мест базы данных Oracle. Процессы DBWn и LGWR могут их применять для эмуляции асинхронного ввода-вывода, а процесс RMAN будет их использовать при записи на магнитную ленту.

Для управления применением подчиненных процессов ввода-вывода служат два параметра.

- **BACKUP_TAPE_IO_SLAVES.** Этот параметр указывает, будет ли процесс RMAN использовать подчиненные процессы ввода-вывода для резервирования, копирования или восстановления данных с магнитной ленты. Поскольку параметр **BACKUP_TAPE_IO_SLAVES** предназначен для работы с ленточными устройствами, а такие устройства могут быть доступны только одному процессу в каждый момент времени, его значение является булевским, а не числовым, равным количеству применяемых подчиненных процессов, как можно было бы предположить. Процесс RMAN запустит столько подчиненных процессов, сколько требуется для конкретного количества используемых физических устройств. Когда параметр **BACKUP_TAPE_IO_SLAVES** установлен в **TRUE**, для записи

или чтения с ленточного устройства используется подчиненный процесс ввода-вывода. Если же он установлен в FALSE (по умолчанию), то подчиненные процессы ввода-вывода при резервном копировании не применяются. Вместо этого к ленточному накопителю будет обращаться процесс выделенного сервера, задействованный в резервном копировании.

- **DBWR_IO_SLAVES.** Этот параметр указывает количество подчиненных процессов ввода-вывода, используемых процессом DBW0. Процесс DBW0 и его подчиненные процессы всегда выполняют запись на диск грязных блоков буферного кеша. По умолчанию значением DBWR_IO_SLAVES является 0 и подчиненные процессы ввода-вывода не применяются. Обратите внимание, что если вы установите этот параметр в ненулевое значение, то процессы LGWR и ARCn будут использовать также собственные подчиненные процессы ввода-вывода — для LGWR и ARCn разрешено иметь до четырех подчиненных процессов ввода-вывода.

Подчиненные процессы ввода-вывода процесса DBWn отображаются с именем I1nn, а подчиненные процессы ввода-вывода процесса LGWR — с именем I2nn.

Pnnn: серверы выполнения параллельного запроса

В версии Oracle 7.1.6 была введена возможность параллельного запроса базы данных. Она позволяет создавать для SQL-оператора, такого как SELECT, CREATE TABLE, CREATE INDEX, UPDATE и т.д., план выполнения, состоящий из *множества* планов выполнения, которые могут быть сделаны одновременно. Результаты выполнения всех этих планов объединяются в один результат большего размера. Цель заключается в сокращении времени выполнения операции по сравнению с ее последовательным выполнением. Например, предположим, что существует действительно большая таблица, распределенная по десяти разным файлам. В нашем распоряжении имеется 16 процессоров и требуется выполнить произвольный запрос к этой таблице. Полезно разбить план выполнения на 32 небольших фрагмента и действительно в полной мере задействовать вычислительные возможности компьютера, вместо того чтобы применять только один процесс для последовательного чтения и обработки всех данных.

Когда используется параллельный запрос, вы увидите процессы с именами Pnnn — это собственно серверы выполнения параллельного запроса. Во время обработки параллельного оператора ваш серверный процесс будет известен как *координатор параллельного запроса*. На уровне операционной системы его название не меняется, но если вы встретите в документации по параллельным запросам ссылки на процесс координатора, то знайте, что это просто ваш исходный серверный процесс.

До выхода Oracle 12c стандартное количество запускаемых серверов параллельного выполнения было нулевым. Это поведение можно было изменить, указывая ненулевое значение для параметра PARALLEL_MIN_SERVERS. Начиная с Oracle 12c, ваш экземпляр будет автоматически создавать несколько серверов параллельного выполнения.

Причина в том, что параметр PARALLEL_MIN_SERVERS установлен в ненулевое значение (выведенное из выражения $CPU_COUNT * PARALLEL_THREADS_PER_CPU * 2$). Например, на моем двухпроцессорном компьютере (параметр CPU_COUNT равен 2 и параметр PARALLEL_THREADS_PER_CPU также равен 2) можно наблюдать следующие восемь функционирующих серверов параллельного выполнения:

```
$ ps -ef | grep ora_p00 | grep -v grep
```

| | | | | | | |
|--------|-------|---|---|---------|----------|-------------------|
| oracle | 31086 | 1 | 0 | Apr06 ? | 00:00:03 | ora_p000_ORA12CR1 |
| oracle | 31088 | 1 | 0 | Apr06 ? | 00:00:05 | ora_p001_ORA12CR1 |
| oracle | 31104 | 1 | 0 | Apr06 ? | 00:00:02 | ora_p002_ORA12CR1 |
| oracle | 31106 | 1 | 0 | Apr06 ? | 00:00:02 | ora_p003_ORA12CR1 |
| oracle | 31108 | 1 | 0 | Apr06 ? | 00:00:02 | ora_p004_ORA12CR1 |
| oracle | 31110 | 1 | 0 | Apr06 ? | 00:00:02 | ora_p005_ORA12CR1 |
| oracle | 31112 | 1 | 0 | Apr06 ? | 00:00:02 | ora_p006_ORA12CR1 |
| oracle | 31114 | 1 | 0 | Apr06 ? | 00:00:02 | ora_p007_ORA12CR1 |

Совет. Параллельная обработка подробно рассматривается в главе 14.

Резюме

В этой главе были раскрыты файлы, применяемые Oracle, начиная с простого, но важного файла параметров и заканчивая файлами данных, файлами журналов повторения и т.д. Мы взглянули на структуры памяти, используемые Oracle как в серверных процессах, так и в области SGA. Мы показали, что разные серверные конфигурации, такие как подключения в режимах выделенного и разделяемого серверов, оказывают сильное влияние на то, как память применяется системой. Наконец, мы рассмотрели процессы (или потоки — в зависимости от операционной системы), которые позволяют СУБД Oracle делать то, что она делает. Теперь можно приступать к исследованию реализации ряда других функциональных средств Oracle, таких как блокировка, управление параллельным доступом и транзакции.

ГЛАВА 6

Блокировка и защелкивание данных

Одна из основных проблем при разработке многопользовательских приложений, управляемых базами данных, связана с доведением до максимума степени параллельного доступа с одновременным обеспечением для каждого пользователя возможностей читать и модифицировать данные в согласованной манере. Механизмы блокировки, которые позволяют добиться этого, являются ключевыми средствами любой базы данных, и СУБД Oracle лучшая в плане их предоставления. Тем не менее, реализация Oracle этих функциональных средств специфична для Oracle — подобно тому, как реализация SQL Server характерна именно для SQL Server — и задача обеспечения того, что приложение правильно использует эти механизмы при манипулировании данными, целиком возлагается на разработчика. Если вы проигнорируете это, то приложение будет вести себя непредсказуемым образом, что неизбежно приведет к нарушению целостности данных (как демонстрировалось в главе 1).

В этой главе мы подробно рассмотрим, как Oracle блокирует данные (например, строки в таблицах) и разделяемые структуры данных (такие как компоненты SGA). Мы исследуем уровень детализации, до которого Oracle блокирует данные, и что это означает для разработчика. Когда это уместно, я буду сравнивать схему блокировки Oracle с другими популярными реализациями, главным образом, чтобы развеять миф о том, что блокировка на уровне строк увеличивает накладные расходы; в действительности накладные расходы добавляются, только если это делает реализация. В следующей главе мы продолжим обсуждение этой темы и исследуем технологии многоверсионности, а также взаимодействие с ними стратегий блокировки.

Понятие блокировки

Блокировка (lock) — это механизм, применяемый для регулирования параллельного доступа к разделяемому ресурсу. Обратите внимание на использование термина “разделяемый ресурс”, а не “строка базы данных”. Действительно, Oracle блокирует данные таблиц на уровне строк, но также применяет блокировки на многих других уровнях для предоставления параллельного доступа к разнообразным ресурсам.

Например, пока хранимая процедура выполняется, она сама блокируется в режиме, который позволяет другим запускать ее, но не разрешает другому пользователю каким-либо образом изменять экземпляр этой хранимой процедуры. Блокировки используются в базе данных для обеспечения параллельного доступа к таким разделяемым ресурсам, в то же самое время поддерживая целостность и согласованность данных.

В базе данных с единственным пользователем применение блокировок не обязательно. По определению есть только один пользователь, изменяющий информацию. Но когда несколько пользователей получают доступ и модифицируют данные или структуры данных, критически важно располагать механизмом, предотвращающим одновременное изменение одного и того же фрагмента информации. Именно этой цели и служит блокировка.

Очень важно понимать, что способов реализации блокировок доступно столько же, сколько существует баз данных. Наличие опыта работы с моделью блокировки в одной системе управления реляционными базами данных (СУРБД) вовсе не означает, что вы знаете о блокировке абсолютно все. Например, прежде чем вплотную заняться Oracle, я использовал другие СУРБД, в том числе Sybase, Microsoft SQL Server и Informix. Все упомянутые СУРБД предоставляют механизмы блокировки для управления параллельной обработкой, но реализация блокировки в каждой из них характеризуется глубокими и фундаментальными отличиями.

Чтобы проиллюстрировать сказанное, я кратко опишу свое продвижение от разработчика SQL Server до пользователя Informix и, наконец, до разработчика Oracle. Это происходило много лет назад, и стойкие приверженцы SQL Server скажут: “Но теперь мы располагаем блокировкой на уровне строк”. Это действительно так: в настоящее время SQL Server допускает применение блокировки на уровне строк, но способ ее реализации совершенно отличается от способа, как это сделано в Oracle. Сравнить их — все равно, что сравнивать яблоки с апельсинами, и это является ключевым моментом.

Мне, как программисту SQL Server, вряд ли пришло бы в голову учитывать возможность одновременной вставки данных таблицу сразу несколькими пользователями. Такая ситуация редко происходила в этой базе данных. В те времена база данных SQL Server предоставляла только блокировку на уровне страниц, и поскольку все данные, как правило, должны были вставляться в последнюю страницу неклассифицированных таблиц, параллельная вставка данных двумя пользователями просто не должна была происходить.

На заметку! Кластеризованная таблица SQL Server (таблица, которая имеет кластеризованный индекс) в определенном отношении похожа, но существенно отличается от кластера Oracle. В SQL Server обычно поддерживалась только блокировка на уровне страниц (блоков), и если бы каждая вставляемая строка помещалась в “конец” таблицы, то в этой базе данных никогда не возникали бы ситуации параллельных вставок или параллельных транзакций. В SQL Server кластеризованный индекс применялся для вставки строк во всех местах таблицы, в порядке сортировки, определенном ключом кластера, что в результате улучшало степень параллелизма в базе данных.

Те же самые проблемы возникали при параллельных обновлениях (из-за того, что в действительности операция UPDATE представляла собой последовательно выполняемые операции DELETE и INSERT). Вероятно, именно по этой причине SQL

Server по умолчанию производит фиксацию или откат немедленно после выполнения каждого отдельного оператора, жертвуя целостностью транзакций в попытке достижения более высокой степени параллельной обработки.

Таким образом, в большинстве случаев из-за блокировки на уровне страниц несколько пользователей не могли одновременно модифицировать одну и ту же таблицу. Ситуация усугублялась тем фактом, что во время изменения таблицы многие запросы к ней также блокировались. Попытка запросить таблицу и получить доступ к странице, заблокированной операцией обновления, приводила к длительному ожиданию. Механизм блокировки был настолько несовершенным, что предоставление поддержки для транзакций, выполняющихся дольше одной секунды, было чрезвычайно опасным — вся база данных оказывалась замороженной. В результате я обзавелся массой скверных привычек. Я усвоил, что транзакции — это “плохо”, и их следует фиксировать как можно скорее, никогда не удерживая блокировки на данных. За параллельную обработку приходилось платить снижением целостности данных. Можно было обеспечить либо их корректность, либо скорость получения. Я пришел к заключению, что добиться того и другого было нельзя.

Когда я перешел на Informix, ситуация улучшилась, но не намного. Если я не забывал создавать таблицу с включенной блокировкой на уровне строк, то два пользователя действительно получали возможность одновременной вставки данных в эту таблицу. К сожалению, за такой параллелизм приходилось платить высокую цену. Реализация блокировок на уровне строк в Informix требовала больших затрат как времени, так и памяти. Их установка и снятие (освобождение) занимало значительное время, и каждая блокировка потребляла реальную память. Кроме того, общее количество доступных в системе блокировок должно было вычисляться до запуска базы данных. Если вы превышали это количество, значит, вам не повезло. Вследствие этого большинство таблиц в любом случае создавались с блокировкой на уровне страниц и подобно SQL Server блокировки на уровне строк и страниц стопорили прохождение запросов. В результате я в очередной раз обнаружил, что хочу выполнять фиксацию как можно быстрее. Плохие привычки, которые я приобрел при использовании SQL Server, просто укрепились, к тому же я научился считать блокировку очень дефицитным ресурсом — чем-то особо желанным. Я усвоил, что во избежание установки слишком большого числа блокировок и затормаживания системы — что случалось многократно — необходимо вручную повышать блокировки с уровня строк до уровня таблиц.

Когда я начал работать с Oracle, то не слишком утруждал себя чтением руководства для выяснения особенностей работы блокировки в этой конкретной базе данных. В конце концов, я уже довольно долго работал с базами данных и считался экспертом в этой области (в дополнение к Sybase, SQL Server и Informix я применял Ingress, DB2, Gupta SQLBase и ряд других СУБД). Я попал в ловушку, будучи уверенным, что коль скоро мне известно, как *должны* работать различные функции, то, конечно же, именно так они и *будут* работать. Я крупно ошибался.

Степень своего заблуждения я выяснил во время выполнения сравнительных тестов. На заре развития СУБД (в 1992/1993 году) перед действительно крупными приобретениями обычно проводилось сравнительное тестирование программного обеспечения, предлагаемого разными поставщиками, чтобы выяснить, какое из них способно решать задачи быстрее и проще, предлагая при этом более широкую функциональность.

Мы сравнивали работу Informix, Sybase, SQL Server и Oracle. Первой испытывалась СУБД Oracle. Специалисты из Oracle прибыли в наш офис, прочли спецификации тестов и начали подготовку. Первым делом я заметил, что специалисты из Oracle собирались использовать таблицу базы данных для фиксации показателей времени, хотя мы были намерены применять при работе десятки подключений, каждому из которых предстояло часто выполнять вставку и обновление данных в этой журнальной таблице. Более того, во время тестирования они собирались также *читать* журнальную таблицу! Будучи тактичным, я отозвал одного специалиста из Oracle в сторону, чтобы поинтересоваться, не сошли ли они с ума — зачем им понадобилось сознательно вводить в систему еще одну точку конкуренции? Разве тестовые процессы не будут стремиться выполнять свои операции с этой единственной таблицей последовательно? Разве результаты сравнительного тестирования не окажутся искаженными из-за попыток чтения из этой таблицы, в то время как другие процессы будут интенсивно ее изменять? Зачем нужно привносить все эти дополнительные блокировки, которыми придется управлять? У меня возникли десятки вопросов вроде: “Почему вы вообще подумали об этом?” Специалисты из Oracle решили, что я немного не в себе. Они оставались в таком убеждении до момента, когда я открыл окно в SQL Server или Informix и продемонстрировал последствия вставки данных в таблицу двумя пользователями или попытки запроса таблицы в то время, когда кто-то другой вставлял в нее строки (запрос вообще не возвращал строк). Отличия между технологиями, которые используются в Oracle и практически во всех остальных базах данных, были разительными — они отличались буквально как день и ночь.

Не стоит и говорить, что ни специалисты Informix, ни специалисты SQL Server не особо стремились применять подход с журнальной таблицей базы данных при проведении своих тестов. Они предпочли записывать полученные показатели времени в плоские файлы, поддерживаемые операционной системой. Специалисты Oracle покинули наш офис, вооруженные знанием того, как преуспеть в конкуренции с Sybase, SQL Server и Informix. Нужно просто поинтересоваться: “Сколько строк в секунду возвращает используемая вами база данных, когда данные заблокированы?” — и этого будет достаточно. Мораль этой истории двояка. Во-первых, *все базы данных являются фундаментально разными*. Во-вторых, при проектировании приложения для новой платформы базы данных не должны делаться какие-либо предположения о том, как эта база данных работает. *К каждой новой СУБД нужно относиться так, словно ранее вы вообще не имели дело с базами данных*. Действия, выполняемые в одной базе данных, окажутся либо необязательными, либо неприемлемыми в другой базе данных.

Применяя Oracle, вы усвоите следующие моменты.

- Транзакции — это то, для чего предназначены базы данных. Они являются положительным аспектом.
- Фиксация должна быть отложена до подходящего момента. Ее не следует выполнять быстро, чтобы избежать перегрузки системы, т.к. длительные или крупные транзакции вовсе ее не перегружают. Правило предусматривает *выполнение фиксации тогда, когда она должна быть сделана, но не раньше*. Размер транзакций должен диктоваться бизнес-логикой. (Интересное стороннее замечание: сегодня утром я сформулировал это правило фиксации — выполнять фиксацию тогда, когда она должна быть сделана, но не раньше — на сайте <http://asktom.oracle.com> вероятно в миллионный раз. Кое-что в этом мире никогда не меняется.)

- Блокировки данных должны удерживаться столько, сколько необходимо. Они являются инструментами, предназначенными для использования, а не тем, чего следует избегать. Блокировки — не дефицитный ресурс. Напротив, данные следует блокировать настолько долго, насколько это нужно. Не являясь дефицитными, блокировки могут предотвратить изменение информации другими сеансами.
- С блокировкой на уровне строк в Oracle не связаны какие-либо накладные расходы — вообще. Независимо от того, сколько есть блокировок строк, 1 или 1 000 000, объем ресурсов, выделенных для блокирования этой информации, будет тем же самым. Конечно, вам придется проделывать намного больше работы, модифицируя 1 000 000 строк вместо 1 строки, но количество ресурсов, необходимых для блокировки 1 000 000 строк, будет таким же, как при блокировке 1 строки; это фиксированная константа.
- Ранг блокировки никогда не должен повышаться (примером может быть изменение блокировки таблицы вместо блокировки строк) из-за того, что так было бы “лучше для системы”. В Oracle это не будет лучше для системы — никакие ресурсы не экономятся. Существуют ситуации для использования блокировок таблиц, такие как пакетный процесс, при выполнении которого хорошо известно, что будет обновляться вся таблица и нежелательно, чтобы другие сеансы блокировали отдельные строки в ней. Но не следует применять блокировку таблицы лишь для того, чтобы облегчить системе работу путем устранения необходимости в выделении блокировок строк; блокировка таблицы используется для обеспечения доступа ко всем ресурсам, в которых нуждается пакетная программа.
- Параллелизма и согласованности можно достигнуть одновременно. В любой момент времени данные можно извлечь быстро и точно. Процессы чтения данных *не* блокируются процессами записи данных. Процессы записи данных *не* блокируются процессами чтения данных. В этом состоит одно из фундаментальных отличий между Oracle и большинством других реляционных баз данных.

По мере изложения материала в этой и следующей главах, я подкреплю эти утверждения дополнительными аргументами.

Проблемы блокировки

Прежде чем мы обсудим разнообразные типы блокировок, применяемые в Oracle, полезно взглянуть на определенные проблемы блокировки, многие из которых являются результатом плохо спроектированных приложений, неправильно использующих (или вообще не использующих) механизмы блокировки базы данных.

Потерянные обновления

Потерянное обновление (lost update) — классическая проблема базы данных. Действительно, эта проблема возникает во всех многопользовательских компьютерных средах. Выражаясь просто, потерянное обновление происходит, когда происходят следующие события в приведенном здесь порядке.

1. Транзакция в сеансе Session1 извлекает (запрашивает) строку данных в локальную память и отображает ее конечному пользователю, User1.
2. Другая транзакция в сеансе Session2 извлекает ту же самую строку, но отображает данные другому конечному пользователю, User2.
3. С помощью своего приложения пользователь User1 изменяет эту строку и вынуждает приложение обновить базу данных и произвести фиксацию. На этом транзакция сеанса Session1 завершается.
4. Пользователь User2 также модифицирует эту строку и вынуждает приложение обновить базу данных и выполнить фиксацию. Транзакция сеанса Session2 завершена.

Описанный сценарий называют *потерянным обновлением*, потому что все изменения, сделанные на шаге 3, будут утеряны. Например, рассмотрим экран обновления сведений о сотрудниках, который позволяет пользователю изменять адрес, номер рабочего телефона и т.д. Само приложение является очень простым: небольшой экран поиска, который генерирует список сотрудников, а затем предоставляет возможность вывода подробных сведений о каждом из них. Выглядит пустячным делом. Таким образом, мы создаем приложение безо всякой блокировки со своей стороны, которое выполняет только простые команды SELECT и UPDATE.

Затем конечный пользователь (User1) переходит на экран подробных сведений, изменяет адрес, щелкает на кнопке сохранения и получает сообщение о том, что обновление было успешно выполнено. Все прекрасно за исключением того, что когда на следующий день пользователь User1 проверяет запись, чтобы отправить декларацию о доходах, в ней по-прежнему указан старый адрес. Как такое могло произойти? К сожалению, подобная ситуация может возникнуть очень легко. В этом случае другой конечный пользователь (User2) запросил ту же самую запись сразу после пользователя User1 — после того, как пользователь User1 прочитал данные, но перед тем, как он их изменил. Затем, после запроса пользователем User2 данных пользователь User1 произвел обновление, получил подтверждение и даже повторил запрос, чтобы удостовериться в успешном обновлении данных. Однако затем пользователь User2 обновил поле номера рабочего телефона и щелкнул на кнопке сохранения, совершенно не подозревая о том, что он перезаписал старыми данными измененное пользователем User1 поле адреса! Причина возникновения такой проблемы *в этом случае* связана с тем, что разработчик реализовал свое приложение так, что при обновлении какого-то одного поля все поля данной записи также обновляются (всего лишь потому, что проще обновить все столбцы, чем выяснять, какие столбцы изменились, и обновлять только их).

Обратите внимание, что для появления описанной ситуации даже не требуется, чтобы пользователи User1 и User2 работали с записью в точности одновременно. Нужно чтобы они работали с записью *приблизительно* в одно и то же время.

Мне неоднократно приходилось наблюдать эту проблему, когда программисты, занимающиеся созданием графических пользовательских интерфейсов, которые имели лишь небольшой опыт работы с базами данных или вообще не имели его, получали задание написать приложение базы данных. Они бегло знакомились с операциями SELECT, INSERT, UPDATE и DELETE и принимались за создание приложения. Когда результирующее приложение ведет себя описанным образом, оно полностью подрывает доверие пользователя к себе, особенно с учетом того, что этот недостаток

проявляется случайно, непредсказуемо и совершенно не поддается воспроизведению в контролируемой среде (вырабатывая у разработчика уверенность в том, что это должна быть ошибка пользователя).

Многие инструменты, подобные Oracle Forms и APEX (Application Express — инструмент, который применялся для создания веб-сайта AskTom), прозрачным образом предотвращают такое поведение, гарантируя невозможность изменения записи с момента ее запроса и блокируя ее вплоть до внесения всех необходимых изменений (*оптимистическая блокировка*). Но многие другие (вроде написанных вручную программ на Visual Basic или Java) этого не делают. Обеспечение защиты данных какими-то средствами, работающими “за кулисами”, или же самими разработчиками сводится к использованию одного из двух типов блокировки: *пессимистической* или *оптимистической*.

Пессимистическая блокировка

Метод пессимистической блокировки будет вступать в действие непосредственно перед тем, как пользователь изменяет значение на экране. Например, блокировка строки будет установлена, как только пользователь сообщит о своем намерении обновить конкретную строку, которая выбрана им и отображается на экране (скажем, щелчком на какой-нибудь кнопке). Эта блокировка строки будет *постоянной* до тех пор, пока приложение не применит изменения, внесенные пользователями, к соответствующей строке и не выполнит фиксацию.

Пессимистическая блокировка приемлема только в среде с сохранением состояния или подключенной среде — т.е. среде, в которой приложение имеет постоянное подключение к базе данных и только один пользователь использует это подключение, по крайней мере, в течение времени существования транзакции. В начале и середине 1990-х годов, когда широкое распространение получили клиент-серверные приложения, этот способ блокировки был преобладающим. Каждое приложение должно было получать прямое подключение к базе данных, применяемое исключительно этим экземпляром приложения. В настоящее время такой метод подключения с сохранением состояния стал менее распространенным (хотя и не исчез полностью), особенно после появления серверов приложений ближе к концу 1990-х годов.

Предполагая, что используется подключение с сохранением состояния, можно иметь приложение, которое запрашивает данные безо всякой блокировки:

```
SCOTT@ORA12CR1> select empno, ename, sal from emp where deptno = 10;
```

| EMPNO | ENAME | SAL |
|-------|--------|------|
| 7782 | CLARK | 2450 |
| 7839 | KING | 5000 |
| 7934 | MILLER | 1300 |

В конце концов, пользователь выбирает строку, которую желает обновить. Представим, что в рассматриваемом случае он решил обновить строку MILLER. В этот момент (перед тем, как пользователь внесет любые изменения на экране, но после того, как строка извлечена из базы данных на некоторое время) приложение привяжет выбранные пользователем значения так, что мы можем запросить базу данных и удостовериться в том, что данные пока не были изменены.

В среде SQL*Plus для эмуляции операций привязки, которые должно выполнить приложение, можно выдать следующий запрос:

```
SCOTT@ORA12CR1> variable empno number
SCOTT@ORA12CR1> variable ename varchar2(20)
SCOTT@ORA12CR1> variable sal number
SCOTT@ORA12CR1> exec :empno := 7934; :ename := 'MILLER'; :sal := 1300;
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.
```

Теперь в дополнение к простому запросу значений и проверки того, что они не были изменены, мы собираемся заблокировать строку с применением FOR UPDATE NOWAIT. Приложение выполнит такой запрос:

```
SCOTT@ORA12CR1> select empno, ename, sal
2   from emp
3  where empno = :empno
4    and decode( ename, :ename, 1 ) = 1
5    and decode( sal, :sal, 1 ) = 1
6    for update nowait
7  /
```

| EMPNO | ENAME | SAL |
|-------|--------|------|
| 7934 | MILLER | 1300 |

На заметку! Почему мы использовали конструкцию `decode(столбец, :переменная_привязки, 1) = 1`? Это просто сокращенный способ выражения `where (столбец = :переменная_привязки ИЛИ (столбец является NULL и :переменная_привязки является NULL))`. Вы можете выбрать любой вариант кода, но вызов `decode()` в данном случае более компактен, и поскольку в SQL результат сравнения `NULL = NULL` никогда не является истинным (и не ложным!), один из этих двух подходов необходим, если любой из столбцов допускает значения `NULL`.

Приложение предоставляет значения переменных привязки на основе данных, отображаемых на экране (7934, MILLER и 1300) и запрашивает ту же строку из базы данных, на этот раз блокируя ее от обновлений другими сеансами; вот почему такой подход называется *пессимистической* блокировкой. Мы блокируем строку до попытки ее обновления, т.к. сомневаемся — *настроены пессимистически* — в том, что иначе она останется неизменной.

Поскольку все таблицы *должны* иметь первичный ключ (предыдущий запрос SELECT извлечет не более одной строки, потому что он содержит первичный ключ EMPNO), и первичные ключи должны быть неизменными (мы никогда не должны их обновлять), мы получим из этого оператора один из трех результатов.

- Если исходные данные не были изменены, мы получим строку MILLER, и эта строка будет заблокирована для обновления (но не чтения) другими пользователями.
- Если другой пользователь находится в процессе изменения этой строки, мы получим сообщение об ошибке ORA-00054: resource busy (ORA-00054: ресурс занят). Нам придется ожидать, пока другой пользователь не завершит работу с этой строкой.

- Если в промежутке между выбором данных и указанием о намерении обновить их кто-то другой уже изменил строку, мы получим обратно ноль строк. Это означает, что данные, отображенные на экране, устарели. Чтобы предотвратить ранее описанный сценарий потеряннного обновления, приложение должно *заново запросить* данные и заблокировать их, перед тем как позволить конечному пользователю вносить изменения. Благодаря пессимистической блокировке, при попытке обновления поля с телефонным номером пользователем User2 приложение будет знать, что поле адреса было изменено и повторно запросит данные. Таким образом, пользователь User2 не перезапишет изменение, внесенное пользователем User1, старыми данными этого поля.

После того как мы успешно заблокировали строку, приложение привяжется к новым значениям, выполнит обновление и зафиксирует изменения:

```
SCOTT@ORA12CR1> update emp
  2   set ename = :ename, sal = :sal
  3   where empno = :empno;

1 row updated.
1 строка обновлена.

SCOTT@ORA12CR1> commit;
Commit complete.
Фиксация выполнена.
```

Итак, мы изменили нужную строку весьма безопасным образом. При этом перезапись изменений какого-то другого пользователя оказывается невозможной, т.к. мы проверили, что данные не изменялись в течение периода между их первоначальным считыванием и блокировкой. Эта проверка гарантирует, что никто не изменит данные раньше нас, а блокировка — что никто не изменит их, пока мы с ними работаем.

Оптимистическая блокировка

Второй метод, который называется *оптимистической* блокировкой, откладывает все действия по блокировке до момента, непосредственно предшествующего выполнению обновления. Другими словами, мы будем изменять информацию на экране, не прибегая к блокировке. Мы *оптимистически* предполагаем, что данные не будут изменены каким-то другим пользователем, поэтому проверка того, верно ли такое предположение, откладывается до самого последнего момента.

Этот метод блокировки работает во всех средах, но увеличивает вероятность потери обновления, выполняемого пользователем. То есть, когда пользователь собирается обновить строку, он обнаруживает, что данные были модифицированы, и придется начинать все сначала.

Одна из популярных реализаций оптимистической блокировки предусматривает хранение в приложении старых и новых значений и применение для обновления данных запроса, подобного показанному ниже:

```
Update таблица
  Set столбец1 = :новый_столбец1, столбец2 = :новый_столбец2, ....
Where первичный_ключ = :первичный_ключ
  And столбец1 = :старый_столбец1
  And столбец2 = :старый_столбец2
...
```


Здесь мы оптимистически предполагаем, что данные не были изменены. В таком случае, если система сообщит об обновлении одной строки, нам повезло: данные не изменялись в промежутке между моментом их чтения и моментом, когда мы отправляем обновление. Если мы обновили ноль строк, произошла потеря: кто-то другой изменил данные, и теперь предстоит решить, каким образом продолжать работу в приложении. Должны ли мы дать конечному пользователю возможность повторения транзакции после запроса новых значений для строки (потенциально вызывая разочарование у пользователя, т.к. есть шанс, что строка снова изменится)? Или мы должны попытаться объединить значения двух обновлений, устраняя конфликт обновления на основе бизнес-правил (большой объем кода)?

Предыдущий оператор `UPDATE` фактически избегает потерянного обновления, но остается возможность его блокирования — оператору придется дожидаться завершения обновления этой строки другим сеансом. Если все ваши приложения используют оптимистическую блокировку, то применение простого оператора `UPDATE` в основном приемлемо, потому что строки блокируются на очень короткий период времени, пока обновления применяются и фиксируются. Тем не менее, если в некоторых приложениях используется пессимистическая блокировка, при которой блокировки будут удерживаться на строках в течение относительно долгих периодов времени, или если есть приложение (вроде пакетного процесса), которое может блокировать строки на длительный период (длительным периодом считается более одной или двух секунд), то вы должны подумать о применении оператора `SELECT FOR UPDATE NOWAIT` вместо выполнения проверки того, что строка не изменялась, и блокировании строки непосредственно перед выдачей операции `UPDATE`, избегая ее блокировки другим сеансом.

Существует много методов реализации оптимистического управления параллельной обработкой. Мы рассмотрели один из них, при котором все исходные образы строки будут храниться в самом приложении. В последующих разделах мы исследуем еще два метода:

- использование специального столбца, поддерживаемого триггером базы данных или кодом приложения, который сообщает “версию” записи;
- применение контрольной суммы или хеш-значения, вычисленного с использованием исходных данных.

Оптимистическая блокировка с использованием столбца версии

Это простая реализация, предусматривающая добавление по одному столбцу в каждую таблицу базы данных, которую нужно защитить от потерянных обновлений. Обычно такой столбец имеет тип `NUMBER` или `DATE/TIMESTAMP`. Как правило, он поддерживается через триггер строки самой таблицы, который отвечает за инкрементирование значения столбца `NUMBER` или обновление столбца `DATE/TIMESTAMP` при каждом изменении строки.

На заметку! Я указал, что столбец версии обычно поддерживается через триггер строки. Однако я не утверждал, что это лучший или самый правильный способ. Лично я предпочитаю, чтобы этот столбец обрабатывался самим оператором `UPDATE`, а не посредством триггера, потому что триггеров, не являющихся абсолютно необходимыми (а этот триггер обязательным назвать нельзя), следует избегать.

На заметку! Предпосылки и причины, по которым я избегаю триггеров, изложены в моей статье "Trouble with Triggers" ("Проблемы с триггерами") в журнале *Oracle Magazine*, которую можно найти в Oracle Technology Network по адресу <http://www.oracle.com/technetwork/issue-archive/2008/08-sep/o58asktom-101055.html>.

Приложение, в котором вы хотите реализовать оптимистическое управление параллельной обработкой, должно будет сохранять только значение этого дополнительного столбца, а не все предварительные образы остальных столбцов. Приложению придется проверять лишь то, что значение в этом столбце в момент запроса обновления соответствует первоначально прочитанному значению. Если эти значения совпадают, то строка не была обновлена.

Давайте взглянем на реализацию оптимистической блокировки с применением копии таблицы SCOTT.DEPT. Для создания таблицы можно воспользоваться следующим кодом DDL (Data Definition Language — язык определения данных):

```
EODA@ORA12CR1> create table dept
  2 ( deptno    number(2),
  3   dname     varchar2(14),
  4   loc       varchar2(13),
  5   last_mod   timestamp with time zone
  6             default systimestamp
  7             not null,
  8   constraint dept_pk primary key(deptno)
  9 )
10 /
```

Table created.

Таблица создана.

Теперь вставим в эту таблицу копию данных таблицы DEPT:

```
EODA@ORA12CR1> insert into dept( deptno, dname, loc )
  2 select deptno, dname, loc
  3   from scott.dept;
4 rows created.
4 строки создано.
```

```
EODA@ORA12CR1> commit;
```

Commit complete.

Фиксация выполнена.

Этот код воссоздает таблицу DEPT, но с дополнительным столбцом LAST_MOD, имеющим тип данных TIMESTAMP WITH TIME ZONE. Мы определили этот столбец как NOT NULL, так что он должен быть заполнен, и его стандартным значением является текущее время в системе.

Тип данных TIMESTAMP поддерживает очень высокую точность в Oracle, обычно вплоть до микросекунды. Для приложений, учитывающих время размышлений пользователя, такой уровень точности типа TIMESTAMP более чем достаточен, и весьма маловероятно, что процесс базы данных, извлекающий строку, и процесс пользователя, который ее просматривает, модифицирует и выдает команду обновления, могут пересечься в рамках доли секунды. И действительно, вероятность того, что два пользователя будут читать и изменять одну и ту же строку в течение одной и той же доли секунды, крайне низка.

Теперь нам требуется метод поддержания этого значения. Существуют две возможности: поддержка столбца `LAST_MOD` может выполняться либо приложением путем установки его значения в `SYSTIMESTAMP` при обновлении записи, либо триггером или хранимой процедурой. Поддержка этого столбца приложением является определенно более производительным подходом, чем прием на основе триггера, т.к. триггер добавит дополнительную обработку сверх той, что уже делается Oracle. Однако это означает, что все приложения должны поддерживать столбец `LAST_MOD` в согласованном состоянии во всех местах, где они модифицируют данную таблицу. Таким образом, если каждое приложение несет ответственность за поддержание столбца `LAST_MOD`, то оно должно постоянно выполнять проверку на предмет отсутствия в нем изменений и устанавливать его в текущее значение `SYSTIMESTAMP`. Например, если приложение запрашивает строку с `DEPTNO=10`:

```
EODA@ORA12CR1> variable deptno      number
EODA@ORA12CR1> variable dname        varchar2(14)
EODA@ORA12CR1> variable loc          varchar2(13)
EODA@ORA12CR1> variable last_mod     varchar2(50)
EODA@ORA12CR1>
EODA@ORA12CR1> begin
  2   :deptno := 10;
  3   select dname, loc, to_char( last_mod, 'DD-MON-YYYY HH.MI.SSXXFF AM TZR' )
  4     into :dname, :loc, :last_mod
  5     from dept
  6     where deptno = :deptno;
  7 end;
  8 /
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.
```

которая в текущий момент выглядит так:

```
EODA@ORA12CR1> select :deptno dno, :dname dname, :loc loc, :last_mod lm
  2   from dual;
```

| DNO | DNAME | LOC | LM |
|-----|------------|----------|---------------------------------------|
| 10 | ACCOUNTING | NEW YORK | 15-APR-2014 07.04.01.147094 PM -06:00 |

то для изменения информации будет использоваться приведенный ниже оператор обновления. Последняя строка делает очень важную проверку, чтобы удостовериться в том, что метка времени не изменилась, и применяет встроенную функцию `TO_TIMESTAMP_TZ` (`TZ` — это аббревиатура от “time zone” (часовой пояс)) для преобразования строки, сохраненной оператором `SELECT`, в подходящий тип данных. Вдобавок строка 3 оператора `UPDATE` обновляет столбец `LAST_MOD` текущим временем, если подлежащая обновлению строка найдена:

```
EODA@ORA12CR1> update dept
  2   set dname = initcap(:dname),
  3     last_mod = systimestamp
  4   where deptno = :deptno
  5     and last_mod = to_timestamp_tz(:last_mod, 'DD-MON-YYYY HH.MI.SSXXFF AM TZR');
1 row updated.
1 строка обновлена.
```

Как видите, была обновлена одна строка — строка, которая нас интересовала. Мы обновили строку по первичному ключу (DEPTNO) и проверили, что столбец LAST_MOD не был изменен каким-то другим сеансом в промежутке между моментом первоначального считывания и моментом выполнения обновления. Если бы мы попытались снова обновить эту же строку, используя ту же самую логику, но без извлечения нового значения LAST_MOD, результат был бы следующим:

```
EODA@ORA12CR1> update dept
  2   set dname = upper(:dname),
  3       last_mod = systimestamp
  4 where deptno = :deptno
  5   and last_mod = to_timestamp_tz(:last_mod, 'DD-MON-YYYY HH.MI.SSXXFF AM TZR');
0 rows updated.
0 строк обновлено.
```

Обратите внимание, что на этот раз сообщается об обновлении 0 строк, т.к. предикат на LAST_MOD не был удовлетворен. Хотя строка со значением 10 в столбце DEPTNO по-прежнему существует, значение на момент времени, когда мы пожелаем выполнить обновление, больше не совпадает со значением метки времени на момент запроса строки. Таким образом, основываясь на том факте, что строки не были модифицированы, приложение знает об изменении данных в базе данных и теперь должно принять решение о своих последующих действиях.

По ряду причин поддержку этого поля нежелательно возлагать на каждое приложение. Прежде всего, это ведет к увеличению объема кода приложения, причем дополнительный код требует повторения и корректной реализации везде, где данная таблица изменяется. В крупном приложении таких мест может быть много. Более того, каждое приложение, которое будет разработано в будущем, должно соответствовать этим же правилам. Высока вероятность упустить из виду такой фрагмент кода в каком-то приложении, в результате чего это поле будет использоваться неправильно. Следовательно, если сам код приложения не несет ответственность за поддержание поля LAST_MOD, то, по моему убеждению, приложение не должно отвечать за проверку этого поля (если оно может делать проверку, то определенно способно выполнять и обновление). Таким образом, в этом случае я предлагаю *инкапсулировать логику обновления в хранимую процедуру* и вообще не разрешать приложению обновлять таблицу напрямую. Если приложению нельзя доверить поддержание значения этого поля, то ему нельзя доверять и проверку. Итак, хранимая процедура будет принимать в качестве входных данных переменные привязки, которые мы применяли в предыдущих обновлениях, и выполнять точно такое же обновление. В случае обнаружения, что обновлено 0 строк, хранимая процедура могла бы генерировать исключение, ставя клиента в известность о том, что фактически обновление не было произведено.

Альтернативная реализация предусматривает использование триггера для поддержания поля LAST_MOD, но для таких простых задач я рекомендую избегать триггера и позволить позаботиться обо всем хранимой процедуре DML. Триггеры приносят умеренный объем накладных расходов, и в этом случае они необязательны. Вдобавок триггер не имеет возможности подтвердить, что строка не была модифицирована (он способен только предоставить значение для столбца LAST_MOD, но не проверять его во время обновления), следовательно, приложение должно быть освещено

домлено об этом столбце и о корректной работе с ним. Таким образом, одного лишь триггера оказывается недостаточно.

Оптимистическая блокировка с использованием контрольной суммы

Этот метод очень похож на предыдущий метод со столбцом версии, но в нем применяются данные самой базы для вычисления “виртуального” столбца версии. Чтобы помочь в прояснении цели и концепций, лежащих в основе контрольной суммы или хеш-функции, ниже приведена выдержка из справочного руководства по пакетам и типам PL/SQL (*Oracle Database PL/SQL Packages and Types Reference*).

Однонаправленная хеш-функция принимает входную строку переменной длины (данные) и преобразует ее в выходную строку фиксированной длины (обычно меньшую по размеру), которая называется хеш-значением. Хеш-значение служит уникальным идентификатором (подобным отпечатку пальца) входных данных. Это хеш-значение можно использовать для проверки, изменялись ли данные.

Обратите внимание, что однонаправленная хеш-функция — это хеш-функция, которую обратить нелегко. Вычислить хеш-значение на основе входных данных очень просто, но сгенерировать данные, приводящие к получению конкретного хеш-значения, довольно трудно.

С готовыми хеш-значениями или контрольными суммами можно работать в той же манере, как и со столбцом версии. Мы просто сравниваем хеш-значение или контрольную сумму, полученную при чтении данных из базы, со значением, вычисленным перед модификацией данных. Если кто-то изменил значения строки после того, как мы их прочитали, но перед тем, как их обновили, то хеш-значение или контрольная сумма почти наверняка будет отличаться.

Существует много способов вычисления хеш-значения или контрольной суммы. Ниже перечислено несколько способов, а далее в разделе приведена демонстрация одного из них. Все эти методы основаны на предоставляемой функциональности базы данных.

- `OWA_OPT_LOCK.CHECKSUM`. Этот метод доступен в Oracle8i версии 8.1.5 и выше. Имеется функция, которая возвращает 16-битную контрольную сумму для переданной ей строки, и функция, которая вычисляет 16-битную контрольную сумму для предоставленного идентификатора строки (ROWID) и одновременно блокирует эту строку. Вероятность конфликта составляет 1 к 65 536 строкам (наибольший шанс ложного совпадения значений).
- `DBMS_OBFUSCATION_TOOLKIT.MDS`. Этот метод доступен в Oracle8i версии 8.1.7 и выше. Он вычисляет 128-битный дайджест сообщения. Вероятность возникновения конфликта составляет около 1 к 3,4028E+38 (очень низкая).
- `DBMS_CRYPTO.HASH`. Этот метод доступен в Oracle 10g Release 1 и последующих версиях. Он способен вычислять дайджесты сообщений по алгоритмам SHA-1 (Secure Hash Algorithm 1 — безопасный алгоритм вычисления хеш-значения) или MD4/MD5. Рекомендуется применять алгоритм SHA-1.
- `DBMS_SQLHASH.GETHASH`. Этот метод доступен в Oracle 10g Release 2 и последующих версиях. Он поддерживает хеш-алгоритмы SHA-1, MD4 и MD5. Как привилегированный пользователь SYSDBA, вы должны выдать право на выпол-

нение этого пакета пользователю, прежде чем он сможет обратиться к нему. Пакет документирован в руководстве по безопасности базы данных Oracle (*Oracle Database Security Guide*).

- **STANDARD_HASH.** Этот метод доступен в Oracle 12c Release 1 и последующих версиях. Он представляет собой встроенную функцию SQL, которая вычисляет хеш-значение для выражения с использованием стандартных хеш-алгоритмов, таких как SHA-1 (по умолчанию), SHA-256, SHA-384, SHA-512 и MD5. Возвращаемое значение имеет тип данных RAW.
- **ORA_HASH.** Данный метод доступен в Oracle 10g Release 1 и последующих версиях. Это встроенная функция SQL, принимающая на входе значение VARCHAR2 и (необязательно) еще два параметра, которые управляют возвращаемым значением. Возвращаемым значением является число, по умолчанию находящееся в пределах от 0 до 4 294 967 295.

На заметку! Массив функций вычисления хеш-значений и контрольных сумм доступен во многих языках программирования, поэтому в вашем распоряжении могут быть и другие функции, реализованные за рамками базы данных. Тем не менее, применяя встроенные возможности базы данных, вы увеличиваете степень переносимости (на новые языки и новые подходы) в будущем.

В следующем примере показано, как использовать встроенную функцию DBMS_CRYPTO в Oracle 10g и последующих версиях для вычисления хеш-значений/контрольных сумм. Такой прием также применим к остальным перечисленным подходам; логика будет отличаться не сильно, но вызываемые API-интерфейсы оказываются очень разными. Мы начнем с того, что избавимся от столбца, который использовался в предыдущем примере:

```
EODA@ORA12CR1> alter table dept drop column last_mod;
Table altered.
Таблица изменена.
```

Затем приложение должно запросить и отобразить информацию об отделе с номером 10. Обратите внимание, что при запрашивании информации мы вычисляем хеш-значение с помощью встроенной функции DBMS_CRYPTO. Это информация о версии, которая сохраняется в приложении. Следующий код выдает запрос и отображает информацию:

```
EODA@ORA12CR1> variable deptno number
EODA@ORA12CR1> variable dname varchar2(14)
EODA@ORA12CR1> variable loc varchar2(13)
EODA@ORA12CR1> variable hash number
EODA@ORA12CR1> begin
2  select deptno, dname, loc,
3      ora_hash( dname || '/' || loc ) hash
4  into :deptno, :dname, :loc, :hash
5  from dept
6  where deptno = 10;
7  end;
8  /
```

PL/SQL procedure successfully completed.

Процедура PL/SQL успешно завершена.

```
EODA@ORA12CR1> select :deptno, :dname, :loc, :hash
2      from dual;
```

| :DEPTNO | :DNAME | :LOC | :HASH |
|---------|------------|----------|------------|
| 10 | Accounting | NEW YORK | 2721972020 |

Как видите, хеш-значение — это просто некоторое число. Именно его мы хотим применять перед обновлением. Чтобы обновить строку, понадобится заблокировать строку в базе данных в текущем состоянии и затем сравнить хеш-значение этой строки с хеш-значением, которое было вычислено при чтении данных из базы. Логика приложения может выглядеть примерно так:

```
EODA@ORA12CR1> exec :dname := lower(:dname);
```

PL/SQL procedure successfully completed.

Процедура PL/SQL успешно завершена.

```
EODA@ORA12CR1> update dept
2      set dname = :dname
3  where deptno = :deptno
4      and ora_hash( dname || '/' || loc ) = :hash
5  /
```

1 row updated.

1 строка обновлена.

```
EODA@ORA12CR1> select dept.*,
2      ora_hash( dname || '/' || loc ) hash
3  from dept
4  where deptno = :deptno;
```

| DEPTNO | DNAME | LOC | HASH |
|--------|------------|----------|------------|
| 10 | accounting | NEW YORK | 2818855829 |

При повторном запросе этих же данных и вычислении хеш-значения заново после обновления можно заметить, что хеш-значение отличается. Если кто-то модифицировал строку до нас, то хеш-значения не совпадали бы. Мы можем увидеть это, попытавшись снова выполнить наше обновление с использованием *старого* хеш-значения, которое было вычислено при чтении данных в первый раз:

```
EODA@ORA12CR1> update dept
2      set dname = :dname
3  where deptno = :deptno
4      and ora_hash( dname || '/' || loc ) = :hash
5  /
```

0 rows updated.

0 строк обновлено.

Как видите, было обновлено 0 строк, потому что хеш-значение не соответствует данным, находящимся в текущий момент в базе.

Для корректной работы этого подхода на основе хеш-значений мы должны обеспечить, чтобы в каждом приложении применялся один и тот же способ вычисления хеша, а именно — конкатенации dname с символом / и loc, причем в указанном

порядке. Чтобы сделать этот подход универсальным, я бы посоветовал добавить к таблице виртуальный столбец (в Oracle 11g Release 1 и последующих версиях) либо использовать для добавления столбца представление, скрыв эту функцию от самого приложения. Вот как выглядит добавление столбца в Oracle Database 11g Release 1 и последующих версиях:

```
EODA@ORA12CR1> alter table dept
 2 add hash as
 3 ( ora_hash(dname || '/' || loc ) );
Table altered.
Таблица изменена.

EODA@ORA12CR1> select *
 2 from dept
 3 where deptno = :deptno;
```

| DEPTNO | DNAME | LOC | HASH |
|--------|------------|----------|------------|
| 10 | accounting | NEW YORK | 2818855829 |

Добавленный столбец является виртуальным, и как таковой не влечет за собой накладные расходы по хранению. Его значение не вычисляется заранее для сохранения на диске. Вместо этого оно вычисляется при извлечении данных из базы.

В приведенном примере продемонстрирована реализация оптимистической блокировки с применением хеш-значения или контрольной суммы. Следует иметь в виду, что вычисление хеш-значения или контрольной суммы — операция, требующая значительной загрузки процессора. Это следует учитывать в тех системах, где процессорное время является дефицитным ресурсом. Однако такой подход значительно более дружелюбен по отношению к сети, потому что передача через сеть сравнительно небольшого хеш-значения вместо начального и конечного образа строки (для выполнения сравнения по столбцам) требует значительно меньшего объема сетевых ресурсов.

Выбор между оптимистической и пессимистической блокировкой

Так какой же метод лучше? Согласно моему опыту, пессимистическая блокировка очень хорошо работает в Oracle (но, возможно, не так хорошо в других СУБД) и обладает многими преимуществами по сравнению с оптимистической блокировкой. Тем не менее, она требует подключения к базе данных с сохранением состояния, подобного клиент-серверному подключению. Причина в том, что блокировки не удерживаются между подключениями. Этот простой факт делает пессимистическую блокировку нереальной в наши дни. В прошлом, когда приходилось иметь дело с клиент-серверными приложениями и парой десятков или сотен пользователей, пессимистическая блокировка была бы моим первым и последним вариантом. Однако в настоящее время для большинства приложений я рекомендую использовать оптимистическое управление параллельной обработкой. Наличие подключения на протяжении всей транзакции — просто чересчур высокая цена.

А что можно сказать о доступных методах? Лично я предпочитаю применять подход со столбцом версии в форме метки времени. Это предоставляет дополнительную информацию об обновлении в долгосрочном смысле. К тому же этот метод требует меньшего объема вычислений, чем использование хеш-значения или контрольной

суммы, и не создает потенциальных проблем, которые могут возникать при вычислении хеш-значения или контрольной суммы для данных типа LONG, LONG ROW, CLOB, BLOB и других очень больших столбцов (типы LONG и LONG RAW устарели; я упоминаю их здесь потому, что они по-прежнему часто применяются в словаре данных Oracle).

Если бы мне понадобилось добавить средства оптимистического управления параллельной обработкой в таблицу, которая должна применяться в схеме пессимистической блокировки (например, таблицу, к которой происходит обращение из клиент-серверных приложений и веб-приложений), то я бы остановил свой выбор на подходе с ORA_HASH. Это связано с тем, что существующие унаследованные приложения могут не воспринять положительно появление нового столбца. Даже если предпринять меры для сокрытия этого дополнительного столбца, приложение может пострадать от накладных расходов, вызванных обязательным в таком случае триггером. В этом отношении прием с ORA_HASH является ненавязчивым и легковесным. Подход с хешированием/контрольной суммой может быть в значительной степени независимым от базы данных, особенно если хеш-значения или контрольные суммы вычисляются вне базы данных. Тем не менее, выполнение вычислений в промежуточном звене, а не в базе данных, ведет к значительному увеличению объема используемых ресурсов процессора и пропускной способности сети.

Блокирование

Блокирование (blocking) происходит, когда один сеанс удерживает блокировку (lock) на ресурсе, который запрашивается другим сеансом. В результате запрашивающий сеанс окажется заблокированным — он зависнет до тех пор, пока удерживающий блокировку сеанс не освободит заблокированный ресурс. Блокирования почти всегда можно избежать. Фактически, если оказывается, что сеанс заблокирован в интерактивном приложении, то вполне вероятно, что приложение содержит также ошибку потерянного обновления, возможно, даже не подозревая об этом. То есть, логика приложения содержит изъян и в этом заключается причина блокирования.

Пятью распространенными операторами DML, которые будут вызывать блокирование в базе данных, являются INSERT, UPDATE, DELETE, MERGE и SELECT FOR UPDATE. Решение для заблокированного оператора SELECT FOR UPDATE тривиально: нужно просто добавить к нему конструкцию NOWAIT и больше он блокироваться не будет. Взамен приложение будет возвращать конечному пользователю сообщение о том, что строка уже заблокирована. Остальные четыре оператора DML более интересны. Мы рассмотрим каждый из них и выясним, когда они не должны приводить к блокированию, и как исправить ситуацию, если это все же происходит.

Заблокированные вставки

Существует несколько ситуаций, когда оператор INSERT будет вызывать блокирование. Наиболее распространенный сценарий предусматривает наличие таблицы с первичным ключом или помещенным на нее уникальным ограничением и двух сеансов, пытающихся вставить строку с одним и тем же значением. Один сеанс будет заблокирован до тех пор, пока второй не выполнит либо фиксацию (в этом случае заблокированный сеанс получит сообщение о наличии дублированного значения), либо откат (в таком случае заблокированный сеанс успешно завершит свою операцию). Еще одна ситуация касается таблиц, связанных друг с другом ограничениями

ссылочной целостности. Вставка строки в дочернюю таблицу может оказаться заблокированной, если строка, от которой она зависит, создается или удаляется.

Обычно заблокированные вставки появляются в приложениях, которые позволяют конечному пользователю генерировать значение первичного ключа/уникального столбца. Предотвратить возникновение такой ситуации легче всего, применяя для генерации значений первичного ключа/уникального столбца последовательность или встроенную функцию `SYS_GUID()`. Последовательности и `SYS_GUID()` были разработаны в качестве методов генерации ключей с высокой степенью параллелизма в многопользовательской среде. В случае если вы не можете использовать один из этих методов, и должны разрешить конечному пользователю генерировать ключи, которые могут дублироваться, можете применить описанный ниже прием, который устраняет проблему посредством устанавливаемых вручную блокировок, реализованных с помощью встроенного пакета `DBMS_LOCK`.

На заметку! В следующем примере показано, как предотвратить блокирование сеансом оператора вставки из-за наличия ограничения первичного ключа или ограничения уникальности. Необходимо подчеркнуть, что демонстрируемый здесь способ устранения проблемы должен рассматриваться как кратковременное решение на период, пока инспектируется сама архитектура приложения. Этот подход добавляет очевидные накладные расходы и требует серьезного отношения к своей реализации. В правильно спроектированном приложении проблема подобного рода возникать не будет (например, в параллельной среде не должны выполняться транзакции, продолжающиеся часами). Такой прием следует считать последним средством и определенно не чем-то таким, что нужно делать в отношении каждой таблицы приложения просто “на всякий случай”.

При выполнении операций вставки не существует строк, которые можно было бы выбрать и заблокировать; нет никакого способа предотвратить вставку строки с таким же значением другими пользователями, что заблокирует наш сеанс и вызовет ожидание на неопределенное время. Именно здесь в игру вступает пакет `DBMS_LOCK`. Для демонстрации этого подхода мы создадим таблицу с первичным ключом и триггером, который будет препятствовать одновременной вставке одних и тех же значений двумя (и более) сеансами. Триггер будет использовать функцию `DBMS_UTILITY.GET_HASH_VALUE` для хеширования первичного ключа в некоторое число между 0 и 1 073 741 823 (разрешенный для применения диапазон значений идентификатора блокировки в Oracle). В этом примере я выбрал хеш-таблицу с размером 1024, т.е. мы будем хешировать первичные ключи в один из 1024 разных идентификаторов блокировок. Затем с помощью функции `DBMS_LOCK.REQUEST` мы установим монопольную блокировку на основе этого идентификатора. В каждый момент времени это может делать только один сеанс, так что если другой пользователь попытается вставить в нашу таблицу строку с тем же первичным ключом, то выданный им запрос на блокировку потерпит неудачу (с сообщением о занятости ресурса).

На заметку! Для успешной компиляции этого триггера привилегия на выполнение `DBMS_LOCK` должна быть выдана непосредственно схеме. Привилегия на выполнение пакета `DBMS_LOCK` может не поступить из назначенной роли.

```
SCOTT@ORA12CR1> create table demo ( x int primary key );
```

Table created.

Таблица создана.

```
SCOTT@ORA12CR1> create or replace trigger demo_bifer
```

```
2 before insert on demo
```

```
3 for each row
```

```
4 declare
```

```
5     l_lock_id number;
```

```
6     resource_busy exception;
```

```
7     pragma exception_init( resource_busy, -54 );
```

```
8 begin
```

```
9     l_lock_id :=
```

```
10         dbms_utility.get_hash_value( to_char( :new.x ), 0, 1024 );
```

```
11     if ( dbms_lock.request
```

```
12         ( id => l_lock_id,
```

```
13         lockmode => dbms_lock.x_mode,
```

```
14         timeout => 0,
```

```
15         release_on_commit => TRUE ) not in (0,4) )
```

```
16     then
```

```
17         raise resource_busy;
```

```
18     end if;
```

```
19 end;
```

```
20 /
```

Trigger created.

Триггер создан.

```
SCOTT@ORA12CR1> insert into demo(x) values (1);
```

1 row created.

Теперь для иллюстрации перехвата этой проблемы блокирующей операции INSERT в одиночном сеансе мы используем AUTONOMOUS_TRANSACTION, чтобы это выглядело, как если бы следующий блок кода выполнялся в другом сеансе SQL*Plus. На самом деле, если применить другой сеанс, поведение будет таким же. Вот что мы сделаем:

```
SCOTT@ORA12CR1> declare
```

```
2     pragma autonomous_transaction;
```

```
3 begin
```

```
4     insert into demo(x) values (1);
```

```
5     commit;
```

```
6 end;
```

```
7 /
```

```
declare
```

```
*
```

ERROR at line 1:

ORA-00054: resource busy and acquire with NOWAIT specified or timeout expired

ORA-06512: at "SCOTT.DEMO_BIFER", line 14

ORA-04088: error during execution of trigger 'SCOTT.DEMO_BIFER'

ORA-06512: at line 4

ОШИБКА в строке 1:

ORA-00054: ресурс занят и запрошен с указанием NOWAIT либо истек тайм-аут

ORA-06512: в SCOTT.DEMO_BIFER, строка 14

ORA-04088: ошибка во время выполнения триггера SCOTT.DEMO_BIFER

ORA-06512: в строке 4

Концепция заключается в извлечении переданного значения первичного ключа из таблицы, защищенной триггером, и помещении его в символьную строку. Затем с помощью функции `DBMS_UTILITY.GET_HASH_VALUE` для строки можно получить по большей части уникальное хеш-значение. До тех пор, пока размер хеш-таблицы меньше 1 073 741 823, это значение можно блокировать в монопольном режиме с использованием пакета `DBMS_LOCK`.

После хеширования мы берем это значение и применяем пакет `DBMS_LOCK` для запроса монопольной блокировки этого идентификатора блокировки с тайм-аутом `ZERO` (что обеспечивает немедленный возврат, если кто-то другой заблокировал данное значение). При истечении тайм-аута или неудаче по любой причине возвращается сообщение об ошибке `ORA-00054: resource busy (ORA-00054: ресурс занят)`. В противном случае ничего не делается: вставка прошла успешно и блокировка не требуется. После выполнения фиксации нашей транзакции все блокировки, включая размещенные этим вызовом `DBMS_LOCK`, будут освобождены.

Разумеется, если первичный ключ таблицы имеет тип `INTEGER`, а его значение не превышает 1 миллиард, то можно не прибегать к хешированию и просто использовать число в качестве идентификатора блокировки.

Во избежание искусственной генерации сообщений о занятости ресурса из-за хеширования различных строк в одно и то же число придется поэкспериментировать с размером хеш-таблицы (1024 в этом примере). Размер хеш-таблицы будет зависеть от приложения (данных), и на него оказывает влияние также количество одновременно выполняемых вставок. Можно также добавить к триггеру флаг, чтобы позволить пользователям включать и отключать проверку. Например, если мы собираемся вставлять сотни или тысячи записей, то проверку имеет смысл отключить.

Блокированные слияния, обновления и удаления

В интерактивном приложении — приложении, которое запрашивает какие-либо данные из базы, предоставляет конечному пользователю возможность манипулирования ими и затем помещает их обратно в базу — заблокированный запрос `UPDATE` или `DELETE` указывает на то, что в коде, вероятно, присутствует проблема потерянного обновления. (Если хотите, я буду называть это ошибкой в коде.) Вы пытаетесь выполнить операцию `UPDATE` для строки, которую уже кто-то обновляет (другими словами, для строки, которую кто-то другой заблокировал). Проблемы блокирования можно избежать, применяя запрос `SELECT FOR UPDATE NOWAIT`, чтобы:

- удостовериться в том, что данные не были изменены с момента их запроса (предотвращая потерянные обновления);
- заблокировать строку (предотвращая блокирование запроса `UPDATE` или `DELETE`).

Как обсуждалось ранее, это можно делать независимо от принятого метода блокировки. И пессимистическая, и оптимистическая блокировки могут использовать запрос `SELECT FOR UPDATE NOWAIT` для проверки отсутствия изменений в строке. Пессимистическая блокировка будет применять оператор `SELECT FOR UPDATE NOWAIT` немедленно после того, как пользователь сообщит о своем намерении модифицировать данные. Оптимистическая блокировка будет использовать этот оператор непосредственно перед обновлением данных в базе. Это не только решит проблему блокирования в приложении, но и устранил проблему целостности данных.

Поскольку запрос MERGE является просто сочетанием INSERT и UPDATE (в Oracle 10g и последующих версиях поддерживается расширенный синтаксис оператора MERGE, включающий также DELETE), оба приема будут применяться одновременно.

Взаимоблокировки

Взаимоблокировка (deadlock) возникает при наличии двух сеансов, каждый из которых блокирует ресурс, требуемый другому сеансу. Например, взаимоблокировку легко продемонстрировать, имея две таблицы А и В, которые содержат по одной строке. Для этого достаточно открыть два сеанса (скажем, два сеанса SQL*Plus). В сеансе А мы будем обновлять таблицу А, а в сеансе В — таблицу В. Теперь при попытке обновления таблицы А в сеансе В этот сеанс окажется заблокированным. Сеанс А уже заблокировал данную строку. Это еще не взаимоблокировка, а простая блокировка. Ситуация взаимоблокировки пока не возникла, т.к. есть шанс, что сеанс А выполнит фиксацию или откат и сеанс В просто продолжит работу с этой точки.

Если теперь вернуться в сеанс А и попробовать обновить таблицу В, возникнет взаимоблокировка. Один из двух сеансов будет выбран в качестве жертвы и вынужден будет произвести откат своего оператора. Например, попытка сеанса В обновить таблицу А может быть отменена с выводом сообщения об ошибке, подобного следующему:

```
update a set x = x+1
      *
```

ERROR at line 1:

ORA-00060: deadlock detected while waiting for resource

ОШИБКА в строке 1:

ORA-00060: обнаружена взаимоблокировка во время ожидания ресурса

Попытка обновления таблицы В в сеансе А остается заблокированной — Oracle не будет производить откат всей транзакции. Откатывается только один из операторов, участвующих во взаимоблокировке. Сеанс В по-прежнему блокирует строку в таблице В, а сеанс А терпеливо дожидается освобождения этой строки. После получения сообщения о взаимоблокировке сеанс В должен принять решение о том, что делать: зафиксировать ожидающую работу в таблице В, произвести ее откат или продолжить двигаться альтернативным путем и выполнить фиксацию позже. Как только этот сеанс сделает фиксацию или откат, другой заблокированный сеанс продолжит работу, будто бы ничего не происходило.

СУБД Oracle считает взаимоблокировки настолько редкими и необычными, что создает трассировочный файл на сервере каждый раз, когда они возникают. Содержимое этого трассировочного файла может выглядеть следующим образом:

```
*** 2014-04-16 18:58:26.602
*** SESSION ID: (31.18321) 2014-04-16 18:58:26.603
*** CLIENT ID: () 2014-04-16 18:58:26.603
*** SERVICE NAME: (SYS$USERS) 2014-04-16 18:58:26.603
*** MODULE NAME: (SQL*Plus) 2014-04-16 18:58:26.603
*** ACTION NAME: () 2014-04-16 18:58:26.603
*** 2014-04-16 18:58:26.603
DEADLOCK DETECTED ( ORA-00060 )
[Transaction Deadlock]
```

The following deadlock is not an ORACLE error. It is a deadlock due to user error in the design of an application or from issuing incorrect ad-hoc SQL. The following information may aid in determining the deadlock:

ОБНАРУЖЕНА ВЗАИМОБЛОКИРОВКА (ORA-00060)

[Взаимоблокировка транзакций]

Следующая взаимоблокировка не является ошибкой ORACLE. Она возникла из-за пользовательской ошибки в проектном решении приложения или по причине выдачи некорректного оператора SQL. Приведенная ниже информация может помочь в определении условий взаимоблокировки:

Очевидно, что Oracle трактует эти взаимоблокировки, возникающие в приложении, как ошибку в какой-то части приложения, и в большинстве случаев это действительно так. В отличие от многих других СУБД, взаимоблокировки в Oracle случаются настолько редко, что практически их можно считать несуществующими. Обычно для их возникновения требуются особые, совершенно искусственные условия.

На основе собственного опыта я полагаю, что главной причиной возникновения взаимоблокировок в базе данных Oracle является наличие неиндексированных внешних ключей. (Вторая по значимости причина — использование битовых индексов в таблицах, подвергающихся одновременным обновлениям — рассматривается в главе 11.) СУБД Oracle поместит блокировку на всю дочернюю таблицу после модификации родительской таблицы в следующих трех сценариях.

- При обновлении первичного ключа родительской таблицы (очень редкая ситуация, если соблюдается правило проектирования реляционных баз данных, в соответствии с которым первичные ключи должны быть неизменяемыми) дочерняя таблица будет заблокирована в случае отсутствия индекса на внешнем ключе.
- При удалении строки родительской таблицы вся дочерняя таблица также будет заблокирована (при отсутствии индекса на внешнем ключе).
- При слиянии с родительской таблицей вся дочерняя таблица также будет заблокирована (при отсутствии индекса на внешнем ключе). Обратите внимание, что это справедливо только в версиях Oracle9i и Oracle 10g, но больше не касается Oracle 11g Release 1 и последующих версий.

В Oracle9i и последующих версиях эти полные блокировки таблиц действуют кратковременно — т.е. их следует принимать во внимание на протяжении выполнения операции DML, а не всей транзакции. Но даже при этом они могут приводить к крупным проблемам блокировки. В качестве иллюстрации первой из названных ситуаций рассмотрим пример. При наличии двух таблиц, созданных следующим образом, вначале не происходит ничего экстраординарного:

```
EODA@ORA12CR1> create table p ( x int primary key );
Table created.
Таблица создана.
```

```
EODA@ORA12CR1> create table c ( x references p );
Table created.
Таблица создана.
```

```

EODA@ORA12CR1> insert into p values ( 1 );
1 row created.
1 строка создана.

EODA@ORA12CR1> insert into p values ( 2 );
1 row created.
1 строка создана.

EODA@ORA12CR1> commit;
Commit complete.
Фиксация выполнена.

EODA@ORA12CR1> insert into c values ( 2 );
1 row created.
1 строка создана.

```

Но при переходе в другой сеанс и попытке удалить первую родительскую запись сеанс немедленно оказывается заблокированным:

```

EODA@ORA12CR1> delete from p where x = 1;

```

Это попытка блокирования всей таблицы C перед тем, как она выполнит удаление. Теперь никакой другой сеанс не может инициировать выполнение операций DELETE, INSERT или UPDATE по отношению к любой строке таблицы C (уже запущенные сеансы могут продолжаться, однако запустить новый сеанс для изменения таблицы C не удастся).

Такая блокировка будет происходить также при обновлении значения первичного ключа. Поскольку в реляционных базах данных обновление первичного ключа — совершенно недопустимая операция, в общем случае эта проблема не должна возникать при обновлениях. Тем не менее, мне приходилось видеть, что обновление первичного ключа становится серьезной проблемой, когда разработчики применяют инструменты автоматической генерации SQL-кода, и эти инструменты обновляют каждый столбец независимо от того, действительно ли он был изменен конечным пользователем. Предположим, например, что мы пользуемся Oracle Forms и создаем стандартный макет для редактирования любой таблицы. По умолчанию Oracle Forms будет генерировать операцию обновления, которая модифицирует каждый столбец таблицы, выбранной для отображения. Если мы строим стандартный макет для таблицы DEPT и включаем в него все три поля, то Oracle Forms будет выполнять следующую команду всякий раз, когда мы изменяем любой столбец таблицы DEPT:

```

update dept set deptno=:1,dname=:2,loc=:3 where rowid=:4

```

В этом случае, если таблица EMP имеет внешний ключ, ссылающийся на DEPT, а индекс на столбце DEPTNO таблицы EMP отсутствует, то вся таблица EMP будет заблокирована во время обновления DEPT. За ситуациями подобного рода необходимо внимательно следить, если вы имеете дело с любыми инструментами, которые генерируют SQL-код автоматически. Несмотря на то что значение первичного ключа не изменяется, дочерняя таблица EMP будет заблокирована после выполнения приведенного ранее SQL-оператора. В случае применения Oracle Forms решение заключается в установке свойства UPDATE CHANGED COLUMNS ONLY (обновлять только измененные столбцы) этой таблицы в YES (да). Тогда Oracle Forms будет генерировать оператор UPDATE, который включает только измененные столбцы (но не первичный ключ).

Проблемы, вытекающие из удаления строки в родительской таблице, встречаются намного чаще. Как уже демонстрировалось, при удалении строки из таблицы Р дочерняя таблица С блокируется на период действия операции DML, предотвращая другие обновления С в течение выполнения транзакции (конечно, предполагается, что никто другой не изменял таблицу С; в такой ситуации операция удаления будет ожидать). Именно здесь возникают проблемы блокирования и взаимоблокировки. Блокировка всей таблицы С значительно снижает показатель параллелизма базы данных — до точки, когда никто ничего не сможет изменить в таблице С. Вдобавок увеличивается вероятность возникновения взаимоблокировки, т.к. один пользователь владеет огромным объемом данных до тех пор, пока не выполнит фиксацию. Теперь вероятность блокирования на таблице С какого-то другого сеанса значительно выше; любой сеанс, пытающийся модифицировать С, будет заблокирован. Следовательно, начинает появляться множество сеансов, удерживающих ряд существующих блокировок на других ресурсах, которые оказываются заблокированными в базе данных. Если любой из этих заблокированных сеансов блокирует ресурс, также требуемый текущему сеансу, возникает взаимоблокировка. В этом случае взаимоблокировка вызвана текущим сеансом, препятствуя доступу к значительно большему числу ресурсов (в рассматриваемой ситуации ко всем строкам одной таблицы), чем требуется в действительности. Когда кто-либо жалуется на взаимоблокировки в базе данных, я советую запустить сценарий, отыскивающий неиндексированные внешние ключи; в 99% случаев удастся найти таблицу, которая нарушает работу. Простая индексация внешнего ключа устраняет взаимоблокировки, равно как и массу других проблем конкуренции. В приведенном ниже примере показано, как использовать этот сценарий для нахождения неиндексированного внешнего ключа в таблице С:

```
EODA@ORA12CR1> column columns format a30 word_wrapped
EODA@ORA12CR1> column table_name format a15 word_wrapped
EODA@ORA12CR1> column constraint_name format a15 word_wrapped
EODA@ORA12CR1> select table_name, constraint_name,
2      cname1 || nvl2(cname2, ',' || cname2, null) ||
3      nvl2(cname3, ',' || cname3, null) || nvl2(cname4, ',' || cname4, null) ||
4      nvl2(cname5, ',' || cname5, null) || nvl2(cname6, ',' || cname6, null) ||
5      nvl2(cname7, ',' || cname7, null) || nvl2(cname8, ',' || cname8, null)
6      columns
7  from ( select b.table_name,
8          b.constraint_name,
9          max(decode( position, 1, column_name, null )) cname1,
10         max(decode( position, 2, column_name, null )) cname2,
11         max(decode( position, 3, column_name, null )) cname3,
12         max(decode( position, 4, column_name, null )) cname4,
13         max(decode( position, 5, column_name, null )) cname5,
14         max(decode( position, 6, column_name, null )) cname6,
15         max(decode( position, 7, column_name, null )) cname7,
16         max(decode( position, 8, column_name, null )) cname8,
17         count(*) col_cnt
18  from (select substr(table_name,1,30) table_name,
19         substr(constraint_name,1,30) constraint_name,
20         substr(column_name,1,30) column_name,
21         position
```



```

22         from user_cons_columns ) a,
23         user_constraints b
24     where a.constraint_name = b.constraint_name
25         and b.constraint_type = 'R'
26     group by b.table_name, b.constraint_name
27 ) cons
28 where col_cnt > ALL
29     ( select count(*)
30       from user_ind_columns i,
31           user_indexes      ui
32     where i.table_name = cons.table_name
33           and i.column_name in (cname1, cname2, cname3, cname4,
34                                cname5, cname6, cname7, cname8 )
35           and i.column_position <= cons.col_cnt
36           and ui.table_name = i.table_name
37           and ui.index_name = i.index_name
38           and ui.index_type IN ( 'NORMAL', 'NORMAL/REV' )
39     group by i.index_name
40 )
41 /

```

| TABLE_NAME | CONSTRAINT_NAME | COLUMNS |
|------------|-----------------|---------|
| ----- | ----- | ----- |
| C | SYS_C0061427 | X |

Сценарий работает с ограничениями внешнего ключа, которые содержат до восьми столбцов (если количество столбцов больше, возможно, проектное решение требует пересмотра). Он начинается с построения встроенного представления CONS в предыдущем запросе. Это встроенное представление преобразует соответствующие имена столбцов ограничения из строк в столбцы, давая в итоге для каждого ограничения строку, которая содержит до восьми столбцов, имеющих имена как у столбцов в ограничении. Кроме того, предусмотрен столбец COL_CNT, содержащий количество столбцов в самом ограничении внешнего ключа. Для каждой строки, возвращенной из встроенного представления, выполняется коррелированный подзапрос, который проверяет все индексы обрабатываемой в данный момент таблицы. Он подсчитывает количество столбцов в индексе, которые соответствуют столбцам в ограничении внешнего ключа, и затем группирует их по имени индекса.

Таким образом, подзапрос генерирует набор чисел, каждое из которых представляет количество совпадающих столбцов в каком-либо индексе обрабатываемой таблицы. Если исходное значение COL_CNT больше всех этих чисел, то таблица не содержит ни одного индекса, поддерживающего это ограничение. Если значение COL_CNT меньше всех этих чисел, то таблица содержит, по крайней мере, один индекс, который поддерживает данное ограничение. Обратите внимание на применение функции NVL2, предназначенной для преобразования списка имен столбцов в список с разделителями-запятыми. Функция NVL2 принимает три аргумента: А, В и С. Если аргумент А не равен null, функция возвращает аргумент В; в противном случае она возвращает аргумент С. В этом запросе предполагается, что владелец ограничения является также владельцем таблицы и индекса. Если таблицу индексировал другой пользователь или если она принадлежит другой схеме (обе ситуации встречаются крайне редко), то сценарий будет работать некорректно.

В предыдущем сценарии также производится проверка, имеет ли индекс тип В-дерева (NORMAL или NORMAL/REV). Причина в том, что битовый индекс на столбце внешнего ключа не предотвращает возникновение проблемы блокировки.

На заметку! В средах хранилищ данных общепринято создавать битовые индексы на столбцах внешних ключей таблицы фактов. Однако в таких средах загрузка данных обычно делается в упорядоченной манере с помощью запланированных процессов ETL (extract, transform, load — извлечение, преобразование, загрузка) и, следовательно, не будет сталкиваться с ситуацией вставки в дочернюю таблицу в одном процессе и параллельного удаления из родительской таблицы в другом процессе (то, что можно встретить в приложении OLTP).

Таким образом, предыдущий сценарий показывает, что таблица С имеет внешний ключ на столбце X, но без индекса. Создавая на столбце X индекс со структурой В-дерева, мы можем полностью устранить эту проблему блокирования. Кроме блокировки этой таблицы неиндексированный внешний ключ может также порождать проблемы в перечисленных ниже случаях.

- При наличии конструкции ON DELETE CASCADE и неиндексированной дочерней таблицы. Например, предположим, что EMP — это дочерняя таблица таблицы DEPT. Оператор DELETE DEPTNO=10 должен выполнить операцию ON DELETE CASCADE в таблице EMP. Если столбец DEPTNO в EMP не индексирован, то каждое удаление строки из таблицы DEPT будет приводить к полному просмотру таблицы EMP. Скорее всего, этот полный просмотр нежелателен, и в случае удаления множества строк из родительской таблицы дочерняя таблица будет просматриваться при каждом удалении строки из родительской таблицы.
- При выполнении запроса дочерней таблицы из родительской. Снова обратимся к примеру с таблицами EMP/DEPT. Запрос таблицы EMP в контексте столбца DEPTNO является довольно распространенной ситуацией. При частом запуске следующего запроса (скажем, при создании отчета), вы обнаружите, что отсутствие индекса будет замедлять выполнение запросов:

```
select * from dept, emp
where emp.deptno = dept.deptno and dept.deptno = :X;
```

А когда не следует индексировать внешний ключ? В общем случае индексация не требуется при удовлетворении перечисленных ниже условий.

- Отсутствие удалений из родительской таблицы.
- Отсутствие обновлений значения уникального/первичного ключа родительской таблицы (при этом следует следить за неумышленными обновлениями первичного ключа со стороны инструментов).
- Отсутствие соединений родительской таблицы с дочерней таблицей (подобных соединению DEPT и EMP).

Если удовлетворены все три условия, можете смело отказываться от индекса — он не нужен. Если соблюдается любое из перечисленных условий, помните о возможных последствиях. Это один из тех редких случаев, когда СУБД Oracle склонна чрезмерно блокировать данные.

Эскалация блокировок

При эскалации блокировок система снижает уровень их детализации. Примером может служить преобразование ста блокировок на уровне строк в единственную блокировку на уровне таблицы. В этом случае мы используем одну блокировку для блокирования всего и, как правило, блокируем значительно больше данных, чем ранее. Эскалация блокировок часто применяется в базах данных, где блокировки считаются дефицитным ресурсом и причиной накладных расходов, которых стремятся избежать.

На заметку! СУБД Oracle никогда не выполняет эскалацию блокировок. Ни при каких условиях.

СУБД Oracle никогда не предпринимает эскалацию блокировок, но практикует преобразование или продвижение блокировок — термины, которые часто путают с эскалацией блокировок.

На заметку! Понятия *преобразование блокировки* и *продвижение блокировки* являются синонимами. В Oracle этот процесс обычно называется преобразованием блокировки (lock conversion).

СУБД Oracle будет устанавливать блокировку на минимально возможном уровне (т.е. использовать наименее ограничивающую блокировку) и при необходимости преобразовывать ее в блокировку более ограничивающего уровня. Например, при выборе строки в таблице с указанием конструкции FOR UPDATE создаются две блокировки. Одна помещается на выбранную строку или строки (и эта блокировка будет монопольной; никто другой не сможет заблокировать эту конкретную строку в монопольном режиме). Вторая блокировка, ROW SHARE TABLE, помещается на саму таблицу. Это предотвратит блокирование таблицы в монопольном режиме другими пользователями, что будет препятствовать, например, изменению ими структуры таблицы. Другой сеанс может модифицировать любую другую строку в этой таблице безо всяких конфликтов. При наличии заблокированной строки в таблице будет разрешено выполнять максимально возможное количество команд.

Эскалация блокировки — это не “функциональная возможность” базы данных. Она не является желательным атрибутом. Тот факт, что база данных поддерживает эскалацию блокировок, подразумевает наличие некоторых накладных расходов, присущих механизму блокирования, и выполнение значительной работы по управлению сотнями блокировок. Накладные расходы, связанные с 1 блокировкой и 1 миллионом блокировок, в Oracle одинаковы — они нулевые.

Типы блокировок

Ниже описаны три основных класса блокировок в Oracle.

- **Блокировки DML.** Здесь под DML понимается *Data Manipulating Language* (язык манипулирования данными). В общем случае это означает операторы SELECT, INSERT, UPDATE, MERGE и DELETE. Блокировки DML представляют собой механизм, который обеспечивает параллельное изменение данных. Блокировками

DML будут, к примеру, блокировки определенных строк данных или блокировка на уровне таблицы, которая блокирует каждую строку в таблице.

- **Блокировки DDL.** Здесь DDL означает *Data Definition Language* (язык определения данных) и охватывает операторы CREATE, ALTER и т.д. Блокировки DDL защищают определение структуры объектов.
- **Внутренние блокировки и защелки.** В Oracle эти блокировки применяются для защиты внутренних структур данных. Например, когда СУБД Oracle выполняет разбор запроса и генерирует оптимизированный план запроса, она будет защелкивать библиотечный кеш, чтобы поместить в него этот план для использования другими сеансами. Защелка (latch) — это легковесный низкоуровневый механизм сериализации, применяемый Oracle, который функционирует подобно блокировке. Пусть термин легковесный не вводит вас в заблуждение; вы увидите, что защелки являются общей причиной конкуренции в базе данных. Они являются легковесными в смысле реализации, но не в смысле воздействия.

Теперь мы более подробно рассмотрим конкретные типы блокировок внутри этих основных классов и последствия их использования. Типов блокировок намного больше, чем можно охватить в одной главе. В последующих разделах раскрываются самые распространенные и удерживаемые в течение продолжительного периода блокировки. Другие типы блокировок обычно устанавливаются на очень короткие промежутки времени.

Блокировки DML

Блокировки DML позволяют гарантировать, что в каждый момент времени изменять строку может только один пользователь, и никто другой не сможет удалить таблицу, с которой производится работа. Во время этой работы Oracle будет самостоятельно помещать эти блокировки более или менее прозрачным образом.

Блокировки TX (транзакций)

Блокировка TX выдается, когда транзакция инициирует свое первое изменение. В этот момент транзакция автоматически начинается (явно запускать транзакцию в Oracle не требуется). Эта блокировка удерживается до тех пор, пока транзакция не выполнит операцию COMMIT или ROLLBACK. Она выступает в качестве механизма организации очереди, чтобы другие сеансы могли дождаться завершения транзакции. Каждая модифицируемая или обрабатываемая с помощью SELECT FOR UPDATE строка в транзакции будет указывать на блокировку TX, связанную с этой транзакцией. Хотя такой подход выглядит затратным в смысле ресурсов, на самом деле это не так. Чтобы все стало понятным, необходимо иметь концептуальное представление о том, где находятся блокировки и как они управляются. В Oracle блокировки хранятся в виде атрибута данных (обзор формата блокировок Oracle представлен в главе 10). В Oracle отсутствует традиционный диспетчер блокировок, который вел бы длинный список всех строк, заблокированных в системе. Многие другие базы данных располагают таким диспетчером, поскольку в их средах блокировки являются дефицитным ресурсом, за расходом которого необходимо следить. Чем больше блокировок применяется, тем большим объемом ресурсов эти системы должны управлять, поэтому вовсе не удивительно, что эти системы беспокоятся о том, чтобы количество используемых блокировок не было слишком большим.

В базе данных с традиционным основанным на памяти диспетчером блокировок процесс блокировки строки предусматривает выполнение следующих действий.

1. Найти адрес строки, которую нужно заблокировать.
2. Установить связь с диспетчером блокировок (требуется сериализация, т.к. он представляет собой общую структуру в памяти).
3. Заблокировать список.
4. Выполнить поиск в списке, чтобы выяснить, не блокирует ли эту строку кто-то другой.
5. Создать в списке новую запись, свидетельствующую том, что вы блокируете данную строку.
6. Разблокировать список.

Теперь, когда строка заблокирована, ее можно модифицировать. Позже, после фиксации изменений, процедура должна быть продолжена, как описано ниже.

1. Снова установить связь с диспетчером блокировок.
2. Заблокировать список.
3. Просмотреть список и освободить все свои блокировки.
4. Разблокировать список.

Как видите, чем больше блокировок получено, тем больше времени тратится на выполнение этой операции, как до, так и после изменения данных. В Oracle процесс выглядит по-другому.

1. Найти адрес строки, которую нужно заблокировать.
2. Перейти к этой строке.
3. Заблокировать строку прямо здесь — по местоположению строки, а не в большом списке где-то в другом месте (с ожиданием завершения блокирующей транзакции, если строка уже заблокирована, если только не применяется опция NOWAIT).

Вот и все. Поскольку блокировка хранится как атрибут данных, Oracle не нуждается в традиционном диспетчере блокировок. Транзакция просто обратится к данным и заблокирует их (если они уже не заблокированы). Интересно, что при обращении к данным они могут выглядеть заблокированными, даже если это не так. При блокировании строк в Oracle строка указывает на копию идентификатора транзакции, хранящуюся с блоком, который содержит данные, а когда блокировка освобождается, этот идентификатор транзакции остается. Он уникален для транзакции и представляет номер сегмента отката, слот и номер последовательности. Обязанность по уведомлению других сеансов о владении этими данными (не всеми данными в блоке, а только одной изменяемой строкой) возлагается на блок, содержащий данные. Когда появляется другой сеанс, он обнаруживает идентификатор блокировки и, учитывая, что этот идентификатор представляет транзакцию, может быстро выяснить, активна ли по-прежнему транзакция, удерживающая блокировку. Если блокировка не активна, сеансу разрешается доступ к данным. Если блокировка все еще активна, сеанс запросит уведомления об освобождении этой блокировки. Таким образом, мы имеем дело с механизмом организации очереди: сеанс, запрашивающий

блокировку, будет помещен в очередь для ожидания завершения этой транзакции, после чего он получит данные.

Ниже приведен небольшой пример, демонстрирующий происходящее посредством трех таблиц V\$.

- V\$TRANSACTION, которая содержит записи для каждой активной транзакции.
- V\$SESSION, отображающая сеансы, в которые был произведен вход.
- V\$LOCK, которая содержит записи для всех находящихся в очереди удерживаемых блокировок, а также для сеансов, ожидающих освобождения блокировок. В этом представлении вы не увидите строк для каждой строки, заблокированной сеансом. Как было указано ранее, общего списка блокировок на уровне строк попросту не существует. Если сеанс имеет одну заблокированную строку в таблице EMP, то в представлении V\$LOCK будет присутствовать одна строка для этого сеанса, отражающая данный факт. Если сеанс имеет миллионы заблокированных строк в таблице EMP, то представление V\$LOCK по-прежнему будет содержать только одну строку для такого сеанса. Представление показывает, какими помещенными в очередь блокировками располагают отдельные сеансы.

Давайте сначала получим копии таблиц EMP и DEPT. Если они уже есть в вашей схеме, замените их следующими определениями:

```
EODA@ORA12CR1> create table dept
  2 as select * from scott.dept;
Table created.
Таблица создана.
```

```
EODA@ORA12CR1> create table emp
  2 as select * from scott.emp;
Table created.
Таблица создана.
```

```
EODA@ORA12CR1> alter table dept
  2 add constraint dept_pk
  3 primary key(deptno);
Table altered.
Таблица изменена.
```

```
EODA@ORA12CR1> alter table emp
  2 add constraint emp_pk
  3 primary key(empno);
Table altered.
Таблица изменена.
```

```
EODA@ORA12CR1> alter table emp
  2 add constraint emp_fk_dept
  3 foreign key (deptno)
  4 references dept(deptno);
Table altered.
Таблица изменена.
```

```
EODA@ORA12CR1> create index emp_deptno_idx
  2 on emp(deptno);
Index created.
Индекс создан.
```

Теперь запустим транзакцию:

```
EODA@ORA12CR1> update dept
2      set dname = initcap(dname);
4 rows updated.
4 строки обновлено.
```

Давайте выясним состояние системы в текущий момент. В этом примере предполагается, что система является однопользовательской; в противном случае представление V\$TRANSACTION может содержать много строк. Но даже в системе с единственным пользователем не удивляйтесь, если увидите в представлении V\$TRANSACTION более одной строки, т.к. многие фоновые процессы Oracle также могут выполнять транзакции.

```
EODA@ORA12CR1> select username,
2      v$lock.sid,
3      trunc(id1/power(2,16)) rbs,
4      bitand(id1,to_number('ffff','xxxx'))+0 slot,
5      id2 seq,
6      lmode,
7      request
8  from v$lock, v$session
9  where v$lock.type = 'TX'
10 and v$lock.sid = v$session.sid
11 and v$session.username = USER;
```

| USERNAME | SID | RBS | SLOT | SEQ | LMODE | REQUEST |
|----------|-------|-------|-------|-------|-------|---------|
| ----- | ----- | ----- | ----- | ----- | ----- | ----- |
| EODA | 22 | 2 | 27 | 21201 | 6 | 0 |

```
EODA@ORA12CR1> select XIDUSN, XIDSLOT, XIDSQN from v$transaction;
```

| XIDUSN | XIDSLOT | XIDSQN |
|--------|---------|--------|
| ----- | ----- | ----- |
| 2 | 27 | 21201 |

Особый интерес здесь представляют описанные ниже моменты.

- В таблице V\$LOCK значение LMODE равно 6, а номер запроса — 0. Согласно определению таблицы V\$LOCK в руководстве Oracle Database Reference (Справочник по базе данных Oracle), значение LMOD=6 представляет монопольную блокировку. Значение 0 запроса свидетельствует об отсутствии запроса и наличии блокировки.
- Эта таблица содержит только одну строку. Таблица V\$LOCK является скорее таблицей очереди, чем таблицей блокировок. Многие ожидают увидеть в ней четыре строки, поскольку есть четыре заблокированных строки. Однако следует помнить, что Oracle не хранит общий список всех заблокированных где-либо строк. Чтобы выяснить, заблокирована ли строка, необходимо перейти к ней.
- Я извлек столбцы ID1 и ID2 и провел над ними ряд манипуляций. СУБД Oracle необходимо было сохранить три 16-битных числа, но для этого имелось только два столбца. Поэтому первый столбец, ID1, содержит два числа из трех. Деление на 2^{16} посредством `trunc(id1/power(2,16)) rbs` и маскирование старших разрядов с помощью `bitand(id1,to_number('ffff','xxxx'))+0 slot` позволяет получить два числа, скрытые в этом одном числе.

- Значения RBS, SLOT и SEQ соответствуют информации из V\$TRANSACTION. Они представляют идентификатор транзакции.

Теперь запустим еще один сеанс с использованием того же самого имени пользователя, обновим несколько строк в таблице EMP, а затем попытаемся обновить таблицу DEPT:

```
EODA@ORA12CR1> update emp set ename = upper(ename);
14 rows updated.
14 строк обновлено.
EODA@ORA12CR1> update dept set deptno = deptno-10;
```

Мы заблокировали этот сеанс. Если снова запустить запросы к V\$, мы увидим следующий вывод:

```
EODA@ORA12CR1> select username,
2      v$lock.sid,
3      trunc(id1/power(2,16)) rbs,
4      bitand(id1,to_number('ffff','xxxx'))+0 slot,
5      id2 seq,
6      lmode,
7      request
8  from v$lock, v$session
9  where v$lock.type = 'TX'
10     and v$lock.sid = v$session.sid
11     and v$session.username = USER;
```

| USERNAME | SID | RBS | SLOT | SEQ | LMODE | REQUEST |
|----------|-----|-----|------|-------|-------|---------|
| EODA | 17 | 2 | 27 | 21201 | 0 | 6 |
| EODA | 22 | 2 | 27 | 21201 | 6 | 0 |
| EODA | 17 | 8 | 17 | 21403 | 6 | 0 |

```
EODA@ORA12CR1> select XIDUSN, XIDSLOT, XIDSQN from v$transaction;
```

| XIDUSN | XIDSLOT | XIDSQN |
|--------|---------|--------|
| 2 | 27 | 21201 |
| 8 | 17 | 21403 |

Здесь мы видим, что началась новая транзакция с идентификатором транзакции (8,17,21403). На этот раз новый сеанс с SID=17 имеет две строки в таблице V\$LOCK. Одна строка представляет блокировки, которыми он владеет (где LMOD=6). К сеансу также относится строка, в которой столбец REQUEST имеет значение 6. Это запрос монопольной блокировки. Интересно отметить, что значения RBS/SLOT/SEQ в данной строке запроса являются идентификатором транзакции *держателя* блокировки. Транзакция с SID=703 блокирует транзакцию с SID=7. Те же сведения можно вывести в более наглядном виде, просто выполнив рефлексивное соединение (self-join) для представления V\$LOCK:

```
EODA@ORA12CR1> select
2      (select username from v$session where sid=a.sid) blocker,
3      a.sid,
4      ' is blocking ',
5      (select username from v$session where sid=b.sid) blockee,
6      b.sid
```



```

7   from v$lock a, v$lock b
8   where a.block = 1
9         and b.request > 0
10        and a.id1 = b.id1
11        and a.id2 = b.id2;

```

| BLOCKER | SID | 'ISBLOCKING' | BLOCKEE | SID |
|---------|-------|--------------|---------|-------|
| ----- | ----- | ----- | ----- | ----- |
| EODA | 22 | is blocking | EODA | 17 |

Если теперь мы выполним фиксацию первоначальной транзакции (SID=22) и повторно запустим запрос блокировки, то обнаружим, что строка запроса исчезла:

```

EODA@ORA12CR1> select username,
2       v$lock.sid,
3       trunc(id1/power(2,16)) rbs,
4       bitand(id1,to_number('ffff','xxxx'))+0 slot,
5       id2 seq,
6       lmode,
7       request
8   from v$lock, v$session
9  where v$lock.type = 'TX'
10     and v$lock.sid = v$session.sid
11     and v$session.username = USER;

```

| USERNAME | SID | RBS | SLOT | SEQ | LMODE | REQUEST |
|----------|-------|-------|-------|-------|-------|---------|
| ----- | ----- | ----- | ----- | ----- | ----- | ----- |
| EODA | 17 | 8 | 17 | 21403 | 6 | 0 |

```

EODA@ORA12CR1> select XIDUSN, XIDSLOT, XIDSQN from v$transaction;

```

| XIDUSN | XIDSLOT | XIDSQN |
|--------|---------|--------|
| ----- | ----- | ----- |
| 8 | 17 | 21403 |

Строка запроса исчезла сразу же после того, как другой сеанс снял свою блокировку. Эта строка запроса играла роль механизма организации очереди. База данных может пробуждать заблокированные сеансы немедленно после завершения транзакции. Существуют значительно более наглядные способы отображения информации с помощью разнообразных инструментов с графическим пользовательским интерфейсом, но при крайней необходимости наличие знаний о таблицах, которые следует просматривать, очень полезно.

Тем не менее, прежде чем можно будет заявлять о хорошем понимании работы блокировки строк в Oracle, нужно рассмотреть последнюю тему: каким образом информация о блокировках и транзакциях управляется с самими данными. Это часть дополнительного пространства блока. В главе 10 мы погрузимся в детали формата блока, но пока достаточно сказать, что в верхней части блока базы данных находится определенное дополнительное пространство, в котором хранится таблица транзакций для этого блока. Таблица транзакций содержит по одной записи для каждой реальной транзакции, которая блокирует какие-то данные в блоке. Размер этой структуры управляется двумя параметрами физических атрибутов в операторе CREATE для объекта.

- **INITTRANS.** Начальный заранее определенный размер этой структуры. По умолчанию значение этого параметра равно 2 для индексов и таблиц.

- **MAXTRANS.** Максимальный размер, до которого может расти эта структура. По умолчанию максимальное значение этого параметра равно 255, а минимальное — 2. В Oracle 10g и последующих версиях *эта настройка объявлена устаревшей*, поэтому она больше не применяется. В этих выпусках значение MAXTRANS всегда равно 255.

По умолчанию каждый блок начинает свое существование с двумя слотами транзакций. Количество одновременно активных транзакций, которые когда-либо могут находиться в блоке, ограничено значением MAXTRANS и доступным пространством в блоке. При отсутствии в блоке пространства, достаточного для увеличения этой структуры, достичь 255 одновременно выполняемых транзакций может быть невозможно.

Мы можем искусственно продемонстрировать, как это работает, путем создания таблицы с множеством строк, упакованных в единственный блок, так чтобы блок оказался практически полным с самого начала; после первоначальной загрузки данных в нем останется очень мало пространства. Присутствие этих строк ограничит пределы роста таблицы транзакций из-за нехватки места. Я использовал размер блока в 8 Кбайт и тестировал этот конкретный пример во всех версиях — от Oracle9i Release 2 до Oracle 12c Release 1 — с теми же самыми результатами (таким образом, если размер блока у вас составляет 8 Кбайт, то вы должны быть способны воспроизвести пример). Начнем с создания упакованной таблицы. Я пробовал разные длины данных, пока не добрался до этого довольно специфичного случая:

```

EODA@ORA12CR1> create table t
2 ( x int primary key,
3   y varchar2(4000)
4 )
5 /
Table created.
Таблица создана.

EODA@ORA12CR1> insert into t (x,y)
2 select rownum, rpad(' ',148,' ')
3   from dual
4  connect by level <= 46;
46 rows created.
46 строк создано.

EODA@ORA12CR1> select length(y),
2         dbms_rowid.rowid_block_number(rowid) blk,
3         count(*), min(x), max(x)
4   from t
5  group by length(y), dbms_rowid.rowid_block_number(rowid);

```

| LENGTH(Y) | BLK | COUNT(*) | MIN(X) | MAX(X) |
|-----------|-------|----------|--------|--------|
| 148 | 23470 | 46 | 1 | 46 |

Итак, таблица имеет 46 строк, причем все находятся в одном и том же блоке. Я выбрал длину 148 символов, поскольку если бы длина была на 1 символ больше, то для хранения тех же 46 записей понадобилось бы два блока. Теперь необходим способ, который позволит посмотреть, что произойдет, когда множество транзакций попытаются одновременно блокировать данные в этом единственном блоке.

Для этого мы снова воспользуемся `AUTONOMOUS_TRANSACTION`, чтобы можно было обойтись одним сеансом и не запускать несколько параллельных сеансов `SQL*Plus`. Наша хранимая процедура заблокирует строку в таблице по первичному ключу, начиная с его значения 1 (первая вставленная запись). Если процедура получит блокировку на этой строке без необходимости ожидания (т.е. не будет блокирована), она просто увеличит значение первичного ключа на 1 и посредством рекурсии сделает все это еще раз. То есть второй вызов попытается заблокировать запись 2, третий вызов — запись 3 и т.д. Если процедуре придется ждать, она сгенерирует ошибку `ORA-00054` (ресурс занят) и будет выведено сообщение “locked out trying to select row <primary key value>” (блокировка при попытке выбора строки <значение первичного ключа>). Это указывает на то, что слоты транзакций в этом блоке исчерпались до блокировки всех строк. С другой стороны, если мы не найдем никакой строки для блокировки, это означает, что мы уже заблокировали каждую строку данного блока; в таком случае выводится сообщение об успешном выполнении (означающее, что таблица транзакций в заголовке блока сможет расти, чтобы вместить все транзакции). Ниже приведен код этой хранимой процедуры:

```

EODA@ORA12CR1> create or replace procedure do_update( p_n in number )
2  as
3      pragma autonomous_transaction;
4      l_rec t%rowtype;
5      resource_busy exception;
6      pragma exception_init( resource_busy, -54 );
7  begin
8      select *
9          into l_rec
10         from t
11        where x = p_n
12        for update NOWAIT;
13
14      do_update( p_n+1 );
15      commit;
16  exception
17  when resource_busy
18  then
19      dbms_output.put_line( 'locked out trying to select row ' || p_n );
20      commit;
21  when no_data_found
22  then
23      dbms_output.put_line( 'we finished - no problems' );
24      commit;
25  end;
26  /
Procedure created.
Процедура создана.

```

Вся магия сосредоточена в строке 14, где процедура рекурсивно вызывает себя с новым значением первичного ключа для многократного блокирования. Если вы запустите эту процедуру после наполнения таблицы 148-символьными строками, то должны будете наблюдать следующее:

```
EODA@ORA12CR1> exec do_update(1);
locked out trying to select row 38
PL/SQL procedure successfully completed.
блокировка при попытке выбора строки 38
Процедура PL/SQL успешно завершена.
```

Этот вывод показывает, что мы смогли заблокировать 37 строк, но на 38-й строке слоты транзакций исчерпались. К данному блоку одновременно могут обращаться максимум 37 транзакций. Если мы повторим пример со строкой чуть меньшего размера, то все пройдет удачно:

```
EODA@ORA12CR1> truncate table t;
Table truncated.
Таблица усечена.

EODA@ORA12CR1> insert into t (x,y)
2  select rownum, rpad('*',147,'*')
3  from dual
4  connect by level <= 46;
46 rows created.
46 строк создано.

EODA@ORA12CR1> select length(y),
2      dbms_rowid.rowid_block_number(rowid) blk,
3      count(*), min(x), max(x)
4  from t
5  group by length(y), dbms_rowid.rowid_block_number(rowid);
```

| LENGTH(Y) | BLK | COUNT(*) | MIN(X) | MAX(X) |
|-----------|-------|----------|--------|--------|
| 147 | 23470 | 46 | 1 | 46 |

```
EODA@ORA12CR1> exec do_update(1);
we finished - no problems
PL/SQL procedure successfully completed..
завершено – проблемы не возникали
Процедура PL/SQL успешно завершена.
```

На этот раз задача решена успешно — разница в единственный байт сработала! В данном случае наличие 46 дополнительных байтов свободного пространства в блоке (каждая из 46 строк на один байт меньше) позволила иметь в блоке, по крайней мере, на 9 активных транзакций больше.

Приведенный пример демонстрирует, что происходит, когда множество транзакций пытаются получить доступ к одному и тому же блоку одновременно — при исключительно высоком числе параллельных транзакций может возникнуть ожидание на таблице транзакций. Блокирование может случиться, если параметр `INITRANS` установлен в низкое значение и пространства в блоке недостаточно для динамического расширения транзакции. В большинстве случаев стандартного значения 2 для `INITRANS` вполне достаточно, т.к. таблица транзакций будет динамически расти (при наличии свободного пространства), но в некоторых средах может требоваться увеличение этого параметра (чтобы зарезервировать больше пространства для слотов) для повышения степени параллелизма и уменьшения времени ожидания.

Примером ситуации, когда имеет смысл увеличить `INITRANS`, может быть наличие таблицы или, даже еще чаще, индекса (поскольку блоки индекса могут содер-

жать значительно больше строк, чем обычно хранит таблица), который интенсивно изменяется и в среднем имеет много строк на блок. В таком случае может понадобиться увеличить значение либо параметра PCTFREE (обсуждается в главе 10), либо параметра INITRANS, чтобы заблаговременно выделить в блоке пространство, достаточное для предполагаемого количества параллельных транзакций. Это особенно справедливо, если ожидается, что с самого начала блоки будут почти полными, т.е. в них не будет места для динамического расширения структуры транзакций.

И последнее замечание о параметре INITRANS. Пару раз я утверждал, что стандартным значением для него является 2. Тем не менее, если вы просмотрите словарь данных после создания таблицы, то заметите, что для INITRANS отображается значение 1:

```
EODA@ORA12CR1> create table t ( x int );
EODA@ORA12CR1> select ini_trans from user_tables where table_name = 'T';

INI_TRANS
-----
1
```

Так чему же равно стандартное количество слотов транзакций — 1 или 2? Хотя словарь данных показывает значение 1, мы можем продемонстрировать, что в действительности им является 2. Проведем следующий эксперимент. Первым делом сгенирируем одну транзакцию для таблицы T, вставив в нее одиночную запись:

```
EODA@ORA12CR1> insert into t values ( 1 );
```

Теперь удостоверимся в том, что таблицей T потребляется один блок:

```
EODA@ORA12CR1> select dbms_rowid.ROWID_BLOCK_NUMBER(rowid) from t;
DBMS_ROWID.ROWID_BLOCK_NUMBER(ROWID)
-----
57715
```

Далее поместим в переменные B и F номер блока и номер файла данных блока, используемого таблицей T:

```
EODA@ORA12CR1> column b new_val B
EODA@ORA12CR1> column f new_val F
EODA@ORA12CR1> select dbms_rowid.ROWID_BLOCK_NUMBER(rowid) B,
2         dbms_rowid.ROWID_TO_ABSOLUTE_FNO( rowid, user, 'T' ) F
3         from t;
```

Создадим дамп блока, применяемого таблицей T:

```
EODA@ORA12CR1> alter system dump datafile &F block &B;
```

Поместим в переменную TRACE местоположение и имя трассировочного файла, содержащего информацию дампа для блока:

```
EODA@ORA12CR1> column trace new_val TRACE
EODA@ORA12CR1> select c.value || '/' || d.instance_name || '_ora_' ||
a.spid || '.trc' trace
2         from v$process a, v$session b, v$parameter c, v$instance d
3         where a.addr = b.paddr
4         and b.audsid = userenv('sessionid')
5         and c.name = 'user_dump_dest';
```

А теперь завершим работу сеанса и откроем трассировочный файл для редактирования:

```
EODA@ORA12CR1> disconnect
EODA@ORA12CR1> edit &TRACE
```

Выполнив поиск в трассировочном файле значения Itl, мы увидим, что были инициализированы два слота транзакций (хотя для этой таблицы была выдана только одна транзакция):

| Itl | Xid | Uba | Flag | Lck | Scn/Fsc |
|------|---------------------|--------------------|------|-----|---------------------|
| 0x01 | 0x0013.00e.000024be | 0x00c000bf.039e.2d | --U- | 1 | fsc 0x0000.01cfa56a |
| 0x02 | 0x0000.000.00000000 | 0x00000000.0000.00 | ---- | 0 | fsc 0x0000.00000000 |

Значение 1 для INITRANS, сообщаемое словарем данных, скорее всего, является унаследованным и в более новых версиях Oracle вместо него должно отображаться 2.

Блокировки TM (помещение в очередь DML)

Блокировки TM используются для гарантирования того, что структура таблицы не изменится во время модификации ее содержимого. Например, при обновлении таблицы вы запрашиваете блокировку TM на ней. Это предотвратит выдачу другим пользователем команд DROP или ALTER в отношении данной таблицы. Если другой пользователь попытается выполнить оператор DDL для этой таблицы в то время, когда на ней размещена блокировка TM, он получит следующее сообщение об ошибке:

```
drop table dept
      *
ERROR at line 1:
ORA-00054: resource busy and acquire with NOWAIT specified
ОШИБКА в строке 1:
ORA-00054: ресурс занят и запрошен без указания NOWAIT
```

На заметку! В Oracle 11g Release 2 и последующих версиях можно устанавливать DDL_LOCK_TIMEOUT, чтобы заставить DDL-оператор ожидать. Обычно это делается командой ALTER SESSION. Например, перед выдачей команды DROP TABLE можно выполнить команду ALTER SESSION SET DDL_LOCK_TIMEOUT=60;. Тогда команда DROP TABLE должна будет ждать 60 секунд, прежде чем возвратить сообщение об ошибке (разумеется, она также может выполниться успешно).

Сообщение ORA-00054 поначалу может сбивать с толку, поскольку прямого метода указать конструкцию NOWAIT или WAIT в DROP TABLE вообще не существует. Это просто универсальное сообщение, отображаемое при попытке выполнения операции, которая будет заблокирована, но не допускает блокирования. Как вы уже видели ранее, то же самое сообщение выводится при выдаче команды SELECT FOR UPDATE NOWAIT применительно к заблокированной строке.

Ниже показано, как эти блокировки будут выглядеть в таблице V\$LOCK:

```
EODA@ORA12CR1> create table t1 ( x int );
Table created.
Таблица создана.
```

```
EODA@ORA12CR1> create table t2 ( x int );
```

```
Table created.
```

Таблица создана.

```
EODA@ORA12CR1> insert into t1 values ( 1 );
```

```
1 row created.
```

1 строка создана.

```
EODA@ORA12CR1> insert into t2 values ( 1 );
```

```
1 row created.
```

1 строка создана.

```
EODA@ORA12CR1> select (select username
```

```
2         from v$session
```

```
3         where sid = v$lock.sid) username,
```

```
4         sid,
```

```
5         id1,
```

```
6         id2,
```

```
7         lmode,
```

```
8         request, block, v$lock.type
```

```
9     from v$lock
```

```
10    where sid = sys_context('userenv','sid');
```

| USERNAME | SID | ID1 | ID2 | LMODE | REQUEST | BLOCK | TY |
|-------------|-----------|---------------|----------|----------|----------|----------|-----------|
| ----- | ---- | ----- | ----- | ----- | ----- | ----- | -- |
| EODA | 22 | 133 | 0 | 4 | 0 | 0 | AE |
| EODA | 22 | 244271 | 0 | 3 | 0 | 0 | TM |
| EODA | 22 | 244270 | 0 | 3 | 0 | 0 | TM |
| EODA | 22 | 1966095 | 152 | 6 | 0 | 0 | TX |

```
EODA@ORA12CR1> select object_name, object_id
```

```
2     from user_objects
```

```
3     where object_id in (244271,244270);
```

```
OBJECT_NAM  OBJECT_ID
```

```
-----
```

```
T2          244271
```

```
T1          244270
```

На заметку! Блокировка AE — это блокировка редакции, доступная в Oracle 11g и последующих версиях. Она является частью средства Edition Based Redefinition (Переопределение на основе редакции), которое в настоящей книге не рассматривается. ID1 — идентификатор объекта редакции, SID которой используется в текущий момент. Эта блокировка редакции защищает соответствующую редакцию от модификации (вроде отбрасывания редакции) подобно тому, как блокировки TM защищают от структурных модификаций таблицы, на которые они указывают.

В то время как для каждой транзакции мы получаем только одну блокировку TX, блокировок TM может быть столько, сколько модифицируемых объектов. Здесь интересно отметить, что в столбце ID1 блокировки TM находится идентификатор объекта, заблокированного для операции DML, поэтому выявление объекта, на котором удерживается блокировка, не представляет труда.

Еще одна интересная особенность блокировок TM: общее количество блокировок TM, разрешенных в системе, допускает конфигурирование (за деталями обращайтесь к определению параметра DML_LOCKS в руководстве *Oracle Database*

Reference). Параметр `DML_LOCKS` в действительности может быть установлен в 0. Это не означает превращение базы данных в доступную только для чтения (без блокировок), а то, что операции DDL не разрешены. Подобная возможность полезна в очень специализированных приложениях, таких как реализации RAC, для сокращения объема работ по координации внутри экземпляра, которые иначе пришлось бы выполнять. Кроме того, с помощью команды `ALTER TABLE <имя_таблицы> DISABLE TABLE LOCK` отключается возможность выдачи блокировок ТМ на уровне отдельных объектов. Это служит быстрым способом минимизации случайного удаления таблицы, т.к. перед удалением таблицы придется заново включить блокировку таблицы. Данный прием можно применять также для обнаружения обсуждаемой ранее блокировки всей таблицы из-за наличия неиндексированного внешнего ключа.

Блокировки DDL

Блокировки DDL автоматически помещаются на объекты во время выполнения операции DDL с целью защиты этих объектов от изменений другими сеансами. Например, при выполнении DDL-операции `ALTER TABLE T` на таблицу `T` в *общем случае* будет помещена монопольная блокировка DDL, которая предотвращает помещение на эту таблицу блокировок DDL и ТМ другими сеансами.

На заметку! В версии Oracle 11g было модифицировано то, что привыкли считать правилом.

В прошлом оператор `ALTER TABLE T` устанавливал на таблицу монопольную блокировку DDL. В этом примере таблица `T` запрещает другим сеансам выполнение DDL-операций и запрашивание блокировок ТМ (используемых для изменения содержимого таблицы). Теперь многие команды `ALTER` могут выполняться оперативно — без предотвращения модификаций.

Блокировки DDL удерживаются на протяжении выполнения оператора DDL и освобождаются немедленно после его завершения. Фактически это достигается за счет помещения операторов DDL внутрь неявных фиксаций (или пары “фиксация/откат”). По этой причине в Oracle операторы DDL всегда фиксируются. Каждый оператор `CREATE`, `ALTER` и т.д. в действительности выполняется так, как показано в следующем псевдокоде:

```
Begin
  Commit;
  ОПЕРАТОР DDL
  Commit;
Exception
  Иначе выполнить откат;
End;
```

Таким образом, оператор DDL фиксируется всегда, даже если он завершился неудачно. Выполнение оператора DDL начинается с фиксации — помните об этом. Вначале выполняется фиксация, чтобы при необходимости отката оператора это не привело к откату транзакции. Запуск оператора DDL делает постоянной любую оставшуюся невыполненной работу, даже когда он не является успешным. Если вам необходимо выполнить оператор DDL, но вы не хотите фиксировать существующую транзакцию, можете воспользоваться автономной транзакцией.

Существуют три типа блокировок DDL.

- *Монопольные блокировки DDL.* Такие блокировки препятствуют другим сеансам самим получать блокировку DDL или TM (DML). Это значит, что во время выполнения операции DDL таблицу можно запрашивать, но не модифицировать ее каким-либо образом.
- *Разделяемые блокировки DDL.* Такие блокировки защищают структуру соответствующего объекта от изменений другими сеансами, но разрешают модификацию данных.
- *Прерываемые блокировки разбора.* Эти блокировки разрешают объекту, такому как план запроса, кешированный в разделяемом пуле, регистрировать свою зависимость от какого-то другого объекта. При выполнении операции DDL в отношении этого объекта Oracle просмотрит список объектов, которые зарегистрировали свою зависимость, и сделает их недействительными. Следовательно, такие блокировки являются прерываемыми — они не препятствуют выполнению операции DDL.

Большинство операций DDL захватывают монопольную блокировку DDL. Если вы запустите оператор вроде следующего:

```
alter table t move;
```

то таблица T будет недоступной для изменений на период его выполнения. В это время таблицу можно запрашивать с применением оператора SELECT, но выполнять большинство других действий, включая все операторы DDL, будет запрещено. Теперь в базе данных Oracle некоторые операции DDL могут выполняться без блокировок DDL. Например, можно выдать следующий оператор:

```
create index t_idx on t(x) ONLINE;
```

Ключевое слово ONLINE модифицирует метод, посредством которого в действительности строится индекс. Вместо того чтобы использовать монопольную блокировку DDL, предотвращающую изменение данных, Oracle будет пытаться получить для таблицы только низкоуровневую блокировку TM (режима 2). Это эффективно воспрепятствует выполнению других операций DDL, но сделает возможным нормальное выполнение операций DML. База данных Oracle достигает этого результата, сохраняя запись об изменениях, которые были произведены в таблице во время выполнения оператора DDL, и применяя эти изменения к новому индексу по завершении операции CREATE. В результате существенно увеличивается доступность данных. Чтобы убедиться в этом, создайте таблицу некоторого размера:

```
EODA@ORA12CR1> create table t as select * from all_objects;  
Table created.  
Таблица создана.
```

```
EODA@ORA12CR1> select object_id from user_objects where object_name = 'T';  
OBJECT_ID  
-----  
244277
```

Затем создайте индекс на этой таблице:

```
EODA@ORA12CR1> create index t_idx on t(owner,object_type,object_name) ONLINE;
```

В то же самое время запустите в другом сеансе приведенный ниже запрос, чтобы увидеть блокировки, установленные на этой вновь созданной таблице (не забывайте, что значение ID1=244277 специфично для моего случая, поэтому указывайте свой идентификатор объекта):

```
EODA@ORA12CR1> select (select username
2         from v$session
3         where sid = v$lock.sid) username,
4         sid,
5         id1,
6         id2,
7         lmode,
8         request, block, v$lock.type
9     from v$lock
10    where id1 = 244277
11 /
```

| USERNAME | SID | ID1 | ID2 | LMODE | REQUEST | BLOCK | TY |
|----------|-----|--------|-----|-------|---------|-------|----|
| EODA | 22 | 244277 | 0 | 3 | 0 | 0 | DL |
| EODA | 22 | 244277 | 0 | 3 | 0 | 0 | DL |
| EODA | 22 | 244277 | 0 | 2 | 0 | 0 | TM |
| EODA | 22 | 244277 | 0 | 4 | 0 | 0 | OD |

Итак, здесь мы видим четыре блокировки, установленные на нашем объекте. Две блокировки DL — это блокировки *прямой загрузки* (direct load). Они используются для предотвращения загрузки в прямом режиме в базовую таблицу, когда создается индекс (естественно, это подразумевает невозможность одновременного выполнения загрузки данных в прямом режиме и создания индекса). Блокировка OD относится к типу блокировок, впервые появившемуся в Oracle 11g (в версиях Oracle 10g и Oracle9i они отсутствуют), который разрешает по-настоящему оперативные команды DDL. В прошлом (в Oracle 10g и предшествующих версиях) оперативные команды DDL, такие как CREATE INDEX ONLINE, не были на 100% оперативными. Они требовали установки блокировки в начале и конце оператора CREATE, предотвращая другие параллельные действия (изменения данных базовой таблицы). Они были *по большей части оперативными*, но не *полностью оперативными*. Начиная с версии Oracle 11g, команда CREATE INDEX ONLINE полностью оперативна; она не требует монопольных блокировок в начале/конце команды. Отчасти реализация этого трюка стала возможной благодаря введению блокировки OD (Online DDL — оперативная команда DDL); она применяется внутренне для обеспечения по-настоящему оперативных операций DDL.

Другие типы операций DDL имеют дело с разделяемыми блокировками DDL. Эти блокировки устанавливаются на зависимые объекты при создании хранимых, компилируемых объектов, таких как процедуры и представления. Например, при выполнении следующего запроса разделяемые блокировки DDL будут помещены на таблицы EMP и DEPT на время обработки команды CREATE VIEW:

```
create view MyView
as
select emp.empno, emp.ename, dept.deptno, dept.dname
   from emp, dept
  where emp.deptno = dept.deptno;
```

Допускается изменять содержимое этих таблиц, но не их структуру.

Последний тип блокировки DDL — прерываемая блокировка разбора. Когда сеанс выполняет разбор оператора, блокировка разбора устанавливается на каждый объект, на который ссылается этот оператор. Эти блокировки получаются для того, чтобы проанализированный, кешированный оператор мог быть объявлен недействительным (очищен) в разделяемом пуле, если объект, на который он ссылается, был удален либо изменен.

При просмотре этой информации особенно полезным является представление DBA_DDL_LOCKS. Подходящих представлений V\$ не существует. Представление DBA_DDL_LOCKS строится на основе более загадочных таблиц X\$ и по умолчанию может быть не установлено в вашей базе данных. Это и другие представления, связанные с блокировкой, можно установить, запустив сценарий catblock.sql, который находится в каталоге [ORACLE_HOME]/rdbms/admin. Указанный сценарий должен выполняться от имени пользователя SYS. После его успешного завершения представление DBA_DDL_LOCKS можно запрашивать. Например, в только что подключенном сеансе я могу получить следующий вывод:

```
EODA@ORA12CR1> connect eoda/foo
```

```
Connected.
```

```
Подключено.
```

```
EODA@ORA12CR1> set linesize 1000
```

```
EODA@ORA12CR1> select session_id sid, owner, name, type,
2         mode_held held, mode_requested request
3   from dba_ddl_locks
4  where session_id = (select sid from v$mystat where rownum=1)
5 /
```

| SID | OWNER | NAME | TYPE | HELD | REQUEST |
|-----|-------|-----------------------|----------------------|-------|---------|
| 22 | SYS | DBMS_OUTPUT | Body | Null | None |
| 22 | SYS | DBMS_OUTPUT | Table/Procedure/Type | Null | None |
| 22 | EODA | EODA | 18 | Null | None |
| 22 | SYS | DBMS_APPLICATION_INFO | Body | Null | None |
| 22 | SYS | PLITBLM | Table/Procedure/Type | Null | None |
| 22 | SYS | DBMS_APPLICATION_INFO | Table/Procedure/Type | Null | None |
| 22 | EODA | EODA | 73 | Share | None |
| 22 | SYS | DATABASE | 18 | Null | None |

```
8 rows selected.
```

```
8 строк выбрано.
```

Все это объекты, блокируемые моим сеансом. Я имею прерываемые блокировки разбора на паре пакетов DBMS_*. Это побочный эффект от применения интерфейса SQL*Plus; например, он может вызывать DBMS_APPLICATION_INFO, когда вы изначально в него входите (чтобы включить/отключить DBMS_OUTPUT через команду SET SERVEROUTPUT). Я могу видеть здесь более одной копии различных объектов; это нормально и просто означает, что разделяемый пул содержит более одного объекта, которые ссылаются на эти объекты. Обратите внимание, что в этом представлении столбец OWNER представляет не владельца блокировки, а владельца заблокированного объекта. Именно потому мы видим множество строк SYS. Пользователь SYS владеет этими пакетами, но все они принадлежат моему сеансу.

Чтобы увидеть прерываемую блокировку разбора в действии, давайте создадим и запустим хранимую процедуру P:

```
EODA@ORA12CR1> create or replace procedure p
2 as
3 begin
4 null;
5 end;
6 /
```

Procedure created.

Процедура создана.

```
EODA@ORA12CR1> exec p
```

PL/SQL procedure successfully completed.

Процедура PL/SQL успешно завершена.

Процедура P теперь будет отображаться в представлении DBA_DDL_LOCKS. Мы имеем на ней блокировку разбора:

```
EODA@ORA12CR1> select session_id sid, owner, name, type,
2 mode_held held, mode_requested request
3 from dba_ddl_locks
4 where session_id = (select sid from v$mystat where rownum=1)
5 /
```

| SID | OWNER | NAME | TYPE | HELD | REQUEST |
|-----|-------|----------|----------------------|------|---------|
| 22 | EODA | P | Table/Procedure/Type | Null | None |
| ... | | | | | |
| 22 | SYS | DATABASE | 18 | Null | None |

9 rows selected.

9 строк выбрано.

Затем мы перекомпилируем процедуру и снова запросим представление:

```
EODA@ORA12CR1> alter procedure p compile;
```

Procedure altered.

Процедура изменена.

```
EODA@ORA12CR1> select session_id sid, owner, name, type,
2 mode_held held, mode_requested request
3 from dba_ddl_locks
4 where session_id = (select sid from v$mystat where rownum=1)
5 /
```

| SID | OWNER | NAME | TYPE | HELD | REQUEST |
|-----|-------|-----------------------|----------------------|-------|---------|
| 22 | SYS | DBMS_OUTPUT | Body | Null | None |
| 22 | SYS | DBMS_OUTPUT | Table/Procedure/Type | Null | None |
| 22 | EODA | EODA | 18 | Null | None |
| 22 | SYS | DBMS_APPLICATION_INFO | Body | Null | None |
| 22 | SYS | PLITBLM | Table/Procedure/Type | Null | None |
| 22 | SYS | DBMS_APPLICATION_INFO | Table/Procedure/Type | Null | None |
| 22 | | EODA | 73 | Share | None |
| 22 | SYS | DATABASE | 18 | Null | None |

8 rows selected.

8 строк выбрано.

Мы видим, что процедура Р исчезла из представления. Блокировка разбора была разрушена.

Представление `DBA_DDL_LOCKS` полезно для разработчика, когда обнаруживается, что какая-то часть кода не компилируется в тестовой системе или в системе разработки — она зависает и в конечном итоге прерывается. Это указывает на то, что ее использует кто-то другой (выполняет в текущий момент), и с помощью данного представления можно выяснить, что это за пользователь. То же самое будет происходить с операторами `GRANT` и другими типами операций `DDL`, запускаемыми применительно к объекту. Например, нельзя выдать полномочия `EXECUTE` для выполняющейся процедуры. Тот же самый метод можно использовать для выявления потенциальных блокирующих и ожидающих пользователей.

На заметку! В Oracle 11g Release 2 и последующих версиях имеется средство переопределения на основе редакции (EBR). С его помощью можно выдавать право `EXECUTE` и/или право на перекомпиляцию кода в базе данных, не мешая пользователям, которые выполняют код в текущий момент. Средство EBR позволяет иметь несколько версий одной и той же хранимой процедуры в схеме одновременно. Это обеспечивает возможность работы с копией процедуры в новой редакции (версии), не конфликтуя с текущей версией процедуры, которая применяется другими пользователями. Однако средство EBR в настоящей книге не рассматривается, а только упоминается там, где оно меняет установленные правила.

Зашелки

Зашелка (latch) — это легковесный механизм сериализации, используемый для координации многопользовательского доступа к разделяемым структурам данных, объектам и файлам.

Зашелки представляют собой блокировки, предназначенные для установки на очень короткие интервалы времени, например, на время, которое требуется для модификации структуры данных в памяти. Они применяются для защиты определенных структур памяти, таких как кеш буферов блоков базы данных или библиотечный кеш в разделяемом пуле. Зашелки обычно запрашиваются внутренне в режиме *добровольного ожидания*. Это означает, что если зашелка не доступна, то запрашивающий сеанс будет приостановлен на короткий промежуток времени и повторит операцию позже. Другие зашелки могут запрашиваться в немедленном режиме, что концептуально подобно выполнению оператора `SELECT FOR UPDATE NOWAIT`, т.е. процесс предпримет какие-то другие действия, например, попытается захватить эквивалентную родственную зашелку, которая может быть свободной, а не будет просто дожидаться, пока конкретная зашелка станет доступной. Поскольку в данный момент зашелку могут ожидать множество запрашивающих процессов, некоторые из них могут ждать дольше других. Зашелки назначаются довольно случайным образом, на основе везения, если уж на то пошло. Любой сеанс, запросивший зашелку сразу после того, как она была освобождена, получит ее. Никакой очереди ожидающих процессов не существует — есть только группа ожидающих процессов, которые постоянно повторяют свои попытки.

Для оперирования зашелками в Oracle используются атомарные инструкции типа “проверить и установить” и “сравнить и поменять местами”. Поскольку инструкции установки и освобождения зашелоков являются атомарными, сама операционная сис-

тема гарантирует, что проверять и устанавливать зашелку будет только один процесс, хотя одновременно она может требоваться нескольким процессам. С учетом того, что инструкция только одна, она может выполняться достаточно быстро (но весь алгоритм зашелкивания представляет собой множество инструкций центрального процессора). Защелки удерживаются в течение коротких периодов и предоставляют механизм для очистки в случае аварийного отказа держателя зашелки. Процесс очистки будет выполняться монитором процессов (PMON).

Рассмотренная ранее организация очередей — еще одно, более развитое средство последовательной обработки, которое применяется, например, при обновлении строк в таблице базы данных. Оно отличается от защепок тем, что позволяет запрашивающему процессу стать в очередь и ожидать ресурса. При запросе зашелки запрашивающий сеанс сразу получает ответ о том, получит он зашелку или нет. При использовании блокировки в порядке очереди запрашивающий сеанс будет блокироваться до тех пор, пока действительно не получит блокировку.

На заметку! В случае применения `SELECT FOR UPDATE NOWAIT` или `WAIT [n]` можно дополнительно отказаться от ожидания предоставления блокировки в порядке очереди, если сеанс окажется заблокированным, но если вы все же решите выполнять блокировку и ожидать, то ожидание будет происходить в очереди.

Сама по себе блокировка с очередью — не настолько быстрое средство, каким может быть зашелка, но по сравнению с ней предоставляет намного больше функциональных возможностей. Блокировки с очередями разрешено получать на разном уровне, поэтому может существовать множество разделяемых блокировок и блокировок с разными уровнями совместного использования.

“Раскручивание” зашелки

В отношении защепок я хотел бы разъяснить один момент: зашелки представляют собой тип блокировки, блокировки — это механизмы сериализации, а механизмы сериализации препятствуют масштабируемости. Если ваша цель заключается в построении приложения, которое хорошо масштабируется в среде Oracle, то вы должны искать подходы и решения, сводящие к минимуму необходимое количество защепок.

Даже на первый взгляд простые действия вроде разбора SQL-оператора получают и освобождают сотни или тысячи защепок на библиотечном кеше и связанных структурах в разделяемом пуле. Если мы располагаем зашелкой, значит, кто-то другой может ожидать ее освобождения. Когда мы обращаемся за зашелкой, то вполне можем ожидать ее предоставления.

Ожидание зашелки может оказаться дорогостоящей в смысле ресурсов операцией. Если зашелка не доступна сразу же и мы готовы ее ждать, как происходит в большинстве случаев, то на многопроцессорной машине сеанс будет выполнять раскрутку (spin), циклически снова и снова повторяя попытки получить зашелку. Такой подход обусловлен тем, что переключение контекста (т.е. отключение от процессора и повторное подключение к нему) является дорогостоящей операцией. Поэтому, если процесс не может получить зашелку немедленно, он остается подключенным к процессору и незамедлительно повторяет попытку, а не просто приостанавливается, освобождая процессор и повторяя попытку позднее при повторном подключении к

процессору в соответствии с установленным расписанием. При этом есть надежда, что держатель зашелки будет занят обработкой на другом процессоре (а поскольку зашелки разработаны так, чтобы удерживаться в течение очень коротких периодов времени, это весьма вероятно), и нам вскоре удастся получить зашелку. Если после раскручивания и непрерывных попыток получить зашелку это все же не удастся, то только тогда процесс перейдет в состояние бездействия, или отключится от процессора и позволит ему выполнять какую-то другую работу. Такой переход в состояние бездействия обычно является результатом одновременного запрашивания многими сеансами одной и той же зашелки. Проблема вовсе не в том, что одиночный сеанс удерживает зашелку долгое время, а в том, что множество сеансов хотят получить ее в тот же самый момент времени и каждый удерживает ее на протяжении короткого периода. Если вы делаете что-то краткое (быстрое) достаточно часто, то это все и объясняет! Псевдокод получения зашелки выглядит следующим образом:

```

Цикл
  для значения i из диапазона 1 .. 2000
  Цикл
    Попытаться получить зашелку
    Если зашелка получена, выполнить возврат
    Если i = 1, то счетчик неудач = счетчик неудач + 1
  Конец цикла
ИНКРЕМЕНТИРОВАТЬ СЧЕТЧИК ОЖИДАНИЯ
Бездействовать
ДОВАВИТЬ ВРЕМЯ ОЖИДАНИЯ
Конец цикла;
```

Логика работы предусматривает попытку получения зашелки и, если сделать это не удастся, то инкрементирование счетчика неудачных попыток, являющегося статистическими сведениями, которые можно видеть в отчете Statspack или путем непосредственного запрашивания представления `V$LATCH`. В случае неудачной попытки процесс циклически повторит попытки определенное количество раз (количеством попыток управляет недокументированный параметр, который обычно установлен в 2000), снова и снова пробуя получить зашелку. Если одна из этих попыток удастся, происходит выход и обработка продолжается. Если все попытки оказываются безуспешными, процесс на короткий период времени перейдет в состояние простоя, предварительно увеличив счетчик ожидания для этой зашелки. После пробуждения процесс повторяет все описанные действия заново. Это означает, что затраты на получение зашелки не ограничиваются только происходящей операцией типа “проверить и установить”, но также включают в себя значительную часть времени процессора, пока предпринимаются попытки получения зашелки. Система будет выглядеть очень загруженной (со значительным потреблением времени процессора), но при этом выполняется не особенно большой объем работ.

Измерение затрат на зашелкивание разделяемого ресурса

В качестве примера мы проанализируем затраты на установку зашелки на разделяемый пул. Мы сравним хорошо написанную программу (использующую переменные привязки) с программой, написанной не настолько удачно (в которой для выполнения каждого действия применяется литеральный, или уникальный, SQL-оператор). Для этого мы воспользуемся небольшой программой на Java, которая

просто входит в базу данных Oracle, отключает функцию автоматической фиксации (как и должны поступать все Java-программы немедленно после подключения к базе данных) и выполняет в цикле 25 000 уникальных операторов INSERT. Мы прогоним два набора тестов: в первом наборе программа не будет использовать переменные привязки, а во втором — будет.

Чтобы оценить эти программы и их поведение в многопользовательской среде, для сбора результатов измерений я предпочел использовать пакет Statspack, как описано ниже.

1. С помощью пакета Statspack получите снимок текущего состояния системы.
2. Запустите N копий программы, заставив каждую программу выполнять операцию INSERT в собственную таблицу базы данных, чтобы избежать конкуренции, связанной с попытками всех программ производить вставку в единственную таблицу.
3. Получите еще один снимок состояния системы немедленно после завершения последней копии программы.

Теперь остается только вывести отчет Statspack и выяснить, сколько времени потребовалось для завершения операции N копиями программы, сколько было потреблено времени процессора, какие основные события ожидания имели место и т.д.

На заметку! Почему бы для проведения этого анализа не применить AWR (Automatic Workload Repository — автоматический репозиторий рабочей нагрузки)? Дело в том, что доступ к пакету Statspack имеют все — действительно все. Он может требовать установки администратором базы данных, но доступ к нему имеет каждый пользователь Oracle. Я хочу представить результаты, которые могли бы быть воспроизведены абсолютно всеми.

Эти тесты были выполнены на двухпроцессорной машине с включенной гиперпоточностью (создавшей иллюзию наличия четырех процессоров). Учитывая наличие двух физических процессоров, здесь можно было ожидать весьма линейного масштабирования — то есть, если один пользователь для обработки своих операций вставки задействует 1 единицу загрузки процессора, то двоим пользователям потребуются 2 таких единицы. Вы обнаружите, что это предположение, хотя и кажется правдоподобным, может быть неточным (как вы убедитесь, степень расхождения зависит от используемой техники программирования). Предположение было бы корректным, если бы выполняемая обработка не нуждалась в разделяемых ресурсах, но наш процесс применяет такой ресурс, а именно — разделяемый пул. Нам необходимо зашелкивать разделяемый пул для выполнения разбора SQL-операторов, и это требуется делать из-за того, что он является совместно используемой структурой данных, и мы не можем его изменять, пока другие пользователи производят чтение из него, или выполнять чтение из него, когда другие его модифицируют.

На заметку! Я прогнал эти тесты с участием Java, PL/SQL, Pro*C и других языков программирования. Каждый раз конечные результаты были в основном такими же. Приведенная демонстрация и обсуждение применимо ко всем языкам и всем интерфейсам, предназначенным для взаимодействия с базой данных. В этом примере я выбрал язык Java, поскольку программы на Java и Visual Basic, скорее всего, не будут использовать переменные привязки при работе с базой данных Oracle.

Настройка теста

Для того чтобы провести тестирование, понадобится схема (набор таблиц), с которой мы будем иметь дело. Мы будем выполнять тестирование с множеством пользователей и хотим измерить в основном степень конкуренции из-за зашелкивания, т.е. мы не заинтересованы в измерении степени конкуренции, вызванной тем, что множество сеансов выполняют вставку в одну и ту же таблицу базы данных. Таким образом, мы создадим по одной таблице на пользователя и назовем их от T1 до T10. Например:

```
SCOTT@ORA12CR1> connect scott/tiger
Connected.
Подключено.

SCOTT@ORA12CR1> begin
2     for i in 1 .. 10
3     loop
4         for x in (select * from user_tables where table_name = 'T' || i )
5         loop
6             execute immediate 'drop table ' || x.table_name;
7         end loop;
8         execute immediate 'create table t' || i || ' ( x int )';
9     end loop;
10 end;
11 /

PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.
```

Мы будем запускать этот сценарий перед каждой итерацией теста, чтобы сбросить схему и инициировать полный разбор, если тест прогоняется более одного раза. Во время тестирования мы будем выполнять следующие шаги.

1. Запуск Statspack.snap.
2. Немедленный запуск N процедур Java, где N варьируется от 1 до 10, представляя от 1 до 10 параллельных пользователей.
3. Ожидание завершения всех N процедур.
4. Запуск Statspack.snap.
5. Генерация отчета Statspack для последних двух идентификаторов Statspack.

Числа, полученные в результате показанных ниже прогонов тестов, были собраны с применением этого приема.

Совет. Сценарии, предназначенные для автоматизации описанного здесь теста, доступны для загрузки на веб-сайте издательства. В каталоге ch06 есть два подкаталога: nobinds и binds. Внутри этих подкаталогов сценарий run.sql вызывает код, требующийся для выполнения этого теста. Вам понадобится изменить код, чтобы отразить актуальную информацию о подключении к базе данных там, где это необходимо. Кроме того, вы также должны, конечно же, скомпилировать Java-программу на своем сервере.

Тестирование без переменных привязки

В первом случае при вставке данных наша программа будет использовать не переменные привязки, а конкатенацию строк (очевидно, что вы должны указать свою строку подключения):

```
import java.sql.*;
public class instest
{
    static public void main(String args[]) throws Exception
    {
        DriverManager.registerDriver(new oracle.jdbc.driver.OracleDriver());
        Connection
            conn = DriverManager.getConnection
                ("jdbc:oracle:thin:@heesta:1521:ORA12CR1","scott","tiger");
        conn.setAutoCommit( false );
        Statement stmt = conn.createStatement();
        for( int i = 0; i < 25000; i++ )
        {
            stmt.execute
                ("insert into "+ args[0] +
                 " (x) values(" + i + ")");
        }
        conn.commit();
        conn.close();
    }
}
```

Я запустил этот тест в однопользовательском режиме (т.е. в отсутствие других активных сеансов базы данных), и пакет Statspack возвратил отчет с такой информацией:

| | | | | |
|--------------------------------|--------------------|---------------------|-------------|----------|
| Elapsed: | 0.25 (mins) | Av Act Sess: | 0.9 | |
| DB time: | 0.22 (mins) | DB CPU: | 0.20 (mins) | |
| Cache Sizes | Begin | End | | |
| ~~~~~ | ----- | ----- | | |
| Buffer Cache: | 2,656M | Std Block Size: | 8K | |
| Shared Pool: | 640M | Log Buffer: | 14,808K | |
| Load Profile | Per Second | Per Transaction | Per Exec | Per Call |
| ~~~~~ | ----- | ----- | ----- | ----- |
| ... | | | | |
| Parses: | 3,342.1 | 25,066.0 | | |
| Hard parses: | 1,667.2 | 12,504.0 | | |
| ... | | | | |
| Instance Efficiency Indicators | | | | |
| ~~~~~ | | | | |
| Buffer Nowait %: | 100.00 | Redo Nowait %: | 100.00 | |
| Buffer Hit %: | 99.99 | Optimal W/A Exec %: | 100.00 | |
| Library Hit %: | 60.05 | Soft Parse %: | 50.12 | |
| Execute to Parse %: | 0.11 | Latch Hit %: | 100.00 | |
| Parse CPU to Parse Elapsd %: | 108.72 | % Non-Parse CPU: | 16.29 | |
| ... | | | | |

| Top 5 Timed Events | | | Avg | %Total |
|-------------------------|-------|-----------|-------|-------------|
| ~~~~~ | | | wait | Call |
| Event | Waits | Time (s) | (ms) | Time |
| ----- | ----- | ----- | ----- | ----- |
| AQPC idle | 1 | 30 | 30010 | 36.1 |
| heartbeat redo informer | 15 | 15 | 1004 | 18.1 |
| LGWR worker group idle | 14 | 15 | 1055 | 17.8 |
| lreg timer | 4 | 12 | 3001 | 14.4 |
| CPU time | | 11 | | 13.4 |

В справочных целях в отчет была включена конфигурация SGA, но в данном случае интерес представляет следующая информация:

- затраченное время (Elapsed time) или время базы данных (DB time), равное приблизительно 15 секундам (0,25 минуты);
- 1667 полных разборов (Hard parses), выполненных в секунду;
- использованное время процессора (CPU time), составляющее 11 секунд.

Теперь если бы мы запустили две таких программы одновременно, то могли бы ожидать увеличения скорости полного разбора примерно до 3300 разборов в секунду (в конце концов, мы располагаем двумя процессорами) и потребляемого времени процессора — до 22 секунд. Давайте посмотрим:

| Elapsed: | 0.27 (mins) | Av Act Sess: | 1.6 | |
|--------------------------------|---------------|---------------------|-------------|-------------|
| DB time: | 0.44 (mins) | DB CPU: | 0.41 (mins) | |
| ... | | | | |
| Load Profile | Per Second | Per Transaction | Per Exec | Per Call |
| ~~~~~ | ----- | ----- | ----- | ----- |
| ... | | | | |
| Parses: | 6,259.8 | 33,385.3 | | |
| Hard parses: | 3,125.6 | 16,669.7 | | |
| ... | | | | |
| Instance Efficiency Indicators | | | | |
| ~~~~~ | | | | |
| Buffer Nowait %: | 100.00 | Redo NoWait %: | 100.00 | |
| Buffer Hit %: | 99.99 | Optimal W/A Exec %: | 100.00 | |
| Library Hit %: | 60.03 | Soft Parse %: | 50.07 | |
| Execute to Parse %: | 0.06 | Latch Hit %: | 98.41 | |
| Parse CPU to Parse Elapsed %: | 96.28 | % Non-Parse CPU: | 15.06 | |
| ... | | | | |
| Top 5 Timed Events | | | Avg | %Total |
| ~~~~~ | | | wait | Call |
| Event | Waits | Time (s) | (ms) | Time |
| ----- | ----- | ----- | ----- | ----- |
| CPU time | | 23 | | 32.7 |
| LGWR worker group idle | 18 | 16 | 876 | 22.8 |
| heartbeat redo informer | 15 | 15 | 1005 | 21.8 |
| lreg timer | 5 | 15 | 3001 | 21.7 |
| latch: shared pool | 15,076 | 0 | 0 | .6 |

Легко заметить, что скорость выполнения полных разборов слегка возросла, а время процессора увеличилось более чем в два раза. В чем дело? Ответ кроется в реализации защелкивания в Oracle. На этой многопроцессорной машине невозможность

немедленного получения зашелки приводит к раскрутке. Сама раскрутка потребляет время процессора. Процесс 1 много раз пытался получить зашелку на разделяемом пуле всего лишь для того, чтобы выяснить, что эту зашелку удерживает процесс 2; в итоге процессу 1 пришлось выполнять раскрутку и ожидать освобождения зашелки (загружая при этом процессор). Для процесса 2 справедливо обратное утверждение: он многократно обнаруживал, что процесс 1 удерживает зашелку на необходимых ему ресурсах. Таким образом, значительная часть процессорного времени была потрачена не на выполнение реальной работы, а на ожидание, пока ресурс станет доступным. Если пролистать отчет Statspack до раздела Latch Sleep Breakdown (Прекращение бездействия зашелки), можно найти следующие сведения:

| Latch Name | Requests | Misses | Sleeps | Gets |
|-------------|-------------|---------|---------|---------|
| ----- | ----- | ----- | ----- | ----- |
| shared pool | 2, 296, 041 | 75, 240 | 15, 267 | 60, 165 |

Обратили внимание на число 15 267 в столбце SLEEPS? Оно довольно близко соответствует количеству ожиданий (указанному в столбце Waits () и строке latch: shared pool (зашелки: разделяемый пул)) в предыдущем отчете Top 5 Timed Events (5 первых событий, касающихся времени).

На заметку! Количество засыпаний (sleep) близко соответствует количеству ожиданий; это может вызвать недоумение. Почему не точно соответствует? Причина в том, что действие по получению снимка не является атомарным; во время его выполнения запускается последовательность запросов, собирающих статистические сведения в таблицы, и каждый запрос начинается в слегка отличающийся момент времени. Таким образом, показатели события ожидания были собраны чуть раньше сведений о зашелкивании.

Наш отчет Latch Sleep Breakdown показывает количество раз, когда мы пытались получить зашелку и не смогли, оказавшись в цикле раскрутки. Это значит, что отчет Top 5 Timed Events отражает только верхушку айсберга, касающегося проблем с зашелками — 75 240 неудачных попыток (означающих раскрутку с повторением попыток получить зашелку) в нем не отображены. После изучения данного отчета можно и не догадаться о наличии проблемы с полным разбором, хотя она является весьма серьезной. Для выполнения 2 единиц работы понадобилось более 2 единиц загрузки процессора. Это всецело обусловлено потребностью в ресурсе совместного использования — разделяемом пуле. Такова природа зашелкивания.

Как видите, диагностирование проблем, связанных с применением зашелок, может оказаться очень сложной задачей, если только не уяснить механику их реализации. При беглом просмотре раздела Top 5 Timed Events отчета Statspack можно упустить из виду тот факт, что мы имеем дело с крупной проблемой масштабирования. Вскрыть существование этой проблемы позволит только более глубокий анализ раздела с информацией о зашелках в отчете Statspack.

Кроме того, обычно невозможно определить, сколько процессорного времени используется системой из-за раскрутки — в этом тесте с двумя пользователями можно выяснить только то, что обработка заняла 23 секунды процессорного времени и 75 240 раз не удалось получить зашелку на разделяемом пуле. Мы не знаем, сколько раз выполнялась раскрутка при каждой неудачной попытке получения зашелки, поэтому нет реального способа измерить процессорное время, затраченное на раскрутку.

ку, и процессорное время, в течение которого производилась полезная обработка. Чтобы вывести такую информацию, необходимо располагать несколькими измерительными точками.

На основе результатов тестирования, поскольку сравнение выполняется с однопользовательским случаем, можно заключить, что около 1 секунды процессорного времени было потрачено на раскрутку для получения зашелки в ожидании освобождения данного ресурса. Мы можем прийти к этому заключению, т.к. знаем, что одному пользователю требуется только 11 секунд процессорного времени, поэтому двум пользователям понадобится 22 секунды, а 23 (общее количество секунд процессорного времени) минус 22 равно 1.

Тестирование с переменными привязки

Теперь рассмотрим ту же ситуацию, что и представленная в предыдущем разделе, но на этот раз с применением программы, которая во время своей работы использует значительно меньше защепок. В этой Java-программе будут задействованы переменные привязки. Для этого мы заменим тип Statement типом PreparedStatement, проведем разбор одиночного оператора INSERT, а затем привяжем и многократно выполним PreparedStatement в цикле:

```
import java.sql.*;
public class instest
{
    static public void main(String args[]) throws Exception
    {
        System.out.println( "start" );
        DriverManager.registerDriver(new oracle.jdbc.driver.OracleDriver());
        Connection
            conn = DriverManager.getConnection
                ("jdbc:oracle:thin:@heesta:1521:ORA12CR1", "scott","tiger");
        conn.setAutoCommit( false );
        PreparedStatement pstmt =
            conn.prepareStatement
                ("insert into "+ args[0] + " (x) values(?)" );
        for( int i = 0; i < 25000; i++ )
        {
            pstmt.setInt( 1, i );
            pstmt.executeUpdate();
        }
        conn.commit();
        conn.close();
        System.out.println( "done" );
    }
}
```

Давайте взглянем на отчеты Statspack для одного и двух пользователей, как это было сделано в примере без переменных привязки. Результаты существенно отличаются. Вот отчет для одного пользователя:

| | | | |
|-----------------|--------------------|---------------------|-------------|
| Elapsed: | 0.07 (mins) | Av Act Sess: | 0.6 |
| DB time: | 0.04 (mins) | DB CPU: | 0.03 (mins) |

| Cache Sizes | Begin | End | | |
|--------------------------------|-------------|---------------------|----------|------------|
| ~~~~~ | ----- | ----- | | |
| Buffer Cache: | 2,656M | Std Block Size: | 8K | |
| Shared Pool: | 640M | Log Buffer: | 14,808K | |
| Load Profile | Per Second | Per Transaction | Per Exec | Per Call |
| ~~~~~ | ----- | ----- | ----- | ----- |
| ... | | | | |
| Parses: | 158.5 | 317.0 | | |
| Hard parses: | 29.8 | 59.5 | | |
| ... | | | | |
| Instance Efficiency Indicators | | | | |
| ~~~~~ | | | | |
| Buffer Nowait %: | 100.00 | Redo NoWait %: | 100.00 | |
| Buffer Hit %: | 98.99 | Optimal W/A Exec %: | 100.00 | |
| Library Hit %: | 96.14 | Soft Parse %: | 81.23 | |
| Execute to Parse %: | 97.72 | Latch Hit %: | 100.00 | |
| Parse CPU to Parse Elapsed %: | 87.10 | % Non-Parse CPU: | 71.58 | |
| ... | | | | |
| Top 5 Timed Events | | | Avg | %Total |
| ~~~~~ | | | wait | Call |
| Event | Waits | Time (s) | (ms) | Time |
| ----- | ----- | ----- | ----- | ----- |
| AQPC idle | 1 | 30 | 30004 | 66.6 |
| lreg timer | 2 | 6 | 3004 | 13.3 |
| heartbeat redo informer | 4 | 4 | 1006 | 8.9 |
| LGWR worker group idle | 12 | 4 | 331 | 8.8 |
| CPU time | | 1 | | 2.1 |

Изменение довольно резкое: уменьшение процессорного времени с 11 секунд в примере без переменных привязки до 1 секунды в текущем примере. Количество полных разборов, выполняемых в секунду, снизилось с 1667 до почти 29 (и, зная внутреннюю работу пакета Statspack, я могу утверждать, что большая их часть приходится именно на него). Даже общее время обработки значительно уменьшилось с почти 15 до 4 секунд (0,07 минуты). Без применения переменных привязки десять одиннадцатых процессорного времени тратилось на разбор SQL-кода (1 секунда против 11). Это было связано не только с зашелками, т.к. без использования переменных привязки значительная часть процессорного времени тратилась на разбор и оптимизацию SQL-запросов. Разбор SQL-операторов требует очень большой загрузки процессора, но вынуждать процессор тратить десять одиннадцатых своего времени на выполнение каких-то действий (разбор), которые в действительности для нас бесполезны — работы, которую делать не нужно — весьма расточительно.

Прогон теста для двух пользователей показывает, что результаты продолжают улучшаться:

| Elapsed: | 0.08 (mins) | Av Act Sess: | 0.9 | |
|---------------------|-------------|-----------------|-------------|----------|
| DB time: | 0.07 (mins) | DB CPU: | 0.07 (mins) | |
| ... | | | | |
| Load Profile | Per Second | Per Transaction | Per Exec | Per Call |
| ~~~~~ | ----- | ----- | ----- | ----- |
| ... | | | | |
| Parses: | 25.6 | 42.7 | | |
| Hard parses: | 0.8 | 1.3 | | |
| ... | | | | |

Instance Efficiency Indicators

```
~~~~~
      Buffer Nowait %: 100.00      Redo NoWait %: 100.00
      Buffer Hit  %:  99.93 Optimal W/A Exec %: 100.00
      Library Hit  %:  99.97      Soft Parse %:  96.88
      Execute to Parse %:  99.74      Latch Hit %:  99.99
Parse CPU to Parse Elapsed %: 100.00      % Non-Parse CPU: 99.66
```

...

| Top 5 Timed Events | | | Avg | %Total |
|-------------------------|-------|----------|-------|--------|
| | | | wait | Call |
| Event | Waits | Time (s) | (ms) | Time |
| ----- | ----- | ----- | ----- | ----- |
| AQPC idle | 1 | 30 | 30012 | 74.6 |
| heartbeat redo informer | 4 | 4 | 1010 | 10.0 |
| lreg timer | 1 | 3 | 3000 | 7.5 |
| CPU time | | 3 | | 7.3 |
| log file parallel write | 22 | 0 | 5 | .3 |

Время загрузки процессора увеличилось примерно в 2–3 раза по сравнению с тестовым сценарием, в котором участвовал один пользователь.

На заметку! Из-за округления 1 секунда процессорного времени в действительности представляет значение из диапазона от 0 до 2, а 3 — значение между 2 и 4.

Более того, процессорное время, потраченное двумя пользователями в случае применения переменных привязки, гораздо меньше половины процессорного времени, которое потребовалось одному пользователю в ситуации без переменных привязки! Когда я просмотрел раздел, посвященный зашелкам, в этом отчете Statspack, то обнаружил, что конкуренция за разделяемый пул и библиотечный кеш была настолько низкой, что она даже не стоила отражения в отчете. На самом деле более тщательное изучение отчета показало, что зашелка разделяемого пула запрашивалась 50 511 раз вместо более 2,2 миллиона раз в ходе выполнения предыдущего теста с двумя пользователями без переменных привязки:

| Latch Name | Requests | Misses | Sleeps | Gets |
|-------------|----------|--------|--------|-------|
| ----- | ----- | ----- | ----- | ----- |
| shared pool | 50,511 | 48 | 1 | 47 |

Сравнение производительности/масштабируемости

В табл. 6.1 приведена сводка по затратам процессорного времени в каждой реализации, а также показатели зашелкивания при увеличении количества пользователей свыше двух. Как видите, решение, в котором используется меньше защепок, будет намного лучше масштабироваться с ростом числа пользователей.

Интересно отметить, что 10 пользователей, применяющих переменные привязки (и в результате очень мало запросов защепок), используют столько же аппаратных ресурсов (центральный процессор), сколько и 1 пользователь, который не применяет переменные привязки (т.е. злоупотребляет зашелкой или обрабатывает больше, чем требуется).

Таблица 6.1. Сравнительные данные по использованию процессора с переменными привязки и без них

| Количество пользователей | Время процессора в секундах/общее время в минутах | | Количество запросов зашелки на разделяемом пуле | | Счетчик ожидания защелок/время ожидания в секундах | |
|--------------------------|---|------------------------|---|------------------------|--|------------------------|
| | Без переменных привязки | С переменными привязки | Без переменных привязки | С переменными привязки | Без переменных привязки * | С переменными привязки |
| 1 | 11/0,22 | 1/0,04 | 0 | 0 | 0/0 | 0/0 |
| 2 | 23/0,44 | 3/0,07 | > 2,2 млн. | > 50 тыс. | 15,0 тыс./0 | 0/0 |
| 3 | 35/0,65 | 4/0,10 | > 3,5 млн. | > 75 тыс. | 19,1 тыс./1 | 0/0 |
| 4 | 47/0,89 | 4/0,13 | > 4,7 млн. | > 95 тыс. | 33,9 тыс./1 | 0/0 |
| 5 | 58/1,13 | 4/0,15 | > 5,9 млн. | > 126 тыс. | 55,8 тыс./3 | 10/0 |
| 6 | 70/1,35 | 6/0,18 | > 7,1 млн. | > 152 тыс. | 61,1 тыс./3 | 22/0 |
| 7 | 82/1,56 | 7/0,21 | > 8,3 млн. | > 176 тыс. | 70,0 тыс./3 | 23/0 |
| 8 | 90/2,55 | 7/0,23 | > 9,5 млн. | > 201 тыс. | 121,3 тыс./42 | 28/0 |
| 9 | 105/2,21 | 8/0,27 | > 10,7 млн. | > 226 тыс. | 111,6 тыс./13 | 25/0 |
| 10 | 115/2,49 | 8/0,28 | > 11,9 млн. | > 252 тыс. | 123,1 тыс./17 | 41/0 |

* **Примечание:** в столбце “Без переменных привязки” значения счетчика ожидания защелок выражается в тысячах, а в столбце “С переменными привязки” — в единицах.

Проанализировав результаты для 10 пользователей, вы увидите, что отказ от переменных привязки приводит к использованию почти в 14 раз больше процессорного времени, а общее время выполнения увеличивается, по крайней мере, в 9 раз по сравнению с решением с переменными привязки. Чем больше со временем добавляется пользователей, тем дольше каждому из них приходится ожидать освобождения защелок. Среднее время ожидания защелок увеличивается с 0,6 секунды на сеанс (3 секунды ожидания на 5 сеансов) для 5 пользователей до 1,7 секунды на сеанс для 10 пользователей. Однако реализация, которая избегает злоупотребления зашелками, не испытывает никакого отрицательного влияния при масштабировании в сторону увеличения числа пользователей.

Семафоры

Семафор (mutex) представляет собой механизм сериализации, во многом подобный зашелке; в действительности *mutex* означает *mutual exclusion* (взаимное исключение). Это еще один инструмент сериализации, применяемый базой данных, который был введен в версии Oracle 10g Release 1 и используется вместо традиционных защелок во многих местах сервера. Семафор отличается от зашелки тем, что его реализация еще более легковесна. Он требует меньше кода для реализации — примерно одну пятую объема инструкций (что дает в результате меньше обращений к процессору), и также требует меньше памяти — приблизительно одну седьмую часть. В дополнение к легковесности семафор в некоторых отношениях чуть менее функционален.

Подобно тому, как блокировка с очередью намного сложнее зашелки, зашелка сложнее семафора. Но, как и при сравнении блокировки с очередью и зашелки, в некоторых случаях зашелка может делать больше, чем семафор (как и блокировка с очередью в ряде ситуаций позволяет сделать больше, чем зашелка). Это значит, что не каждая зашелка может или должна быть заменена семафором, равно как и не каждая блокировка с очередью — зашелкой.

Читая о семафорах в разных источниках, всегда помните, что они являются легковесным механизмом сериализации. Они делают возможной более высокую масштабируемость, чем зашелка (подобно тому, как зашелки обеспечивают большую масштабируемость, чем блокировки с очередью), но все равно остаются *механизмом сериализации*. Если вы можете обойтись без действия, требующего применения семафора, то обычно должны так и поступать по той же самой причине, по которой, где только возможно, следует избегать запрашивания зашелки.

Блокировка вручную и блокировки, определенные пользователем

До сих пор мы рассматривали главным образом блокировки, которые Oracle устанавливает автоматически. При обновлении таблицы Oracle помещает на нее блокировку TM, чтобы предотвратить удаление таблицы другими сеансами (или, фактически, выполнять большинство операций DDL). Существуют блокировки TX, которые остаются на различных модифицируемых нами блоках, что позволяет другим пользователям выяснить, какими данными мы владеем. База данных задействует блокировки DDL для защиты объектов от изменения, пока мы сами их изменяем. Она внутренне использует зашелки и блокировки для защиты собственной структуры.

Давайте посмотрим, как можно вмешаться в некоторые из этих действий по блокировке. Мы располагаем перечисленными ниже возможностями:

- блокировка данных вручную с помощью SQL-оператора;
- создание собственных блокировок посредством пакета DBMS_LOCK.

В последующих разделах кратко обсуждаются причины, по которым могут требоваться эти возможности.

Блокировка вручную

Фактически мы уже видели пару случаев, когда может понадобиться ручная блокировка. Оператор `SELECT...FOR UPDATE` является доминирующим методом блокирования данных вручную. Мы применяли его в предшествующих примерах для предотвращения проблемы потерянных обновлений, когда один сеанс мог бы перезаписать изменения, внесенные другим сеансом. Он использовался в качестве метода сериализации доступа к записям подробностей, чтобы обеспечить соблюдение бизнес-правил (скажем, в примере графика резервирования ресурсов в главе 1).

Данные можно также блокировать вручную с помощью оператора `LOCK TABLE`. Этот оператор применяется редко, из-за крупнозернистости такой блокировки. Он просто блокирует всю таблицу, а не строки в ней. Если вы начинаете изменять строки, они будут заблокированы обычным образом. Следовательно, это не метод, позволяющий сэкономить ресурсы (как это может быть в других СУБД). Оператор `LOCK TABLE IN EXCLUSIVE MODE` можно использовать при реализации объемного пакетного обновления, которое будет воздействовать на большинство строк в конк-

ретной таблице, и нужна уверенность в том, что никто не сможет вас заблокировать. Блокируя таблицу в такой манере, вы можете гарантировать, что обновление сможет выполнить всю свою работу, не оказавшись заблокированным другими транзакциями. Тем не менее, приложения, содержащие оператор `LOCK TABLE`, встречаются довольно редко.

Создание собственных блокировок

В действительности через пакет `DBMS_LOCK` СУБД Oracle предоставляет разработчикам механизм блокирования в порядке очереди, который она применяет внутренне. Часто возникает вопрос, зачем может понадобиться создание собственных блокировок. Как правило, ответ зависит от конкретного приложения. Например, этот пакет может требоваться для обеспечения последовательного доступа к ряду внешних ресурсов Oracle. Предположим, что вы используете процедуру `UTL_FILE`, которая позволяет выполнять запись в файл в файловой системе сервера. Вы могли разработать общую процедуру подготовки сообщения, которую каждое приложение вызывает для записи сообщений. Поскольку файл является внешним, Oracle не будет координировать множество пользователей, которые одновременно пытаются его изменить. В этой ситуации пригодится пакет `DBMS_LOCK`. Теперь, прежде чем открывать, выполнять запись и закрывать файл, вы будете запрашивать монопольную блокировку от имени файла, а после закрытия файла — вручную освобождать эту блокировку. В результате в каждый момент времени только один пользователь сможет записывать сообщение в этот файл. Запросы от всех остальных пользователей будут помещены в очередь. Пакет `DBMS_LOCK` позволяет освобождать блокировку вручную по окончании работы с объектом, автоматически при выполнении фиксации или даже сохранять ее на все время нахождения в системе.

Резюме

Эта глава содержит огромный объем материала, который временами способен довести до головной боли. Хотя блокирование само по себе прямолинейно, подобное нельзя сказать о некоторых его побочных эффектах. Однако очень важно понимать эти проблемы. Например, если вы не знаете о блокировке таблицы, применяемой Oracle для поддержки отношения внешнего ключа, когда внешний ключ не индексирован, то ваше приложение будет страдать от низкой производительности. Если вы не понимаете, каким образом просматривать словарь данных, чтобы увидеть, кто кого блокирует, то можете никогда не выяснить этот вопрос. Вы просто придете к выводу, что иногда база данных зависает. Иногда я жалею, что не получал по доллару каждый раз, когда мне удавалось решать неустранимую проблему зависания, просто запуская запрос для выявления неиндексируемых внешних ключей и подсказывая о необходимости их индексации. Я стал бы богачом.

Параллелизм и многоверсионность

Как констатировалось в предыдущей главе, одной из ключевых проблем при разработке многопользовательских приложений, работающих с базами данных, является доведение до максимума уровня параллелизма доступа, но в то же самое время и обеспечение для каждого пользователя возможностей чтения и модификации данных в согласованной манере. В этой главе мы детально рассмотрим, каким образом Oracle достигает *многоверсионной согласованности чтения*, и что это означает для разработчика. Будет также представлен новый термин *согласованность записи*, который используется для описания работы Oracle не только в среде чтения с согласованностью прочитанных данных, но также в смешанной среде чтения/записи.

Понятие управления параллелизмом

Управление параллелизмом — это коллекция функций, предоставляемых базой данных для разрешения многим пользователям одновременного доступа к данным и их модификации. В предыдущей главе отмечалось, что *блокировка* — это один из основных механизмов, посредством которого Oracle регулирует параллельный доступ к разделяемым ресурсам базы данных и предотвращает взаимное влияние друг на друга параллельно выполняющихся транзакций базы данных. Подводя краткие итоги, можно сказать, что Oracle применяет разнообразные блокировки, включая перечисленные ниже.

- **Блокировки ТХ (транзакции).** Эти блокировки устанавливаются на протяжении транзакции, модифицирующей данные.
- **Блокировки ТМ (помещение в очередь DML) и блокировки DDL.** Эти блокировки гарантируют, что структура объекта не изменится во время модификации его содержимого (блокировка ТМ) или структуры самого объекта (блокировка DDL).
- **Зашелки и семафоры.** Это внутренние блокировки, используемые Oracle при посредничестве доступа к разделяемым структурам данных. На зашелки и семафоры в этой главе мы будем ссылаться просто как на зашелки, хотя в зависимости от версии Oracle они могут быть реализованы в виде семафоров операционной системы.

В каждом случае с установкой блокировки связаны минимальные накладные расходы. Блокировки ТХ исключительно масштабируемы в отношении как производительности, так и кардинальности. Блокировки ТМ и DDL применяются в наименее ограничивающем режиме всякий раз, когда это возможно. Зашелки и очереди являются очень легковесными и быстрыми (из них очереди немного сложнее, хотя обладают более развитыми возможностями). Проблемы возникают только из-за неудачно спроектированных приложений, которые удерживают блокировки *дольше, чем необходимо*, и вызывают блокирование в базе данных. Если вы аккуратно проектируете свой код, то механизмы блокировки Oracle позволят строить масштабируемые приложения с высокой степенью параллелизма.

На заметку! Выше используется формулировка “дольше, чем необходимо”. Это не значит, что вы должны пытаться зафиксировать (завершить) транзакцию как можно скорее. Транзакции должны длиться столько, сколько нужно, но не дольше. То есть транзакция является единицей работы — либо все, либо ничего. Вы должны осуществлять фиксацию, когда единица работы завершена, но не раньше и не позже!

Однако поддержка параллелизма в Oracle не ограничивается эффективной блокировкой. В Oracle реализована архитектура *многоверсионности* (представленная в главе 1), которая обеспечивает управляемый, но с высокой степенью параллелизма доступ к данным. Многоверсионность характеризует способность Oracle параллельно материализовывать множество версий данных и является механизмом, с помощью которого Oracle обеспечивает согласованные по чтению представления данных (т.е. согласованные результаты на определенный момент времени). Довольно приятный побочный эффект многоверсионности состоит в том, что процесс чтения данных никогда не блокируется процессом записи данных. Другими словами, операции записи не блокируют операции чтения. В этом заключается одно из фундаментальных отличий Oracle от других баз данных. Запрос, который только читает информацию, в Oracle никогда не будет блокироваться, он никогда не попадет в состояние взаимоблокировки с другим сеансом и никогда не выдаст ответ, который в базе данных отсутствует.

На заметку! Во время обработки распределенной *двухфазной фиксации* (Two Phase Commit) существует короткий период, когда Oracle будет предотвращать доступ по чтению к информации. Поскольку эта обработка представляет собой редкое и необычное явление (проблема касается только запросов, которые запускаются между фазами подготовки и фиксации и пытаются прочитать данные до поступления команды фиксации), подробно здесь она не рассматривается.

Модель многоверсионности Oracle для обеспечения согласованности чтения по умолчанию применяется на *уровне операторов* (для каждого запроса) и также может применяться на *уровне транзакций*. Это значит, что каждый оператор SQL, отправленный базе данных, видит согласованное по чтению представление базы — по крайней мере. Если такое согласованное по чтению состояние базы данных необходимо получить на уровне транзакции (набора SQL-операторов), то это также можно сделать, как будет показано в разделе “Уровень изоляции *SERIALIZABLE*” далее в главе.

Основным назначением транзакции является перевод базы из одного согласованного состояния в другое. В стандарте ISO языка SQL описаны разнообразные *уровни изоляции транзакций*, которые определяют, насколько *чувствительна* одна транзакция к изменениям, произведенным другой транзакцией. Чем выше уровень чувствительности, тем выше степень изоляции, которую база данных должна обеспечивать между транзакциями, выполняемыми вашим приложением. В следующем разделе мы посмотрим, как через архитектуру многоверсионности с абсолютно минимальным блокированием Oracle может поддерживать каждый из определенных уровней изоляции.

Уровни изоляции транзакций

Стандарт ANSI/ISO языка SQL определяет четыре уровня изоляции транзакций с разными возможными исходами для одного и того же сценария транзакций. То есть одна и та же работа, выполненная в той же манере с теми же самыми входными данными, в зависимости от выбранного уровня изоляции может приводить к получению разных ответов. Эти уровни изоляции определены в терминах трех “феноменов”, которые либо разрешены, либо запрещены на заданном уровне изоляции.

- **Грязное чтение (dirty read).** Смысл этого термина так же плох, как он звучит. Вам разрешено читать незафиксированные, или грязные, данные. Вы достигаете такого эффекта, просто открывая файл операционной системы, куда кто-то другой производит запись или чтение, какие бы данные там не оказались. При этом не обеспечивается целостность данных, нарушаются внешние ключи и игнорируются ограничения уникальности.
- **Невоспроизводимое чтение (nonrepeatable read).** Это означает, что если вы читаете строку в момент времени T1 и затем считываете ее повторно в момент времени T2, то за данный промежуток времени строка может измениться, исчезнуть, обновиться и т.д.
- **Фантомное чтение (phantom read).** Это значит, что если вы запустили запрос в момент времени T1, а затем выполните его повторно в момент времени T2, то в базе данных могут появиться дополнительные строки, что повлияет на полученные результаты. Отличие фантомного чтения от невозпроизводимого чтения состоит в том, что уже прочитанные данные не изменяются, но критерию запроса удовлетворяет больший объем данных, чем было до того.

На заметку! Стандарт ANSI/ISO языка SQL определяет характеристики уровня *транзакции*, а не характеристики на уровне индивидуальных операторов. В последующих разделах мы будем исследовать изоляцию уровня транзакций, а не только изоляцию уровня операторов.

Уровни изоляции SQL определены на основе того, допускают ли они каждый из описанных ранее феноменов. Интересно отметить, что стандарт SQL не навязывает специфичную схему блокирования и не регламентирует определенное поведение, а взамен описывает уровни изоляции в терминах этих феноменов, делая возможным существование разнообразных механизмов блокирования/параллелизма (табл. 7.1).

Таблица 7.1. Уровни изоляции ANSI

| Уровень изоляции | Грязное чтение | Невоспроизводимое чтение | Фантомное чтение |
|------------------|----------------|--------------------------|------------------|
| READ UNCOMMITTED | Разрешено | Разрешено | Разрешено |
| READ COMMITTED | — | Разрешено | Разрешено |
| REPEATABLE READ | — | — | Разрешено |
| SERIALIZABLE | — | — | — |

В Oracle явно поддерживаются уровни изоляции `READ COMMITTED` и `SERIALIZABLE` в том виде, как они определены стандартом. Однако это еще не все. Стандарт SQL пытался установить уровни изоляции, которые допускают разнообразные степени согласованности запросов, выполняемых на каждом уровне. `REPEATABLE READ` — это уровень изоляции, при котором стандарт SQL требует гарантии согласованного по чтению результата, получаемого из запроса. В определении этого стандарта `READ COMMITTED` не дает согласованные результаты, а `READ UNCOMMITTED` является уровнем, используемым для обеспечения неблокирующих чтений.

Тем не менее, в Oracle уровень `READ COMMITTED` обладает всеми атрибутами, необходимыми для достижения согласованных по чтению запросов. Во многих других базах данных запросы `READ COMMITTED` могут и будут возвращать ответы, которые никогда не существовали в базе ни в какой момент времени. Более того, Oracle также поддерживает дух уровня `READ UNCOMMITTED`. Цель предоставления грязного чтения заключается в обеспечении неблокирующего чтения, когда запрос не блокируется и не блокирует обновления некоторых данных. Однако для достижения этой цели Oracle не нуждается в грязных чтениях и не поддерживает их. Грязные чтения вынуждены реализовывать другие СУБД, чтобы предоставить возможность неблокирующих чтений.

Вдобавок к перечисленным выше четырем уровням изоляции SQL в Oracle предлагается еще один уровень — `READ ONLY`. Транзакция `READ ONLY` является эквивалентом транзакции `REPEATABLE READ` или `SERIALIZABLE`, которая не выполняет никаких модификаций в SQL.

Транзакция, применяющая уровень изоляции `READ ONLY`, видит только те изменения, которые были зафиксированы *на момент ее начала*, но вставки, обновления и удаления в этом режиме не разрешены (другие сеансы могут обновлять данные, но не транзакция `READ ONLY`). Используя этот режим, можно достичь уровней изоляции `REPEATABLE READ` и `SERIALIZABLE`.

Давайте теперь посмотрим, как именно многоверсионность и согласованность чтения вписывается в схему изоляции, и каким образом базы данных, не поддерживающие многоверсионность, добиваются тех же результатов. Эта информация будет полезна всем, кто работает с другими базами данных и уверен, что понимает, как должны работать уровни изоляции. Интересно также увидеть, каким образом стандарт, предназначенный для устранения различий между базами данных — ANSI/ISO SQL — в действительности позволяет делать это. Стандарт, хотя он и очень детальный, может быть реализован совершенно разными способами.

Уровень изоляции **READ UNCOMMITTED**

Уровень изоляции **READ UNCOMMITTED** допускает грязные чтения. В Oracle грязные чтения не применяются и даже не разрешены. Основная цель уровня изоляции **READ UNCOMMITTED** — предоставить основанное на стандартах определение, которое обслуживает неблокирующие чтения. Как уже было показано, Oracle обеспечивает неблокирующие чтения по умолчанию. Вы бы с трудом решились на блокировку запроса **SELECT** в базе данных (как отмечалось ранее, имеется специальный случай распределенной транзакции). Каждый одиночный запрос, будь то **SELECT**, **INSERT**, **UPDATE**, **MERGE** или **DELETE**, выполняется в согласованной по чтению манере. Причисление оператора **UPDATE** к запросам может показаться забавным, но это так и есть. Операторы **UPDATE** содержат два компонента: компонент чтения, определяемый конструкцией **WHERE**, и компонент записи, определяемый конструкцией **SET**. Операторы **UPDATE** выполняют чтение и запись в базу данных; такой возможностью обладают все операторы **DML**. Случай однострочного оператора **INSERT**, использующего конструкцию **VALUES**, является единственным исключением — операторы подобного рода не имеют компонента чтения, а только компонент записи.

В главе 1 метод обеспечения согласованного чтения Oracle был продемонстрирован на примере простого однотабличного запроса, который извлекал строки, удаленные *после* открытия курсора. Теперь мы обратимся к реальному примеру и посмотрим, что происходит в Oracle за счет применения многоверсионности, а также каким образом ведут себя в аналогичной ситуации другие базы данных.

Начнем с той же самой элементарной таблицы и запроса:

```
create table accounts
( account_number number primary key,
  account_balance number not null
);
select sum(account_balance) from accounts;
```

Предположим, что перед началом запроса мы имеем данные, приведенные в табл. 7.2.

Таблица 7.2. Таблица **ACCOUNTS перед модификацией**

| Строка | Номер счета | Баланс счета |
|---------|-------------|--------------|
| 1 | 123 | \$500.00 |
| 2 | 456 | \$240.25 |
| ... | ... | ... |
| 342 023 | 987 | \$100.00 |

Теперь наш оператор **SELECT** начинает свое выполнение и читает строку 1, строку 2 и т.д.

На заметку! В этом примере я вовсе не предполагаю, что строки размещаются на диске в каком-нибудь физическом порядке. На самом деле в таблице не существует первой, второй или последней строки. Есть просто набор строк. Здесь мы предполагаем, что строка 1 в действительности означает “первую прочитанную строку”, строка 2 — “вторую прочитанную строку” и т.д.

В какой-то момент посреди запроса другая транзакция переводит \$400 со счета 123 на счет 987. Эта транзакция делает два обновления, но не фиксируется. Таблица теперь выглядит так, как показано в табл. 7.3.

Таблица 7.3. Таблица ACCOUNTS во время модификации

| Строка | Номер счета | Баланс счета | Блокирована? |
|---------|-------------|---|--------------|
| 1 | 123 | Значение (\$500.00) изменилось на \$100.00 | X |
| 2 | 456 | \$240.25 | — |
| ... | ... | ... | ... |
| 342 023 | 987 | Значение (\$100.00) изменилось на \$500.00 | X |

Итак, заблокированы две строки. Если кто-то попытается обновить их, эта попытка будет заблокирована. До сих пор наблюдаемое нами поведение является более или менее согласованным во всех базах данных. Отличия начнут проявляться в том, что произойдет, когда запрос доберется до заблокированных данных.

Когда функционирующий запрос достигает блока данных, содержащего заблокированную строку (строку 342 023) в нижней части таблицы, он обнаружит, что данные в этой строке изменились с момента его запуска на выполнение. Для предоставления согласованного (корректного) ответа Oracle создаст в этой точке копию блока, в котором находится данная строка, *в том виде, в каком он существовал на момент начала запроса*. То есть будет прочитано значение \$100.00 — значение, которое было при запуске запроса. В сущности Oracle обходит модифицированные данные, в процессе чтения реконструируя их из сегмента отмены (который также называется *сегментом отката* и детально обсуждается в главе 9). Согласованный и корректный ответ возвращается обратно без ожидания, пока транзакция будет зафиксирована.

Теперь база данных, допускающая грязное чтение, просто возвратит значение, которое обнаружит в строке для счета 987 в момент ее чтения — в данном случае это \$500.00. Запрос учтет переведенную сумму \$400 дважды. Таким образом, он не только выдаст неправильный ответ, но также вернет общую сумму, которая никогда не существовала в базе данных ни в один зафиксированный момент времени. В многопользовательской базе грязное чтение может быть опасной характеристикой, и лично я совершенно не видел в нем хоть какой-нибудь пользы. Пусть вместо перевода транзакция просто добавила \$400 на счет 987. Грязное чтение приняло бы во внимание \$400.00 и выдало бы “правильный” ответ, не так ли? Хорошо, предположим, что произошел откат незафиксированной транзакции. Мы получаем сумму \$400.00, которой в действительности никогда не было в базе данных.

Вывод из этого такой: грязное чтение — это не средство, а источник неприятностей. В Oracle оно просто не нужно. Вы получаете все преимущества грязного чтения (без блокирования) безо всяких некорректных результатов.

Уровень изоляции READ COMMITTED

Уровень изоляции READ COMMITTED устанавливает, что транзакция может читать только те данные, которые были зафиксированы в базе данных. Грязные чтения отсутствуют. Могут возникать невоспроизводимые чтения (т.е. повторные чтения той

же строки могут давать разные ответы в одной и той же транзакции) и фантомные чтения (т.е. вновь вставленные и зафиксированные строки становятся видимыми запросу, который не видел их ранее в транзакции). Вероятно, READ COMMITTED является наиболее часто используемым уровнем изоляции во всех приложениях баз данных, будучи стандартным режимом в базах данных Oracle. Другие уровни изоляции встречаются редко.

Однако обеспечение уровня изоляции READ COMMITTED не так просто, как может показаться. Если вы заглянете в табл. 7.1, то все выглядит простым. Очевидно, с учетом изложенных выше правил запрос, выполняемый в любой базе данных с применением уровня изоляции READ COMMITTED, будет вести себя одинаково? *Нет, не будет.* Если запросить несколько строк в одном операторе, то почти во всех других базах данных в зависимости от реализации изоляция READ COMMITTED будет настолько же плохой, как и грязное чтение.

В Oracle за счет использования многоверсионности и согласованных по чтению запросов ответ, получаемый из запроса таблицы ACCOUNTS, для уровня изоляции READ COMMITTED будет таким же, как в примере с READ UNCOMMITTED. СУБД Oracle воссоздает модифицированные данные в том виде, в каком они были на момент начала запроса, возвращая ответ, который содержался в базе данных, когда запрос только запускался.

Давайте теперь посмотрим, как предшествующий пример может работать в режиме READ COMMITTED в других базах данных. Ответ вас удивит. Обратимся к примеру в точке, описанной в предыдущей таблице.

- Мы находимся в середине таблицы. Мы прочитали и просуммировали первые *N* строк.
- Другая транзакция перевела \$400.00 со счета 123 на счет 987.
- Транзакция пока не зафиксирована, поэтому строки, содержащие информацию о счетах 123 и 987, блокированы.

Мы знаем, что происходит в СУБД Oracle, когда она добирается до счета 987: она обходит модифицированные данные, выясняет, что они должны выглядеть как \$100.00, и завершает обработку. В табл. 7.4 показано, какой ответ может выдать другая база данных, функционирующая в стандартном режиме READ COMMITTED.

Первое, что следует отметить — другая база данных, добравшись до счета 987, блокирует запрос. Сеанс должен ожидать, находясь на этой строке, пока не будет зафиксирована транзакция, удерживающая монопольную блокировку. Это одна из причин, почему многие имеют плохую привычку фиксировать каждый оператор вместо того, чтобы обрабатывать тщательно продуманные транзакции, которые состоят из всех операторов, необходимых для перевода базы данных из одного согласованного состояния в другое. *Обновления служат препятствием для чтений в большинстве других баз данных.* Но что действительно плохо в таком сценарии — то, что мы вынуждаем конечного пользователя ожидать *неправильного* ответа.

Мы получаем ответ, который не существовал в базе данных ни в какой момент времени, как было бы при грязном чтении, но на этот раз еще и заставляем пользователя ждать этого некорректного ответа. В следующем разделе мы посмотрим, что необходимо предпринять в других базах данных, чтобы получить согласованные по чтению, правильные результаты.

Таблица 7.4. Временная шкала в отличной от Oracle базе данных, использующей изоляцию READ COMMITTED

| Момент времени | Запрос | Транзакция перевода между счетами |
|----------------|---|--|
| T1 | Читает строку 1, счет 123, значение = \$500. К этому моменту сумма = \$500.00 | – |
| T2 | Читает строку 2, счет 456, значение = \$240.25. К этому моменту сумма = \$740.25 | – |
| T3 | – | Обновляет строку 1 (счет 123) и помещает на нее монопольную блокировку, предотвращая другие обновления и чтения. Строка 1 содержала \$500.00, теперь она содержит \$100.00 |
| T4 | Читает строку N. Сумма = ... | – |
| T5 | – | Обновляет строку 342 023 (счет 987) и помещает на нее монопольную блокировку. Строка 342 023 содержала \$100.00, теперь она содержит \$500.00 |
| T6 | Пытается прочитать строку 342 023, счет 987. Обнаруживает, что она заблокирована. Сеанс ожидает снятия блокировки с этой строки. <i>Вся обработка в этом запросе останавливается</i> | – |
| T7 | – | Фиксирует транзакцию |
| T8 | Читает строку 342 023, счет 987. Видит в ней \$500.00 и предоставляет окончательный ответ, в котором дважды учтена сумма \$400.00 | – |

Важный урок, который следует здесь извлечь, заключается в том, что разные базы данных, реализуя один и тот же очевидно безопасный уровень изоляции, в совершенно одинаковых обстоятельствах могут и будут возвращать значительно отличающиеся ответы. Важно понимать, что неблокирующие чтения в Oracle обеспечиваются не за счет корректности ответов. Иногда можно не только иметь мед, но и есть его ложкой.

Уровень изоляции REPEATABLE READ

Целью REPEATABLE READ является обеспечение такого уровня изоляции, который дает согласованные, корректные ответы и предотвращает потерянные обновления. Мы рассмотрим примеры того и другого, увидим, что нужно делать в Oracle для достижения этих целей, и исследуем происходящее в других системах.

Получение согласованного ответа

Если установлен уровень изоляции REPEATABLE READ, то результаты заданного запроса должны быть согласованными на определенный момент времени.

Большинство баз данных (не Oracle) добиваются воспроизводимых чтений через применение разделяемых блокировок чтения на уровне строки. Это, конечно же, снижает степень параллелизма. Для обеспечения согласованных по чтению ответов в Oracle была выбрана более параллельная, многоверсионная модель.

В Oracle с использованием многоверсионности мы получаем ответ, согласованный на момент времени, когда запрос начал свое выполнение. В других базах данных с применением разделяемых блокировок чтения мы получаем ответ на момент завершения запроса — т.е. тогда, когда вообще можно получить ответ (об этом речь пойдет далее в главе).

В системе, в которой для обеспечения воспроизводимых чтений задействована разделяемая блокировка чтения, мы можем наблюдать в таблице строки, заблокированные процессом, который их обрабатывает. Таким образом, в приведенном ранее примере при чтении таблицы ACCOUNTS запрос оставляет разделяемые блокировки чтения на каждой строке (табл. 7.5).

Таблица 7.5. Временная шкала 1 в отличной от Oracle базе данных, использующей изоляцию READ REPEATABLE

| Момент времени | Запрос | Транзакция перевода между счетами |
|----------------|--|---|
| T1 | Читает строку 1. К этому моменту сумма = \$500.00. Устанавливает разделяемую блокировку чтения на строку 1 | — |
| T2 | Читает строку 2. К этому моменту сумма = \$740.25. Устанавливает разделяемую блокировку чтения на строку 2 | — |
| T3 | — | Пытается обновить строку 1, но она заблокирована. Транзакция приостанавливается до тех пор, пока не сможет получить монополную блокировку |
| T4 | Читает строку N. Сумма = ... | — |
| T5 | Читает строку 342 023, видит \$100.00 и возвращает окончательный ответ | — |
| T6 | Фиксирует транзакцию | — |
| T7 | — | Обновляет строку 1 и помещает монополную блокировку на эту строку. В строке теперь \$100.00 |
| T8 | — | Обновляет строку 342 023 и помещает монополную блокировку на эту строку. В строке теперь \$500.00. Фиксирует транзакцию |

В табл. 7.5 легко заметить, что мы теперь получаем корректный ответ, но за счет физического блокирования одной транзакции и выполнения двух транзакций последовательно. Так проявляется один из побочных эффектов разделяемых блокировок чтения для согласованных ответов: *процессы чтения данных будут блокировать*

процессы записи данных. Это еще и в дополнение к тому факту, что в таких системах процессы записи данных будут блокировать процессы чтения. Только представьте, если бы банкоматы работали подобным образом в реальной жизни.

Итак, вы видите, что разделяемые блокировки чтения могут подавлять параллелизм, но они также могут вызывать появление фиктивных ошибок. В табл. 7.6 мы начинаем с исходной таблицы ACCOUNTS, но на этот раз цель заключается в переводе \$50.00 со счета 987 на 123.

Таблица 7.6. Временная шкала 2 в отличной от Oracle базе данных, использующей изоляцию READ REPEATABLE

| Момент времени | Запрос | Транзакция перевода между счетами |
|----------------|--|--|
| T1 | Читает строку 1. К этому моменту сумма = \$500.00. Устанавливает разделяемую блокировку чтения на строку 1 | — |
| T2 | Читает строку 2. К этому моменту сумма = \$740.25. Устанавливает разделяемую блокировку чтения на строку 2 | — |
| T3 | — | Обновляет строку 342 023 и помещает монопольную блокировку на строку 342 023, предотвращая другие обновления и установку разделяемых блокировок чтения. Строка теперь содержит \$50.00 |
| T4 | Читает строку N. Сумма = ... | — |
| T5 | — | Пытается обновить строку 1, но она заблокирована. Транзакция приостанавливается до тех пор, пока не сможет получить монопольную блокировку |
| T6 | Пытается прочитать строку 342 023, но не может, т.к. установлена монопольная блокировка | — |

Мы только что получили классическое условие взаимоблокировки. Наш запрос удерживает ресурсы, в которых нуждается обновление, и наоборот. Транзакция нашего запроса взаимно заблокирована транзакцией обновления. Одна из них должна быть выбрана в качестве жертвы и уничтожена. Мы потратили много времени и ресурсов лишь для того, чтобы в конце потерпеть неудачу и произвести откат. Это второй побочный эффект разделяемых блокировок чтения: *процессы чтения и процессы записи данных могут и часто будут попадать в ситуацию взаимоблокировки друг друга.*

В Oracle обеспечивается согласованность чтения на уровне операторов без блокирования операций записи операциями чтения или взаимоблокировок. В Oracle никогда не применяются разделяемые блокировки чтения — *вообще никогда*. Разработчики Oracle избрали трудную в реализации, но обеспечивающую намного более высокую степень параллелизма схему многоверсионности.

Потерянные обновления: еще одна проблема переносимости

Распространенное использование REPEATABLE READ в базах данных, эксплуатирующих разделяемые блокировки чтения, связано с предотвращением потерянных обновлений.

На заметку! Обнаружение потерянных обновлений и решения проблемы потерянных обновлений обсуждаются в главе 6.

При наличии уровня изоляции REPEATABLE READ в базе данных, которая *действует разделяемые блокировки чтения* (и не поддерживает многоверсионность), ошибки потерянных обновлений происходить не могут. Причина отсутствия потерянных обновлений в этих базах данных связана с тем, что простая выборка данных *оставляет на них блокировку, и данные, однажды прочитанные нашей транзакцией, не могут быть модифицированы никакой другой транзакцией*. Теперь, если в приложении предполагалось, что REPEATABLE READ означает “потерянные обновления не возникнут”, то вы столкнетесь с неприятным сюрпризом при переносе приложения в базу данных, в которой не применяются разделяемые блокировки чтения в качестве лежащего в основе механизма управления параллелизмом.

На заметку! В среде, не поддерживающей состояние, такой как веб-приложение, потерянные обновления, скорее всего, станут причиной проблем — даже при уровне изоляции REPEATABLE READ. Это объясняется тем, что единственный сеанс базы данных используется многими клиентами через пул подключений, а блокировки между вызовами не удерживаются. Уровень изоляции REPEATABLE READ предотвращает потерянные обновления только в среде с поддержкой состояния вроде той, которую можно наблюдать в клиент-серверном приложении.

Хотя это звучит неплохо, вы должны помнить, что оставление разделяемых блокировок чтения на всех данных после их чтения, разумеется, серьезно ограничит возможности параллельных операций чтения и модификации. Таким образом, хотя указанный уровень изоляции в этих базах данных предназначен для предотвращения потерянных обновлений, это достигается полным исключением возможности выполнения параллельных операций! Иметь мед и есть его ложкой можно не всегда.

Уровень изоляции SERIALIZABLE

Обычно SERIALIZABLE считается наиболее ограничивающим уровнем изоляции транзакций, но он предлагает максимальную степень изоляции. Транзакция SERIALIZABLE работает в среде, которая позволяет сделать вид, будто бы в базе данных нет других пользователей, изменяющих данные. Любая прочитанная строка гарантированно будет той же самой при ее повторном чтении, и любой выполняемый запрос гарантированно возвратит те же самые результаты на протяжении времени жизни транзакции. Например, если выполнить приведенные ниже запросы, то ответы, полученные из таблицы T, будут одинаковыми, хотя между ними прошло 24 часа (или может возникнуть ошибка ORA-01555: snapshot too old (ORA-01555: устаревший снимок), о которой речь пойдет в главе 8):

```
select * from T;
begin dbms_lock.sleep( 60*60*24 ); end;
select * from T;
```

Уровень изоляции `SERIALIZABLE` гарантирует, что эти два запроса будут всегда возвращать те же самые результаты. Побочные эффекты (изменения), сделанные другими транзакциями, для этого запроса остаются невидимыми независимо от того, сколько времени он выполнялся.

В Oracle транзакция `SERIALIZABLE` реализована так, что согласованность чтения, которую мы обычно получаем на уровне оператора, расширяется до уровня транзакции.

На заметку! Как отмечалось ранее, в Oracle также существует уровень изоляции, обозначенный как `READ ONLY`. Он обладает всеми характеристиками уровня изоляции `SERIALIZABLE`, но запрещает модификацию. Следует отметить, что пользователь `SYS` (или пользователь, подключенный с привилегией `SYSDBA`) не может иметь транзакцию `READ ONLY` или `SERIALIZABLE`. В этом отношении пользователь `SYS` является особенным.

Вместо обеспечения согласованности результатов на момент запуска оператора это делается для момента начала транзакции. Другими словами, Oracle применяет сегменты отката для воссоздания данных в том виде, в каком они существовали, когда начиналась транзакция, а не оператор.

Здесь скрыта довольно глубокая мысль: база данных уже знает ответ на любой вопрос, который вы можете задать, еще до того, как вы спросите.

Уровень изоляции `SERIALIZABLE` обходится не бесплатно, и ценой является возможность возникновения следующей ошибки:

```
ERROR at line 1:
ORA-08177: can't serialize access for this transaction
ОШИБКА в строке 1:
ORA-08177: не удастся сериализовать доступ для этой транзакции
```

Вы будете получать это сообщение при каждой попытке обновления строки, которая изменилась с момента начала транзакции.

На заметку! СУБД Oracle пытается делать это исключительно на уровне строки, но вы можете получить ошибку `ORA-08177`, даже когда строка, которую вы хотите модифицировать, не была изменена. Ошибка `ORA-08177` может возникнуть из-за того, что была модифицирована какая-то другая строка (или строки) из блока, в котором находится интересующая вас строка.

В Oracle принят оптимистический подход к сериализации — ставка делается на тот факт, что данные, которые ваша транзакция собирается обновить, не будут обновлены какой-то другой транзакцией. Обычно так и происходит, поэтому ставка выигрывает, особенно в системах типа OLTP с быстрыми транзакциями. Если никто другой не обновляет данные во время вашей транзакции, то этот уровень изоляции, который обычно будет снижать степень параллелизма в остальных системах, предоставит здесь ту же самую степень параллелизма, что и без транзакций `SERIALIZABLE`.

Недостаток заключается в том, что можно получить ошибку `ORA-08177`, если ставка оказалась проигрышной. Однако если подумать, то такой риск вполне оправдан. Если вы используете транзакции `SERIALIZABLE`, то не должны ожидать обновлений информации, с которой работают другие транзакции. В противном случае необходимо применять `SELECT...FOR UPDATE`, как было описано в главе 1, и это сериализирует доступ. Таким образом, использование уровня изоляции `SERIALIZABLE` достижимо и эффективно, если соблюдаются следующие условия.

- Высока вероятность того, что никто другой не модифицирует те же самые данные.
- Требуется согласованность чтения на уровне транзакции.
- Транзакции будут короткими (что поможет сделать реальным первое условие).

В Oracle считают этот метод достаточно масштабируемым для прогона всех их эталонных тестов TPC-C (эталонный тест OLTP, соответствующий промышленным стандартам; подробные сведения доступны по адресу <http://www.tpc.org>). Во многих других реализациях вы обнаружите, что те же самые цели достигаются посредством разделяемых блокировок чтения с присущими им взаимоблокировками и блокированием. В Oracle мы не столкнемся с каким-либо блокированием, но получим ошибку `ORA-08177`, если другие сеансы будут модифицировать данные, которые мы также хотим изменять. Однако эта ошибка возникает не так часто, как взаимоблокировки и блокирование в других системах.

Но — всегда есть какое-то “но” — вы должны четко понимать разницу между уровнями изоляции и их реализациями. Помните, что при установке уровня изоляции в `SERIALIZABLE` вы не увидите никаких изменений, внесенных в базу данных после запуска вашей транзакции — вплоть до ее фиксации. Приложения, которые пытаются поддерживать собственные ограничения целостности данных, такие как планировщик ресурсов, описанный в главе 1, в связи с этим должны предпринимать дополнительные меры предосторожности. Если помните, в главе 1 была описана проблема, которая заключалась в невозможности принудительного применения ограничения целостности в многопользовательской системе из-за того, что мы не могли видеть изменения, сделанные незафиксированными сеансами. В случае использования уровня `SERIALIZABLE` мы по-прежнему не увидим незафиксированных изменений, но также не увидим и зафиксированных изменений, произведенных после начала нашей транзакции!

Наконец, имейте в виду, что уровень изоляции `SERIALIZABLE` *вовсе не* означает, что все транзакции, запущенные пользователями, будут вести себя так, как будто они выполняются друг за другом в последовательной манере. Это не подразумевает наличие какого-то последовательного упорядочения транзакций, что приведет к тому же самому результату. Феномены, ранее описанные в стандарте SQL, не позволяют такому случиться. Последнее утверждение отражает часто неправильно понимаемую концепцию, которую поможет прояснить небольшая демонстрация. В табл. 7.7 представлены два сеанса, выполняющие работу на протяжении некоторого времени. Изначально пустые таблицы базы данных A и B создаются следующим образом:

```

EODA@ORA12CR1> create table a ( x int );
Table created.
Таблица создана.

```



```
EODA@ORA12CR1> create table b ( x int );
Table created.
Таблица создана.
```

Мы имеем последовательность событий, описанную в табл. 7.7.

Таблица 7.7. Пример транзакции **SERIALIZABLE**

| Момент времени | Действия в сеансе 1 | Действия в сеансе 2 |
|----------------|---|---|
| T1 | alter session set isolation_level=serializable; | - |
| T2 | - | alter session set isolation_level=serializable; |
| T3 | insert into a select count(*) from b; | - |
| T4 | - | insert into b select count(*) from a; |
| T5 | commit; | - |
| T6 | - | commit; |

Теперь, когда все показанное сделано, таблицы А и В будут содержать строку со значением 0 каждая. Если бы существовало некоторое последовательное упорядочение транзакций, присутствие значения 0 в обеих таблицах было бы невозможным. Если бы сеанс 1 выполнялся *во всей своей полноте* перед сеансом 2, то таблица В имела бы строку со значением 1 в ней. Если бы сеанс 2 выполнялся *полностью* перед сеансом 1, то таблица А содержала бы строку со значением 1. Однако на самом деле мы получили строки со значением 0 в обеих таблицах. Каждая транзакция выполнялась так, как будто в конкретный момент времени она была единственной в базе данных. Независимо от того, сколько раз сеанс 1 выдавал запросы к таблице В, и независимо от состояния фиксации сеанса 2, счетчик получит значение, которое было зафиксировано в базе данных на момент времени T1. Аналогично, не играет роли, сколько запросов к таблице А производил сеанс 2 — счетчик имеет значение, которое было в момент T2.

Уровень изоляции **READ ONLY**

Транзакции READ ONLY очень похожи на транзакции SERIALIZABLE с единственным отличием: они не допускают модификаций, поэтому не восприимчивы к ошибке ORA-08177. Транзакции READ ONLY предназначены для поддержки потребностей, возникающих при формировании отчетов, когда содержимое отчета должно быть согласовано с состоянием на определенный момент времени. В других системах вам придется применять транзакцию REPEATABLE READ и страдать от последствий, вызванных разделяемой блокировкой чтения. В Oracle вы будете использовать транзакцию READ ONLY. В этом режиме вывод, генерируемый в отчете, который применяет 50 операторов SELECT для сбора данных, будет согласованным по состоянию на единственный момент времени — момент начала транзакции. Вы получите возможность делать это, не блокируя ни единой порции данных.

Эта цель достигается использованием той же самой многоверсионности, которая применялась для индивидуальных операторов. Данные реконструируются из сегментов отката по мере необходимости и предоставляются в том виде, в каком они пребывали на начало формирования отчета. Тем не менее, транзакции READ ONLY не лишены недостатков. В то время как в транзакциях SERIALIZABLE можно столкнуться с ошибкой ORA-08177, в транзакциях READ ONLY следует ожидать ошибки ORA-01555: snapshot too old (ORA-01555: устаревший снимок). Это будет происходить в системах, где другие пользователи активно модифицируют информацию, которую вы читаете. Изменения, вносимые в эту информацию, сохраняются в сегментах отката. Но сегменты отката используются в циклическом режиме — почти так же, как журнальные файлы повторения действий.

Чем дольше требуется времени на формирование отчета, тем выше вероятность, что какая-то операция отката, необходимая для воссоздания ваших данных, уже не сможет найти необходимых сведений в сегменте отката. Сегмент отката начинает перезаписываться другой транзакцией, и нужная вам часть информации утрачивается. В этот момент вы получаете ошибку ORA-01555 и вынуждены начинать заново.

Единственное решение этой неприятной проблемы предусматривает корректную установку размера табличного пространства отката для системы. Я постоянно сталкиваюсь с людьми, которые пытаются сэкономить несколько мегабайт дискового пространства, создавая табличное пространство отката с минимально возможным размером (“Зачем тратить пространство на то, в чем я действительно не нуждаюсь?” — примерно так они думают).

Проблема в том, что табличное пространство отката является ключевым компонентом при работе базы данных, и если его размер определен неправильно, вы обязательно столкнетесь с упомянутой выше ошибкой. В течение многих лет работы с версиями Oracle 6, 7, 8, 9, 10, 11 и 12 я никогда не сталкивался с ошибкой ORA-01555 за пределами тестовой системы или системы разработки. В случае возникновения такой ошибки вы должны знать, что некорректно задали размер табличного пространства отката, и это необходимо исправить. Мы вернемся к этой теме в главе 9.

Последствия многоверсионной согласованности чтения

До настоящего времени мы видели, каким образом многоверсионность обеспечивает возможность неблокирующего чтения, и я подчеркивал, что это хорошая вещь: согласованные (корректные) ответы при высокой степени параллелизма. Что же в этом может быть плохого? Если вы четко не понимаете весь механизм и его последствия, то вполне вероятно, что спроектируете некоторые из своих транзакций неправильно.

Вспомните пример с планированием ресурсов из главы 1, в котором мы пытались задействовать ряд приемов ручной блокировки (через SELECT FOR UPDATE для сериализации модификаций записей о ресурсах в таблице SCHEDULES). Но может ли это затронуть нас другими путями? Однозначно может. В последующих разделах приведены соответствующие детали.

Распространенный прием организации хранилищ данных, который не работает

Распространенный прием организации хранилищ данных, применяемый многими людьми, выглядит следующим образом.

1. Они используют триггер для поддержки столбца `LAST_UPDATED` в исходной таблице, что очень похоже на метод, описанный в разделе “Оптимистическая блокировка” в предыдущей главе.
2. Для начального наполнения таблицы хранилища они запоминают текущее время, выбирая `SYSDATE` в исходной системе. Например, предположим, что прямо сейчас ровно 9:00.
3. Затем они извлекают все строки из транзакционной системы, т.е. выполняют полный запрос `SELECT * FROM <ТАБЛИЦА>`, чтобы провести начальное заполнение хранилища данных.
4. Для обновления хранилища данных они опять запоминают текущее время. Например, предположим, что прошел час и теперь в исходной системе 10:00. Они запоминают этот факт. Затем они извлекают все изменения, накопившиеся после 9:00 (момента перед началом первого извлечения), и объединяют их.

На заметку! При таком подходе одна и та же запись может быть извлечена дважды в двух последовательных обновлениях. Это неизбежно из-за степени детализации часов. На операцию `MERGE` данный факт (т.е. обновление существующей записи в хранилище или вставка новой записи) не повлияет.

Эти люди верят в то, что после описанных выше действий имеют в хранилище данных все записи, которые были модифицированы после первоначального их извлечения. Они действительно могут иметь все записи, но точно так же могут и не иметь. Такой прием работает с рядом других баз данных — теми, что применяют систему блокирования, посредством которой операции чтения блокируются операциями записи и наоборот. Но в системе с неблокирующими чтениями подобная логика ущербна.

Чтобы увидеть слабую сторону в этом примере, вполне достаточно предположить, что в 9:00 оставалась, по крайней мере, одна открытая, незафиксированная транзакция. В 8:59:30 она обновила в таблице строку, которую мы скопировали. В 9:00, когда мы начинаем извлекать данные, читая их из этой таблицы, мы не увидим модификации этой строки, а только ее последнюю зафиксированную версию. Если она заблокирована в момент ее обработки запросом, то мы прочитаем ее в обход блокировки. Если она зафиксирована ко времени, когда мы ее достигли, мы по-прежнему будем читать ее в обход блокировки, поскольку согласованность чтения позволяет производить чтение только данных, которые были зафиксированы в базе к моменту начала нашего оператора. Мы *не* прочитаем ни новую версию строки во время начальной загрузки в 9:00, ни модифицированную версию во время обновления хранилища в 10:00. В чем причина? При обновлении в 10:00 извлекаются только записи, модифицированные после 9:00, но эта запись была модифицирована в 8:59:30. Мы никогда не получим эту измененную строку!

Во многих других базах данных, где операции чтения блокируются операциями записи и реализованы зафиксированные, но несогласованные чтения, этот процесс обновления будет работать прекрасно. Если в 9:00 при выполнении начального извлечения данных мы столкнулись с этой строкой, и она оказалась заблокированной, мы будем ожидать ее освобождения и прочитаем зафиксированную версию. Если строка не заблокирована, то мы просто читаем то, что есть в строке и зафиксировано.

Итак, означает ли это, что описанная выше логика не может использоваться? Нет, это значит, что мы должны получать “текущее” время слегка по-другому. Необходимо запросить представление `V$TRANSACTION` и выяснить, что было раньше — текущий момент времени или время, записанное в столбце `START_TIME` этого представления. Понадобится извлечь все записи, измененные после времени начала самой старой транзакции (или текущего значения `SYSDATE`, если активных транзакций нет):

```
select nvl( min(to_date(start_time,'mm/dd/rr hh24:mi:ss')),sysdate)
from v$transaction;
```

На заметку! Приведенный выше запрос работает независимо от наличия каких-либо данных в `V$TRANSACTION`. То есть, даже если представление `V$TRANSACTION` оказывается пустым (из-за того, что в настоящий момент транзакции отсутствуют), этот запрос возвратит запись. Запрос с агрегированием без конструкций `WHERE` всегда возвращает, *по крайней мере, одну строку и не более одной строки.*

В рассматриваемом примере этим моментом будет 8:59:30, т.е. запуск транзакции, модифицирующей эту строку. Когда мы переходим к обновлению данных в 10:00, то извлекаем все изменения, которые произошли после того момента; в результате слияния их с хранилищем данных мы получим все необходимые записи.

Объяснение неожиданно высокой активности ввода-вывода в “горячих” таблицах

Еще одна ситуация, где жизненно важно понимать согласованность чтения и многоверсионность, связана с запросом, который в производственной среде с высокой нагрузкой инициирует намного больше операций ввода-вывода, чем вы наблюдали в тестовой системе или системе разработки, и вы не имеете никакой возможности объяснить это. Вы просматриваете сведения по вводу-выводу, выполненному запросом, и замечаете, что он значительно интенсивнее, чем вы когда-либо видели — намного выше, чем казалось возможным. Вы восстанавливаете производственный экземпляр базы данных в тестовой системе и обнаруживаете, что ввод-вывод снизился. Однако в производственной среде он по-прежнему очень высок (но варьируется: иногда выше, иногда ниже, а иногда где-то посередине). Причина, как мы увидим, заключается в том, что в тестовой системе, в изоляции, вам не приходится отменять изменения, вносимые другими транзакциями. Тем не менее, в производственной среде при чтении определенного блока может возникать необходимость в отмене (откате) изменений, произведенных многими транзакциями, и каждый откат предусматривает выполнение операций ввода-вывода для извлечения данных отмены и их применения.

Вероятно, это может быть запрос к таблице, которая подвержена многочисленным параллельным модификациям. Вы можете видеть это по факту чтений сегмента отката — работе, которую Oracle выполняет для восстановления блока в состоянии, в котором он пребывал, когда ваш запрос начинался. Эти эффекты легче всего увидеть в единственном сеансе, просто чтобы понять, что происходит. Начнем с очень маленькой таблицы:

```
EODA@ORA12CR1> create table t ( x int );
Table created.
Таблица создана.

EODA@ORA12CR1> insert into t values ( 1 );
1 row created.
1 строка создана.

EODA@ORA12CR1> exec dbms_stats.gather_table_stats( user, 'T' );
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.

EODA@ORA12CR1> select * from t;

      X
-----
      1
```

Теперь настроим сеанс на использование уровня изоляции `SERIALIZABLE`, так что независимо от того, сколько раз мы запускаем запрос в сеансе, результат будет таким, как на момент старта транзакции:

```
EODA@ORA12CR1> alter session set isolation_level=serializable;
Session altered.
Сеанс изменен.
```

Выполним запрос к маленькой таблице и понаблюдаем за объемом ввода-вывода:

```
EODA@ORA12CR1> set autotrace on statistics
EODA@ORA12CR1> select * from t;

      X
-----
      1

Statistics
-----
          0 recursive calls
          0 db block gets
          7 consistent gets
...

```

Итак, этому запросу для завершения потребовалось семь операций ввода-вывода. Далее мы многократно модифицируем эту таблицу в *другом сеансе*:

```
EODA@ORA12CR1> begin
2  for i in 1 .. 10000
3  loop
4  update t set x = x+1;
5  commit;
6  end loop;
7  end;
8  /
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.
```

Возвратившись в сеанс `SERIALIZABLE`, запустим тот же самый запрос еще раз:

```
EODA@ORA12CR1> select * from t;
```

```
      X
-----
      1
```

```
Statistics
```

```
-----
          0 recursive calls
          0 db block gets
    10004 consistent gets
...
```

На этот раз понадобилось 10 004 операций ввода-вывода — разница заметна. Откуда взялись все эти операции? Дело в том, что Oracle выполняет откат изменений, внесенных в этот блок данных. Когда мы запускаем второй запрос, Oracle известно, что все блоки, извлеченные и обработанные запросом, должны быть такими, какими они были в момент запуска транзакции. Когда мы обращаемся к буферному кешу, то обнаруживаем, что блок в кеше просто “слишком новый” — другой сеанс модифицировал его 10 000 раз. Наш запрос не видит изменений, поэтому он начинает просматривать информацию отката и отменять последнее изменение. Затем он выясняет, что этот блок после отката по-прежнему является слишком новым, и производит еще один откат блока. Процедура повторяется до тех пор, пока, наконец, не будет найдена версия блока, которая была зафиксирована в базе в начале транзакции. Этот блок можно использовать — и он используется.

На заметку! Интересно отметить, что если вы повторно запустите `SELECT * FROM T`, то количество операций ввода-вывода, скорее всего, снизится до 7 или около того; оно уже не будет составлять 10 004. В чем причина? СУБД Oracle обладает способностью хранить в буферном кеше несколько версий одного и того же блока. Когда вы отменили изменения в этом блоке для запроса, который выполнял 10 004 операции ввода-вывода, то оставили в кеше данную версию, и при последующих запусках запроса к ней будет производиться доступ.

Сталкиваемся ли мы с этой проблемой только в случае применения уровня изоляции `SERIALIZABLE`? Вовсе нет. Рассмотрим запрос, который выполняется в течение пяти минут. За это время он извлекает блоки из буферного кеша. При каждом извлечении блока запрос выясняет следующее: “Этот блок слишком новый? Если да — откатить его!”. И помните: чем дольше работает запрос, тем выше шансы, что нужный ему блок со временем окажется модифицированным.

Теперь база данных ожидает, что такая проверка может произойти (т.е. проверка, не является ли блок “слишком новым”, с последующим откатом изменений), и по этой причине буферный кеш может действительно содержать множество версий одного и того же блока в памяти. Таким образом, есть шанс, что требуемая версия будет находиться в буферном кеше, полностью готовая к использованию, поэтому ее не придется материализовывать с применением информации отката. Для просмотра этих блоков можно использовать запрос, подобный показанному ниже:

```
select file#, block#, count(*)
from v$dbh
group by file#, block#
having count(*) > 3
order by 3
/
```

Вообще говоря, в любой момент времени вы найдете не более примерно шести версий блока в кеше, но эти версии могут применяться любым запросом, который в них нуждается.

Обычно с такими маленькими *горячими* таблицами часто возникает проблема завышенного количества операций ввода-вывода из-за обеспечения согласованности чтения. Другими запросами, наиболее часто подверженными такой проблеме, являются длительно работающие запросы в отношении изменчивых таблиц. Чем дольше они функционируют, тем больше они работают, потому что со временем им может понадобиться выполнять все больше действий по извлечению блока из буферного кеша.

Согласованность записи

До сих пор мы рассматривали согласованность чтения: способность Oracle использовать информацию отката для обеспечения неблокирующих запросов и согласованных (корректных) чтений. Мы разобрались с тем, что когда СУБД Oracle читает блоки для запросов из буферного кеша, она гарантирует, что версия блока достаточно “старая”, чтобы быть видимой данному запросу.

Но это подводит к следующим вопросам. Как насчет записи/модификации? Что произойдет в ситуации, если вы запустили приведенный ниже оператор UPDATE, но во время его выполнения кто-то обновил строку, которую этот оператор должен прочесть, изменив Y=5 на Y=6 и тут же зафиксировав изменение?

```
update t set x = 2 where y = 5;
```

То есть, когда оператор UPDATE начинается, некоторая строка имеет значение Y=5. Поскольку оператор UPDATE читает таблицу с применением согласованных чтений, он видит, что эта строка имела Y=5 на момент начала его работы. Но текущее значение Y теперь равно 6, а не 5, и перед обновлением X база данных Oracle проверит, что Y по-прежнему имеет значение 5. Что случится? Как это повлияет на обновления?

Очевидно, мы не можем модифицировать старую версию блока; когда мы собираемся изменять строку, то должны модифицировать текущую версию блока. Вдобавок Oracle не может просто пропустить данную строку, т.к. это было бы несогласованным чтением и непредсказуемой ситуацией. Как мы вскоре убедимся, в таких случаях Oracle будет перезапускать модификацию посредством записи с нуля.

Согласованные чтения и текущие чтения

При обработке модифицирующего оператора Oracle использует два способа получения блока:

- согласованные чтения — при “поиске” строк для модификации;
- текущие чтения — при получении блока для действительного обновления интересующей строки.

Мы можем легко увидеть это с помощью утилиты TKPROF. Рассмотрим следующий небольшой однострочный пример, в котором производится чтение и обновление единственной строки в ранее представленной таблице T:

```

EODA@ORA12CR1> exec dbms_monitor.session_trace_enable
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.

```

```

EODA@ORA12CR1> select * from t;

```

```

      X
-----
10001

```

```

EODA@ORA12CR1> update t t1 set x = x+1;
1 row updated.
1 строка обновлена.

```

```

EODA@ORA12CR1> update t t2 set x = x+1;
1 row updated.
1 строка обновлена.

```

Запустив TKPROF, мы увидим примерно такие результаты (обратите внимание, что из этого отчета исключены столбцы ELAPSED, CPU и DISK):

```

select * from t

```

| call | count | query | current | rows |
|---------|-------|-------|---------|------|
| Parse | 1 | 0 | 0 | 0 |
| Execute | 1 | 0 | 0 | 0 |
| Fetch | 2 | 7 | 0 | 1 |
| total | 4 | 7 | 0 | 1 |

```

update t t1 set x = x+1

```

| call | count | query | current | rows |
|---------|-------|-------|---------|------|
| Parse | 1 | 0 | 0 | 0 |
| Execute | 1 | 7 | 3 | 1 |
| Fetch | 0 | 0 | 0 | 0 |
| total | 2 | 7 | 3 | 1 |

```

update t t2 set x = x+1

```

| call | count | query | current | rows |
|---------|-------|-------|---------|------|
| Parse | 1 | 0 | 0 | 0 |
| Execute | 1 | 7 | 1 | 1 |
| Fetch | 0 | 0 | 0 | 0 |
| total | 2 | 7 | 1 | 1 |

Итак, во время вполне нормального запроса мы имеем семь *получений в режиме (согласованного) запроса*. Для первого оператора UPDATE мы имеем те же семь операций ввода-вывода (поисковый компонент операции обновления в данном случае предусматривает поиск всех строк, находящихся в таблице на момент начала обновления) и три *получения в текущем режиме*. Получения в текущем режиме (current mode gets) выполняются для того, чтобы извлечь *табличный блок* в том виде, в каком он существует прямо сейчас, вместе с нужной строкой в нем. Это необходимо для получения *блока сегмента отмены*, чтобы начать транзакцию, и *блока отмены*.

Второй оператор UPDATE имеет в точности одно получение в текущем режиме; так как выполнять снова работу по откату не требуется, есть только одно текущее получение для блока со строкой, подлежащей обновлению. Именно присутствие получений в текущем режиме говорит о том, что имела место модификация определенного рода. Прежде чем СУБД Oracle модифицирует блок с применением новой информации, она должна получить наиболее свежую его копию.

Каким же образом согласованность чтения влияет на модификацию? Представьте, что вы запускаете следующий оператор UPDATE в отношении какой-то таблицы базы данных:

```
update t set x = x+1 where y = 5;
```

Мы понимаем, что компонент WHERE Y=5 (согласованная по чтению фаза запроса) будет обработан с использованием согласованного чтения (получения в режиме запроса в отчете TKPROF). Набор записей WHERE Y=5, зафиксированный в таблице в начале выполнения оператора — это те записи, которые он будет видеть (предполагается применение уровня изоляции READ COMMITTED; в случае уровня SERIALIZABLE это был бы набор записей WHERE Y=5, которые существовали, когда начиналась транзакция). Это означает, что если бы оператору UPDATE понадобилось пять минут на обработку от начала до конца, и кто-то добавит в таблицу и зафиксирует новую запись со значением 5 в столбце Y, то оператор UPDATE ее не увидит, поскольку ее не должно видеть согласованное чтение. Поведение является ожидаемым и нормальным. Но вопрос в том, что произойдет, если два сеанса выполняют следующие операторы по порядку:

```
update t set y = 10 where y = 5;  
update t set x = x+1 where y = 5;
```

Временная шкала приведена в табл. 7.8.

Таблица 7.8. Последовательность обновлений

| Момент времени | Сеанс 1 | Сеанс 2 | Комментарии |
|----------------|---------------------------------|----------------------------------|--|
| T1 | update t set y=10 where y=5; | | Обновит одну строку, которая соответствует критерию |
| T2 | | update t set x=x+1 where y=5; | Используя согласованные чтения, найдет запись, модифицированную в сеансе 1, но не сможет обновить ее, т.к. сеанс 1 ее заблокировал. Сеанс 2 будет заблокирован и ожидать освобождения этой строки |
| T3 | commit; | | Освобождает сеанс 2: сеанс 2 становится разблокированным. Наконец, можно выполнить текущее чтение блока, содержащего эту строку, где Y было равно 5, когда сеанс 1 начинал обновление. Текущее чтение покажет, что Y теперь равно 10, а не 5 |

Таким образом, запись, в которой значение Y было равно 5 на момент начала оператора UPDATE, больше не содержит $Y=5$. Компонент согласованного чтения оператора UPDATE говорит: “вы хотите обновить эту запись, потому что Y было равно 5, когда мы начинали”, но текущая версия блока заставляет подумать: “нет, я не могу обновить эту строку, поскольку Y больше не равно 5 — это было бы неправильно”.

Если в этой точке мы просто пропустим данную строку и проигнорируем ее, то получим недетерминированное обновление. Это может нарушить согласованность и целостность данных исчезнет. Исход обновления (сколько строк было обновлено и какие именно) будет зависеть от порядка, в котором строки посещались в таблице, и какие другие действия происходили в ходе работ. Вы можете взять один и тот же набор строк в двух разных базах данных, в которых выполняются транзакции в том же самом порядке, и получить отличающиеся результаты лишь потому, что строки были по-разному расположены на диске.

В этом случае Oracle решит *перезапустить* обновление. Обнаружив, что строка, которая в начале имела $Y=5$, теперь содержит $Y=10$, СУБД Oracle молча выполнит откат обновления (только этого обновления, но никакой другой части транзакции) и перезапустит его заново, предполагая использование уровня изоляции READ COMMITTED. В случае применения уровня изоляции SERIALIZABLE в этот момент возникнет ошибка ORA-08177: can't serialize access for this transaction (ORA-08177: не удастся сериализовать доступ для этой транзакции). В режиме READ COMMITTED после отката транзакцией обновления база данных перезапустит обновление (т.е. изменит точку во времени, которой соответствует обновление), а затем вместо повторного обновления данных войдет в режим SELECT FOR UPDATE и попытается заблокировать для сеанса все строки, имеющие $Y=5$. Сделав это, база данных запустит оператор UPDATE на заблокированном наборе данных, тем самым гарантируя, что на этот раз она сможет завершить операцию без перезапуска.

Продолжая думать в стиле “что случится, если...”, возникает вопрос: что произойдет, если после перезапуска обновления и перехода в режим SELECT FOR UPDATE (который связан с теми же получениями блока операциями согласованного чтения и текущего чтения при выполнении обновления) обнаружится, что строка, которая имела $Y=5$ при запуске SELECT FOR UPDATE, содержит $Y=11$ в своей текущей версии? Этот оператор SELECT FOR UPDATE будет перезапущен и цикл начнет-ся сначала.

Здесь возникает множество интересных вопросов. Можем ли мы это наблюдать? Можем ли мы выяснить, что это действительно произошло? И если да, что из этого следует? Что это означает для разработчиков? Давайте ответим на эти вопросы по очереди.

Наблюдение за перезапуском

Перезапуск увидеть проще, чем может сначала казаться. Фактически его можно наблюдать, используя простую однострочную таблицу. Вот таблица, которая будет применяться при тестировании:

```
EODA@ORA12CR1> create table t ( x int, y int );
Table created.
Таблица создана.
```

```

EODA@ORA12CR1> insert into t values ( 1, 1 );
1 row created.
1 строка создана.

EODA@ORA12CR1> commit;
Commit complete.
Фиксация завершена.

```

Для наблюдения за перезапуском понадобится всего лишь триггер, выводящий некоторую информацию. Мы будем использовать триггер BEFORE UPDATE FOR EACH ROW, чтобы выводить образ строки до и после обновления:

```

EODA@ORA12CR1> create or replace trigger t_bufer
2 before update on t for each row
3 begin
4     dbms_output.put_line
5     ( 'old.x = ' || :old.x ||
6       ', old.y = ' || :old.y );
7     dbms_output.put_line
8     ( 'new.x = ' || :new.x ||
9       ', new.y = ' || :new.y );
10 end;
11 /
Trigger created.
Триггер создан.

```

Теперь обновим строку:

```

EODA@ORA12CR1> set serveroutput on
EODA@ORA12CR1> update t set x = x+1;
old.x = 1, old.y = 1
new.x = 2, new.y = 1
1 row updated.
1 строка обновлена.

```

До сих пор все идет вполне ожидаемо: триггер запустился один раз, и мы видим старые и новые значения. Однако обратите внимание, что мы пока еще не произвели фиксацию — строка по-прежнему заблокирована. Выполним в другом сеансе следующее обновление:

```

EODA@ORA12CR1> set serveroutput on
EODA@ORA12CR1> update t set x = x+1 where x > 0;

```

Конечно, второй сеанс будет немедленно заблокирован, т.к. первый сеанс имеет блокировку на этой строке. Если теперь мы вернемся в первый сеанс и зафиксируем транзакцию, то во втором сеансе увидим такой вывод (для ясности оператор обновления повторен):

```

EODA@ORA12CR1> update t set x = x+1 where x > 0;
old.x = 1, old.y = 1
new.x = 2, new.y = 1
old.x = 2, old.y = 1
new.x = 3, new.y = 1
1 row updated.
1 строка обновлена.

```

Как видите, триггер строки увидел здесь две версии строки. Триггер был активизирован дважды: первый раз с исходной версией строки, а также данными, на которые мы пытались изменить эту версию, и второй раз с финальной строкой, которая действительно была обновлена. Поскольку это был триггер BEFORE FOR EACH ROW, база данных Oracle видит согласованную по чтению версию строки и модификации, которые мы хотели в ней провести. Однако Oracle извлекает блок в текущем режиме для выполнения обновления *после* срабатывания триггера BEFORE FOR EACH ROW. База данных ожидает окончания работы триггера, чтобы получить блок в текущем режиме, т.к. триггер может модифицировать значения :NEW. Поэтому Oracle не может изменять блок до тех пор, пока триггер не завершит работу, а это может занять довольно длительное время. Поскольку удерживать блокировку в текущем режиме может только один сеанс за раз, Oracle приходится ограничивать время пребывания в этом режиме.

После запуска триггера Oracle извлекает блок в текущем режиме и отмечает, что столбец, применявшийся для поиска этой строки (X), был модифицирован. Так как для нахождения этой записи использовался столбец X, и он был изменен, база данных решает перезапустить наш запрос. Обратите внимание, что изменение значения X с 1 на 2 не выносит эту строку за рамки области действия; мы по-прежнему будем обновлять ее с помощью этого оператора UPDATE. Взамен мы имеем ситуацию, при которой для поиска строки применялся столбец X, и полученное посредством согласованного чтения значение этого столбца (1) отличается от значения, прочитанного в текущем режиме (2). Теперь, после перезапуска, триггер видит значение X=2 (вследствие модификации в другом сеансе) как :OLD и значение X=3 — как :NEW.

Итак, мы убедились, что случаи перезапуска были. Благодаря триггеру мы увидели их в действии; по-другому их обычно *обнаружить невозможно*. Это не значит, что вы не заметите другие симптомы, такие как крупный оператор UPDATE, выполняющий откат работы после обновления многих строк и затем находящий строку, которая стала причиной его перезапуска; просто трудно сказать определенно, что какой-то конкретный симптом обусловлен именно перезапуском.

Интересным наблюдением является то, что триггеры сами могут стать причиной перезапуска, даже если собственно оператор не гарантирует этого. Обычно столбцы, упоминаемые в конструкции WHERE оператора UPDATE или DELETE, используются для определения того, нуждается ли модификация в перезапуске. СУБД Oracle выполнит согласованное чтение с применением таких столбцов и при извлечении блока в текущем режиме перезапустит оператор, если обнаружит, что любой из этих столбцов изменился. Как правило, другие столбцы в строке не инспектируются. Например, давайте просто перезапустим предыдущий пример и воспользуемся конструкцией WHERE Y>0, чтобы найти строки в обоих сеансах. В первом сеансе (блокированном) мы увидим следующий вывод:

```
EODA@ORA12CR1> update t set x = x+1 where y > 0;
old.x = 1, old.y = 1
new.x = 2, new.y = 1
old.x = 2, old.y = 1
new.x = 3, new.y = 1
1 row updated.
1 строка обновлена.
```

Так почему же база данных Oracle дважды запускала триггер, когда искала значение Y? Проверяла ли она целиком всю строку? Как видно в выводе, обновление действительно было перезапущено и триггер снова сработал дважды, хотя мы выполняли поиск по условию $Y > 0$ и вообще не модифицировали столбец Y. Но если переписать триггер так, чтобы он просто сообщал о факте своего запуска, не ссылаясь на значения :OLD и :NEW, а затем снова войти во второй сеанс и запустить обновление, то мы будем наблюдать, что оно заблокировано (естественно):

```

EODA@ORA12CR1> create or replace trigger t_bufer
2 before update on t for each row
3 begin
4     dbms_output.put_line( 'fired' );
5 end;
6 /

```

Trigger created.

Триггер создан.

```

EODA@ORA12CR1> update t set x = x+1;
fired
1 row updated.
1 строка обновлена.

```

После фиксации блокирующего сеанса мы увидим такой результат:

```

EODA@ORA12CR1> update t set x = x+1 where y > 0;
fired
1 row updated.
1 строка обновлена.

```

На этот раз триггер сработал только один раз, а не два. Таким образом, значения столбцов :OLD и :NEW, когда они присутствуют в триггере, также задействованы Oracle при проверке необходимости перезапуска. Когда мы ссылаемся в триггере на :OLD.X и :NEW.X, значения согласованного чтения и текущего чтения X сравниваются и оказываются разными. Перезапуск произойдет. Когда мы удалили ссылку на этот столбец из триггера, перезапуска не будет.

Итак, правило заключается в том, что набор столбцов в конструкции WHERE, используемых для нахождения строк, плюс столбцы, участвующие в триггерах строк, будут сравниваться. Версия строки, полученная при согласованном чтении, будет сравниваться с версией, полученной с помощью текущего чтения, и если они отличаются, то оператор модификации перезапустится.

На заметку! Эта информация способствует дальнейшему пониманию причин, по которым применение триггера AFTER FOR EACH ROW более эффективно, чем триггера BEFORE FOR EACH ROW. Триггер AFTER не обеспечит тот же самый результат — к этому моменту блок уже извлечен в текущем режиме.

Все это приводит к вопросу: почему нас это заботит?

Почему перезапуск важен для нас?

Первая реакция должна быть такой: наш триггер сработал дважды! У нас есть таблица с единственной строкой и триггером BEFORE FOR EACH ROW. Мы обновили одну строку, а триггер вызывался два раза.

Подумайте о потенциальных последствиях этого факта. Если триггер делает что-нибудь нетранзакционное, то это может стать довольно серьезной проблемой. Например, рассмотрим триггер, который отправляет обновление, когда в теле сообщения электронной почты присутствует фраза “This is what the data used to look like. It has been modified to look like this now.” (“Так выглядят используемые данные. Они были модифицированы, чтобы теперь выглядеть подобным образом.”). Если вы отправляете сообщение электронной почты прямо из триггера с помощью пакета UTL_SMTP в Oracle9i или пакета UTL_MAIL в Oracle 10g и последующих версиях, то пользователь получит два сообщения, одно из которых уведомит об обновлении, которое на самом деле никогда не происходило.

Перезапуск будет оказывать влияние на все, что делается в триггере и не является транзакционным. Рассмотрим перечисленные ниже последствия.

- Представим триггер, который поддерживает ряд глобальных переменных PL/SQL, таких как количество обработанных строк. Когда производится откат перезапускающегося оператора, модификации в переменных PL/SQL откату не подвергаются.
- Почти любую функцию, которая начинается с UTL_ (UTL_FILE, UTL_HTTP, UTL_SMTP и т.д.) следует считать восприимчивой к перезапуску оператора. При перезапуске оператора функции UTL_FILE не смогут отменить то, что они записали в файл.
- Любой триггер, являющийся частью автономной транзакции, должен ставить-ся под сомнение. Когда оператор перезапускается и происходит откат, автономные транзакции не могут быть отменены.

Все эти последствия должны обрабатываться осторожно с надеждой на то, что триггеры могут быть инициированы более одного раза на строку или запущены для строки, которая вообще не будет обновлена оператором.

Вторая причина, по которой вы должны заботиться о потенциальном перезапуске, связана с производительностью. Мы продемонстрировали однострочный пример, но что случится, если вы начнете крупное пакетное обновление, и оно перезапустится после обработки первых 100 000 записей? Ему придется выполнить откат 100 000 изменений строк, перезапуститься в режиме SELECT FOR UPDATE и затем снова обновить эти 100 000 строк.

Вы можете заметить, что после ввода в действие такого простого триггера аудита (который читает значения :OLD и :NEW) производительность стала намного хуже, чем можно было бы объяснить, хотя кроме добавления триггера ничего не изменилось. Это может происходить из-за перезапуска запросов, которые никогда ранее не использовались. Или же добавление крошечной программы, обновляющей единственную строку в нескольких местах, приводит к тому, что пакетный процесс, ранее выполнявшийся за час, неожиданно стал работать много часов из-за перезапуска, который раньше никогда не случался.

Это не новая особенность Oracle — она присутствует в базе данных, начиная с версии 4.0, когда была введена согласованность чтения. До лета 2003 года я сам не

вполне ориентировался, каким образом она работает, и после того как разобрался, на что эта особенность влияет, то смог ответить на массу вопросов типа “как это могло случиться?” из собственной практики. Я дал зарок почти никогда не применять автономные транзакции в триггерах, и мне пришлось пересмотреть способ реализации некоторых своих ранних приложений. Например, я ни в коем случае не стану отправлять сообщения электронной почты напрямую из триггера; вместо этого я всегда буду использовать DBMS_JOB или подобный пакет, чтобы производить отправку после фиксации транзакции. Это сделает процедуру отправки электронной почты *транзакционной*, т.е. если оператор, инициировавший запуск триггера и отправку электронной почты, будет перезапущен, то откат, который при этом произойдет, заодно отменит и запрос DBMS_JOB. Почти все нетранзакционные действия, которые выполнялись в триггерах, были модифицированы, чтобы делать всю работу после факта фиксации, обеспечивая согласованность в отношении транзакций.

Резюме

В настоящей главе мы представили множество сведений, которые иногда могли быть не особенно очевидными. Тем не менее, понимать все это очень важно. Например, если вы не осведомлены о перезапуске на уровне операторов, то не сможете разобраться, что служит причиной ряда загадочных ситуаций. То есть вы не сможете объяснить некоторые ежедневные эмпирические наблюдения. Фактически, если вы не знаете о перезапуске, то можете строить неверные предположения о причинах сбоев или ошибках конечных пользователей. Это может оказаться одной из невоспроизводимых проблем, когда случаются многие события, причем в специфичном порядке.

Мы ознакомились со смыслом уровней изоляции согласно формулировкам в стандарте SQL и посмотрели, как они реализованы в Oracle; временами мы проводили сравнение их реализации в Oracle с реализациями в других СУБД. Было показано, что в других реализациях (т.е. тех, которые задействуют блокировки чтения для обеспечения согласованности данных) существует крупный компромисс между параллелизмом и согласованностью. Чтобы получить доступ к данным с высокой степенью параллелизма, приходится понижать требования к согласованности ответов. Чтобы получать согласованные и корректные ответы, необходимо смириться с меньшей степенью параллелизма. Благодаря средству многоверсионности Oracle это не касается. В табл. 7.9 подведены итоги того, что можно ожидать от базы данных, в которой применяется блокировка чтения, в сравнении с подходом многоверсионности Oracle.

Управление параллелизмом и особенности его реализации в базе данных определенно следует понимать. Я восхвалял многоверсионность и согласованность чтения, но, как и все в этом мире, они имеют оборотную сторону. Если вы не понимаете многоверсионность и ее работу, то неизбежно будете допускать ошибки при проектировании приложений. Вспомните пример планировщика ресурсов из главы 1. В базе данных, не обладающей многоверсионностью и ассоциированными с ней неблокирующими чтениями, исходная логика программы должна быть очень тщательно проработана. Однако эта логика потерпела бы неудачу при реализации в среде Oracle. Она могла бы привести к нарушению целостности данных. Если вы не знаете, как работает многоверсионность, то будете писать программы, которые разрушают данные. Это так просто.

**Таблица 7.9. Сравнение уровней изоляции транзакций и поведения блокировки
в Oracle и базах данных, в которых используется блокировка чтения**

| Уровень изоляции | Реализация | Запись блокирует чтение | Чтение блокирует запись | Чувствитель- ное к взаимо- блокировкам чтение | Некор- ректные результаты запросов | Потерянные обновления | Эскалация блокировок или ограничения |
|------------------|------------|-------------------------------|-------------------------------|--|---|--------------------------|---|
| READ UNCOMMITTED | He Oracle | Нет | Нет | Нет | Да | Да | Да |
| READ COMMITTED | He Oracle | Да | Нет | Нет | Да | Да | Да |
| READ COMMITTED | Oracle | Нет | Нет | Нет | Нет | Нет* | Нет |
| REPEATABLE READ | He Oracle | Да | Да | Да | Нет | Нет | Да |
| SERIALIZABLE | He Oracle | Да | Да | Да | Нет | Нет | Да |
| SERIALIZABLE | Oracle | Нет | Нет | Нет | Нет | Нет | Нет |

* С оператором SELECT FOR UPDATE NOWAIT.

ГЛАВА 8

Транзакции

Транзакции — это одно из средств, которые отличают базы данных от файловых систем. Если вы находитесь на середине записи файла и операционная система аварийно завершается, то этот файл, скорее всего, будет поврежден. Надо сказать, что существуют файловые системы с ведением журналов, которые могут восстановить файл на какой-то момент времени в прошлом. Тем не менее, если нужно удерживать два файла в синхронизированном состоянии, то такие системы в этом не помогут — если вы обновите один файл, а в системе произойдет отказ, прежде чем будет завершено обновление второго файла, файлы не будут синхронизированы.

В этом и заключается основное предназначение *транзакций* — они переводят базу данных из одного согласованного состояния в другое. Это их функция. Когда вы фиксируете работу в базе данных, то можете быть уверены, что будут сохранены либо все внесенные изменения, либо ни одного из них. Более того, гарантируется реализация разнообразных правил и проверок, защищающих целостность данных.

В предыдущей главе мы обсуждали транзакции в терминах управления параллелизмом и говорили о том, как в результате использования многоверсионной согласованной по чтению модели Oracle транзакции каждый раз могут обеспечить непротиворечивость данных в условиях доступа к данным с высокой степенью параллелизма. Транзакции в Oracle удовлетворяют всем обязательным характеристикам ACID:

- атомарность (atomicity) — выполняются либо все действия, произведенные в транзакции, либо ни одного;
- согласованность (consistency) — транзакция перемещает базу данных из одного согласованного состояния в другое;
- изоляция (isolation) — результаты транзакции могут быть не видны другим транзакциям до тех пор, пока она не будет зафиксирована;
- постоянство (durability) — как только транзакция зафиксирована, она остается постоянной.

В частности, в предыдущей главе было показано, каким образом Oracle обеспечивает *согласованность* и *изоляцию*. Здесь мы сосредоточим основное внимание на концепции *атомарности* и ее применении в Oracle.

В настоящей главе речь пойдет о последствиях атомарности и о том, каким образом она влияет на операторы в Oracle. Мы рассмотрим операторы управления транзакциями, такие как COMMIT, SAVEPOINT и ROLLBACK, и обсудим способы прину-

длительного применения ограничений целостности в транзакциях. Кроме того, мы посмотрим, почему могут вырабатываться определенные плохие привычки в отношении транзакций, если вы занимались разработкой приложений в других СУБД. Мы также взглянем на распределенные транзакции и двухфазную фиксацию (Two Phase Commit — 2PC). Наконец, мы исследуем автономные транзакции — что они собой представляют, и какую играют роль.

Операторы управления транзакциями

В Oracle нет необходимости в операторе типа “начать транзакцию”. Транзакция начинается неявно с первым оператором, который модифицирует данные (первым оператором, получающим блокировку TX). Транзакцию можно начать явно с использованием оператора SET TRANSACTION или пакета DBMS_TRANSACTION, но в Oracle этот шаг не обязателен в отличие от ряда других СУБД. Оператор COMMIT или ROLLBACK явно завершает транзакцию.

На заметку! Не все операторы ROLLBACK создаются идентичными. Следует отметить, что команда ROLLBACK TO SAVEPOINT не завершает транзакцию! Это делает только надлежащий полный оператор ROLLBACK.

Вы должны всегда явно завершать свои транзакции с помощью оператора COMMIT или ROLLBACK, иначе применяемый инструмент или среда самостоятельно выберет тот или другой оператор. Если вы нормально выходите из сеанса SQL*Plus без фиксации или отката, то SQL*Plus предполагает, что вы хотите зафиксировать проделанную работу, и делает это. С другой стороны, если вы выходите из программы на Pro*C, то производится неявный откат. Никогда не полагайтесь на неявное поведение, потому что в будущем оно может измениться. Всегда завершайте явно свои транзакции с помощью COMMIT или ROLLBACK.

На заметку! В качестве примера того, что меняется с течением времени: SQL*Plus в Oracle 11g Release 2 и последующих версиях содержит настройку `exitcommit`, которая управляет тем, какую команду SQL*Plus выдаст при выходе — COMMIT или ROLLBACK. Таким образом, когда вы используете версию Oracle 11g Release 2, стандартное поведение, которое было характерно для средства SQL*Plus с момента его появления, вполне может оказаться другим!

Транзакции в Oracle являются *атомарными* — в том смысле, что каждый оператор, составляющий транзакцию, фиксируется (делается постоянным) или же происходит откат всех операторов. Эта защита распространяется также на индивидуальные операторы. Весь оператор либо полностью успешен, либо полностью неудачен. Обратите внимание на сказанное: производится откат “оператора”. Отказ одного оператора не вызывает автоматического отката ранее выполненных операторов. Их работа предохраняется и должна быть либо зафиксирована, либо отменена (путем отката) вами. Прежде чем мы углубимся в детали того, что в точности означает атомарность для оператора и транзакции, давайте посмотрим на разнообразные операторы управления транзакциями, имеющиеся в нашем распоряжении.

- **COMMIT.** Простейшая форма этого оператора выглядит просто как COMMIT. Можно применить более многословную форму и выдать COMMIT WORK, но эти две формы эквивалентны. Оператор COMMIT завершает транзакцию и делает любые изменения постоянными (сохраняет их). Существуют расширения оператора COMMIT, используемые в распределенных транзакциях, которые дают возможность снабжать COMMIT (помечать транзакцию) некоторым осмысленным комментарием и принудительно фиксировать сомнительную распределенную транзакцию. Есть также расширения, которые позволяют выполнять асинхронную фиксацию — фиксацию, которая в действительности нарушает концепцию постоянства. Мы слегка коснемся этого момента и посмотрим, когда подобное может пригодиться.
- **ROLLBACK.** Простейшая форма этого оператора выглядит как ROLLBACK. Опять-таки, существует более многословный вариант ROLLBACK WORK, но эти формы эквивалентны. Откат (rollback) завершает транзакцию и отменяет все незафиксированные изменения. Это делается путем чтения информации, хранящейся в сегментах отката/отмены (далее я буду называть их сегментами отмены (undo segments), как это принято в терминологии Oracle 10g и последующих версий) и восстановления блоков базы данных в состояние, в котором они пребывали до начала транзакции.
- **SAVEPOINT.** Оператор SAVEPOINT позволяет создать маркированную точку внутри транзакции. В одной транзакции можно иметь множество таких точек.
- **ROLLBACK TO <SAVEPOINT>.** Этот оператор применяется вместе с командой SAVEPOINT. Транзакцию можно откатить к желаемой маркированной точке, не откатывая работу, выполненную до нее. Таким образом, можно выдать два оператора UPDATE, за которыми следует SAVEPOINT, и затем два оператора DELETE. Если во время выполнения операторов DELETE происходит ошибка либо исключительная ситуация определенного вида, а вы перехватываете исключение и выдаете команду ROLLBACK TO SAVEPOINT, то транзакция будет подвержена откату к именованной точке сохранения, отменяя всю работу, которая была выполнена операторами DELETE, но оставляя незатронутым то, что сделано операторами UPDATE.
- **SET TRANSACTION.** Этот оператор позволяет устанавливать различные атрибуты транзакции, такие как уровень изоляции и то, предназначена она только для чтения или для чтения-записи. Оператор SET TRANSACTION можно также использовать для инструктирования транзакции о необходимости работы со специфическим сегментом отмены, когда применяется ручное управление отменой, но поступать так не рекомендуется. Автоматическое и ручное управление отменой более подробно обсуждается в главе 9.

Выше были перечислены все операторы управления транзакциями. Наиболее часто используемыми из них являются COMMIT и ROLLBACK. Оператор SAVEPOINT имеет до некоторой степени специальное назначение. Внутри Oracle он применяется часто; фактически он используется всякий раз, когда вы запускаете любой оператор SQL или PL/SQL, и вы также можете найти ему применение в собственных приложениях.

Атомарность

Теперь мы готовы посмотреть, что понимается под атомарностью оператора, процедуры и транзакции.

Атомарность на уровне оператора

Взгляните на следующий оператор:

```
Insert into t values ( 1 );
```

Отчетливо видно, что если этот оператор завершится неудачей из-за нарушенного ограничения, то строка не должна быть вставлена. Однако рассмотрим следующий пример, где оператор INSERT или DELETE для таблицы T инициирует запуск триггера, который соответствующим образом корректирует значение столбца CNT в таблице T2:

```
EODA@ORA12CR1> create table t2 ( cnt int );
Table created.
Таблица создана.

EODA@ORA12CR1> insert into t2 values ( 0 );
1 row created.
1 строка создана.

EODA@ORA12CR1> commit;
Commit complete.
Фиксация завершена.

EODA@ORA12CR1> create table t ( x int check ( x>0 ) );
Table created.
Таблица создана.

EODA@ORA12CR1> create trigger t_trigger
2 before insert or delete on t for each row
3 begin
4   if ( inserting ) then
5     update t2 set cnt = cnt +1;
6   else
7     update t2 set cnt = cnt -1;
8   end if;
9   dbms_output.put_line( 'I fired and updated ' ||
10                        sql%rowcount || ' rows' );
11 end;
12 /
Trigger created.
Триггер создан.
```

В такой ситуации уже менее ясно, что должно произойти. Если ошибка случится *после* инициирования триггера, должны ли результаты работы триггера сохраниться? То есть, если триггер активизировался и обновил таблицу T2, но строка не была вставлена в таблицу T, то каким должен быть исход? Ясно, что мы не хотим, чтобы значение столбца CNT в T2 инкрементировалось, если строка не была действительно вставлена в T. К счастью, в Oracle исходный оператор, введенный клиентом — в этом случае INSERT INTO T — завершается либо полностью успешно, либо полностью неудачно. Данный оператор является атомарным.

Это можно подтвердить следующим образом:

```

EODA@ORA12CR1> set serveroutput on
EODA@ORA12CR1> insert into t values (1);
I fired and updated 1 rows
Я запустился и обновил 1 строку
1 row created.
1 строка создана.

EODA@ORA12CR1> insert into t values (-1);
I fired and updated 1 rows
Я запустился и обновил 1 строку
insert into t values (-1)
*
ERROR at line 1:
ORA-02290: check constraint (EODA.SYS_C0061484) violated
ОШИБКА в строке 1:
ORA-02290: нарушено проверочное ограничение целостности (EODA.SYS_C0061484)
EODA@ORA12CR1> select * from t2;

      CNT
-----
      1

```

На заметку! Для того чтобы увидеть срабатывание триггера в SQL*Plus в среде Oracle9i Release 2 и предшествующих версий, после второго оператора INSERT понадобится добавить строку кода EXEC NULL. Причина в том, что в этих выпусках SQL*Plus не извлекает и не отображает информацию DBMS_OUTPUT после отказа оператора DML. В Oracle 10g и последующих версиях все отображается.

Итак, в таблицу T была успешно вставлена одна строка, и мы надлежащим образом получили сообщение I fired and updated 1 rows (Я запустился и обновил 1 строку). Следующий оператор INSERT нарушает ограничение целостности, имеющееся в таблице T. Наличие сообщения DBMS_OUTPUT свидетельствует о том, что триггер в таблице T действительно сработал. Он успешно произвел свои обновления таблицы T2. Можно было ожидать, что T2 теперь содержит значение 2, но на самом деле оно равно 1. СУБД Oracle сделала *исходный* оператор INSERT атомарным — первоначальный INSERT INTO T является оператором, и любой сторонний эффект, порожденный этим оператором, считается его частью.

СУБД Oracle достигает такой атомарности уровня оператора за счет того, что молча помещает каждое обращение к базе данных внутрь SAVEPOINT. Предыдущие два оператора INSERT в действительности трактовались так:

```

Savepoint оператор1;
  Insert into t values ( 1 );
Если возникла ошибка, то откатить оператор1;
Savepoint оператор2;
  Insert into t values ( -1 );
Если возникла ошибка, то откатить оператор2;

```

Для программистов, работавших с Sybase или SQL Server, поначалу это может сбивать с толку. В упомянутых базах данных *справедливо в точности обратное*.

В этих системах триггеры выполняются независимо от вызвавшего их оператора. Если возникает ошибка, триггеры должны явно откатить собственную работу и затем сгенерировать другую ошибку, чтобы произвести откат оператора, который их инициировал. В противном случае работа, сделанная триггером, может сохраниться, даже если запустивший его оператор или другая часть этого оператора, в конечном счете, отказал.

В Oracle такая атомарность на уровне оператора распространяется настолько глубоко, насколько это необходимо. Если в предыдущем примере оператор INSERT INTO T запустит триггер, который обновляет другую таблицу, и эта другая таблица имеет триггер, осуществляющий удаление в третьей таблице (и т.д. и т.п.), то либо *вся* работа завершится успешно, либо не будет выполнено *ни одного* действия. Чтобы обеспечить это, не придется писать какой-то специальный код — просто именно так все работает.

Атомарность на уровне процедуры

Интересно отметить, что Oracle рассматривает блоки PL/SQL также в качестве операторов. Взгляните на следующую хранимую процедуру и сброс примеров таблиц:

```
EODA@ORA12CR1> create or replace procedure p
2 as
3 begin
4     insert into t values ( 1 );
5     insert into t values (-1 );
6 end;
7 /
```

Procedure created.

Процедура создана.

```
EODA@ORA12CR1> delete from t;
```

0 rows deleted.

0 строк удалено.

```
EODA@ORA12CR1> update t2 set cnt = 0;
```

1 row updated.

1 строка обновлена.

```
EODA@ORA12CR1> commit;
```

Commit complete.

Фиксация завершена.

```
EODA@ORA12CR1> select * from t;
```

no rows selected

нет выбранных строк

```
EODA@ORA12CR1> select * from t2;
```

CNT

0

Таким образом, мы располагаем процедурой, которая, как мы знаем, завершится неудачей, и в этом случае второй оператор INSERT всегда даст сбой. Давайте посмотрим, что произойдет, если мы запустим эту хранимую процедуру:

```

EODA@ORA12CR1> begin
  2     p;
  3 end;
  4 /
I fired and updated 1 rows
I fired and updated 1 rows
Я запустился и обновил 1 строку
Я запустился и обновил 1 строку
begin
*
ERROR at line 1:
ORA-02290: check constraint (EODA.SYS_C0061484) violated
ORA-06512: at "EODA.P", line 5
ORA-06512: at line 2
ОШИБКА в строке 1:
ORA-02290: нарушено проверочное ограничение целостности (EODA.SYS_C0061484)
ORA-06512: в EODA.P, строка 5
ORA-06512: в строке 2

EODA@ORA12CR1> select * from t;
no rows selected
нет выбранных строк

EODA@ORA12CR1> select * from t2;
      CNT
-----
      0

```

Как видите, Oracle трактует вызов хранимой процедуры как атомарный оператор. Клиент отправляет блок кода `BEGIN P; END;`, и Oracle помещает его в оболочку `SAVEPOINT`. Поскольку процедура `P` терпит отказ, Oracle восстанавливает базу данных в состояние, которое она имела в точке, непосредственно предшествовавшей вызову `P`.

На заметку! Описанное поведение — атомарность на уровне оператора — полагается на то, что в процедуре PL/SQL не выполняются какие-либо фиксации или откаты. По моему мнению, операторы `COMMIT` и `ROLLBACK` вообще не должны использоваться в PL/SQL; только коду, вызывающему хранимую процедуру PL/SQL, известно, когда транзакция завершается. Выдача `COMMIT` или `ROLLBACK` в разрабатываемых процедурах PL/SQL является плохой практикой программирования.

Теперь, если мы отправим слегка отличающийся блок, то получим совершенно другие результаты:

```

EODA@ORA12CR1> begin
  2     p;
  3 exception
  4     when others then
  5         dbms_output.put_line( 'Error!!!! ' || sqlerrm );
  6 end;
  7 /
I fired and updated 1 rows
I fired and updated 1 rows
I fired and updated 1 rows

```



```
Error!!!! ORA-02290: check constraint (EODA.SYS_C0061484) violated
PL/SQL procedure successfully completed.
```

Ошибка!!!! ORA-02290: нарушено проверочное ограничение целостности (EODA.SYS_C0061484)

Процедура PL/SQL успешно завершена.

```
EODA@ORA12CR1> select * from t;
```

```
      X
-----
      1
```

```
EODA@ORA12CR1> select * from t2;
```

```
      CNT
-----
      1
```

```
EODA@ORA12CR1> rollback;
```

```
Rollback complete.
```

Откат выполнен.

Здесь мы запустили блок кода, который игнорирует все возможные ошибки, и отличия в выводе оказались огромными. В то время как первый вызов Р не привел к каким-то изменениям, на этот раз первый оператор INSERT выполняется успешно и столбец CNT в T2 соответствующим образом инкрементируется.

СУБД Oracle считает “оператором” блок кода, отправленный клиентом. Этот оператор завершается успешно за счет самостоятельного перехвата и игнорирования ошибки, так что конструкция “Если возникла ошибка, то откатить...” не вступает в действие и после выполнения Oracle не производит откат к точке сохранения. Следовательно, частичная работа, сделанная Р, сохраняется. Причина предохранения этой частичной работы в первую очередь объясняется наличием атомарности на уровне оператора внутри процедуры Р: каждый оператор в Р является атомарным. Процедура Р становится клиентом Oracle, когда отправляет свои два оператора INSERT. Каждый оператор INSERT либо целиком успешен, либо нет. Об этом свидетельствует тот факт, что мы можем видеть, что триггер таблицы Т запускался два раза и дважды обновил T2, хотя счетчик в T2 отражает только один оператор UPDATE. Второй оператор INSERT, выполненный в Р, находится внутри неявной оболочки SAVEPOINT.

Конструкция WHEN OTHERS

Я считаю ошибочным практически любой код, который содержит обработчик исключений WHEN OTHERS, не включающий также оператор RAISE или вызов процедуры RAISE_APPLICATION_ERROR с целью повторной генерации исключения. Он молча игнорирует ошибку и меняет семантику транзакции. Перехват WHEN OTHERS и трансляция исключения в возвращаемый код в старом стиле изменяет ожидаемое поведение базы данных.

В действительности, когда версия Oracle 11g Release 1 находилась еще на стадии проектирования, у меня появилась возможность отправить три запроса для новых средств в PL/SQL. Я получил замечательный шанс, и мое первое предложение состояло в том, чтобы вообще исключить из языка конструкцию WHEN OTHERS. Приводимая мною аргументация была проста: самой распространенной причиной вносимых разработчиками ошибок — именно *самой распространенной* — было отсутствие за WHEN OTHERS оператора RAISE или обращения к RAISE_APPLICATION_ERROR. Я чувствовал, что мир станет безопаснее без этого языкового средства. Команда, занимающаяся

реализацией PL/SQL, конечно же, не могла на это пойти, но они сделали другую полезную вещь. Они обеспечили генерацию компилятором PL/SQL предупреждающего сообщения в случае, если в коде встречалась конструкция WHEN OTHERS, за которой не следовал вызов RAISE или RAISE_APPLICATION_ERROR. Например:

```
EODA@ORA12CR1> alter session set
  2  PLSQL_Warnings = 'enable:all'
  3  /
```

Session altered.

Сеанс изменен.

```
EODA@ORA12CR1> create or replace procedure some_proc( p_str in varchar2 )
  2  as
  3  begin
  4      dbms_output.put_line( p_str );
  5  exception
  6      when others
  7      then
  8          -- вызов процедуры log_error()
  9          null;
 10  end;
 11  /
```

SP2-0804: Procedure created with compilation warnings

SP2-0804: процедура создана с предупреждениями компилятора

```
EODA@ORA12CR1> show errors procedure some_proc
```

Errors for PROCEDURE P:

LINE/COL ERROR

```
-----
1/1  PLW-05018: unit SOME_PROC omitted optional AUTHID clause;
      default value DEFINER used
1/1  PLW-05018: в модуле SOME_PROC опущена необязательная конструкция AUTHID;
      используется стандартное значение DEFINER
6/10 PLW-06009: procedure "SOME_PROC" OTHERS handler does not end in
      RAISE or RAISE_APPLICATION_ERROR
6/8  PLW-06009: обработчик OTHERS процедуры SOME_PROC не завершается
      RAISE или RAISE_APPLICATION_ERROR
```

Таким образом, если вы включаете в код конструкцию WHEN OTHERS, но за ней не следует оператор RAISE или вызов RAISE_APPLICATION_ERROR, то имейте в виду, что вы почти наверняка столкнетесь с ошибкой в разработанном коде — ошибкой, которую внесли собственноручно.

Отличие между двумя блоками кода, один из которых содержит конструкцию WHEN OTHERS, а другой нет, является довольно тонким, и вы должны учитывать его в своих приложениях. Добавление обработчика исключений к блоку кода PL/SQL может радикально изменить его поведение. Существует другой способ кодирования этого, при котором восстанавливается атомарность на уровне оператора для всего блока PL/SQL:

```
EODA@ORA12CR1> begin
  2  savepoint sp;
  3  p;
  4  exception
  5      when others then
  6          rollback to sp;
```

```

7          dbms_output.put_line( 'Error!!!! ' || sqlerrm );
8 end;
9 /
I fired and updated 1 rows
I fired and updated 1 rows
Error!!!! ORA-02290: check constraint (EODA.SYS_C0061484) violated
PL/SQL procedure successfully completed.
Ошибка!!! ORA-02290: нарушено проверочное ограничение целостности
(EODA.SYS_C0061484)
Процедура PL/SQL успешно завершена.

EODA@ORA12CR1> select * from t;
no rows selected
нет выбранных строк

EODA@ORA12CR1> select * from t2;
      CNT
-----
      0

```

Внимание! В предыдущем коде был продемонстрирован пример исключительно плохой практики. В общем случае вы никогда не должны перехватывать WHEN OTHERS, равно как и явно кодировать то, что уже обеспечивает Oracle — до тех пор, пока речь идет о семантике транзакции.

Имитируя здесь то, что Oracle обычно делает с помощью SAVEPOINT, мы можем восстановить исходное поведение, по-прежнему перехватывая и “игнорируя” ошибку. Этот пример был приведен только в целях иллюстрации; он представляет собой в высшей степени неудачную практику.

Атомарность на уровне транзакции

Общая цель транзакции — набора операторов SQL, выполняемых вместе в качестве единицы работы, — состоит в переводе базы данных из одного согласованного состояния в другое. Для достижения этой цели транзакции также сделаны атомарными, т.е. полный набор успешных действий, произведенных в транзакции, либо целиком фиксируется и становится постоянным, либо откатывается и отменяется. Как и оператор, транзакция является атомарной единицей работы. По сообщению базы данных об “успехе” после фиксации транзакции вы знаете о том, что вся работа, выполненная в транзакции, стала постоянной.

Операторы DDL и атомарность

Полезно отметить, что в Oracle есть определенный класс операторов, которые являются атомарными, но только на уровне оператора. Операторы языка определения данных (Data Definition Language — DDL) реализованы в следующем стиле.

1. Они начинаются с фиксации любой невыполненной работы, завершая любую транзакцию, которая могла быть в наличии.
2. Они выполняют операцию DDL, такую как CREATE TABLE.
3. Они фиксируют операцию DDL, если она была успешной, либо производят ее откат в противном случае.

Это значит, что всякий раз, выдавая оператор DDL, такой как CREATE, ALTER и т.д., вы должны ожидать, что существующая транзакция будет немедленно зафиксирована и последующая команда DDL, подлежащая выполнению, либо фиксируется и ее результаты делаются постоянными, либо откатывается в случае возникновения любой ошибки. Операторы DDL никоим образом не нарушают концепции ACID, но о том факте, что они производят фиксацию, определенно следует знать.

Постоянство

Обычно, когда транзакция зафиксирована, внесенные ею изменения являются постоянными; вы можете рассчитывать на то, что они будут находиться в базе данных, даже если в базе данных произойдет катастрофический отказ немедленно после завершения фиксации. Однако есть два специфичных случая, когда это не так.

- Вы применяете расширения WRITE (доступные в Oracle 10g Release 2 и последующих версиях), предназначенные для оператора COMMIT.
- Вы выдаете операторы COMMIT в нераспределенном (обращающемся только к одной базе данных, без связей баз данных) блоке кода PL/SQL.

Давайте рассмотрим эти случаи по очереди.

Расширения WRITE оператора COMMIT

Начиная с версии Oracle 10g Release 2, к операторам COMMIT можно добавлять конструкцию WRITE. Эта конструкция позволяет фиксировать либо с вариантом WAIT, чтобы сгенерированные данные повторения были записаны на диск (по умолчанию), либо с вариантом NOWAIT — без ожидания записи данных повторения действий. Опция NOWAIT — это средство, которое должно использоваться осторожно, продуманно и с четким пониманием его смысла.

Обычно COMMIT является синхронным процессом. Приложение вызывает COMMIT и затем *ожидает*, пока обработка COMMIT не будет полностью завершена (более детально это описано в главе 9). Таким было поведение COMMIT во всех выпусках, предшествующих Oracle 10g Release 2, и это стандартное поведение в Oracle 10g Release 2 и последующих версиях.

В текущих выпусках базы данных вместо ожидания завершения фиксации, что может потребовать заметного времени, т.к. фиксация предусматривает физическую запись — физический ввод-вывод — в файлы журналов повторения, хранящиеся на диске, можно провести фиксацию в фоновом режиме без ожидания. С таким приемом связан побочный эффект: *вы не можете иметь уверенность в том, что фиксация стала постоянной*. То есть ваше приложение может получить ответ от базы данных о приеме отправленной асинхронной фиксации, другие сеансы могут быть в состоянии видеть изменения, но позже обнаружится, что транзакция, которую вы считали зафиксированной, на самом деле таковой не стала. Такая ситуация будет возникать очень редко и всегда будет связана с серьезным отказом оборудования или программного обеспечения. Она предполагает аварийное завершение базы данных, препятствующее сохранению асинхронной фиксации, а это значит, что экземпляр базы данных или компьютер, на котором этот экземпляр функционирует, полностью вышел из строя.

Итак, если транзакции задуманы быть постоянными, то каково потенциальное применение средства, которое вполне может сделать их непостоянными? Ответ: низкоуровневое управление производительностью. Выдавая в приложении `COMMIT`, вы тем самым просите процесс `LGWR` взять сгенерированные данные повторения и удостовериться, что они записаны в файлы оперативного журнала повторения. Физический ввод-вывод, который при этом происходит, выполняется довольно медленно; он требует относительно долгого времени для записи данных на диск. Таким образом, оператор `COMMIT` может работать дольше, чем операторы `DML` в самой транзакции! Если вы превратите `COMMIT` в асинхронный процесс, то устранили необходимость ожидания физического ввода-вывода в клиентском приложении, возможно, сделав так, что оно будет выглядеть функционирующим быстрее — особенно, если оно выполняет множество операций `COMMIT`.

Это может наводить на мысль о том, что `COMMIT WRITE NOWAIT` имеет смысл использовать все время; в конце концов, разве производительность — не самая важная вещь на свете? Вовсе нет. Большую часть времени требуется постоянство, по умолчанию обеспечиваемое `COMMIT`. Когда вы выдаете `COMMIT` и сообщаете конечному пользователю о выполненной фиксации, то должны быть уверены в том, что изменение является постоянным. Оно будет записано в базу данных, даже если сразу после `COMMIT` произойдет отказ базы данных и/или оборудования. Если вы уведомляете конечного пользователя о том, что заказ 12352 был размещен, то должны гарантировать, что данный заказ 12352 действительно размещен и сохранен. Поэтому почти для каждого приложения `COMMIT WRITE WAIT` — единственный корректный вариант (обратите внимание, что достаточно записать просто `COMMIT`, т.к. по умолчанию принимается настройка `WRITE WAIT`).

Когда в таком случае может понадобиться это средство фиксации без ожидания? На ум приходят три сценария.

- Специальная программа загрузки данных. Она должна быть специальной, поскольку будет иметь дополнительную логику, которая учитывает тот факт, что фиксация может не сохраниться из-за отказа системы.
- Приложение, обрабатывающее действующий источник данных определенного вида, скажем, источник биржевых котировок, которое подразумевает запись в базу больших объемов чувствительной ко времени информации. Когда база данных становится недоступной, то поток данных продолжает поступать и данные, сгенерированные во время отказа системы, никогда не будут обработаны (ведь торги на бирже не останавливаются только из-за того, что ваша база данных аварийно завершила свою работу). То, что такие данные не будут обработаны, вполне нормально, поскольку биржевая информация чувствительна ко времени и по прошествии нескольких секунд все равно будет перезаписана новыми данными.
- Приложение, реализующее собственный механизм “организации очереди”, например, предусматривающий наличие для данных таблицы со столбцом `PROCESSED_FLAG` (признак обработки). При поступлении новых данных они вставляются со значением `PROCESSED_FLAG='N'` (не обработаны). Другая процедура занимается чтением записей с `PROCESSED_FLAG='N'`, выполнением небольшой, быстрой транзакции и обновлением `PROCESSED_FLAG` с 'N' на 'Y'. Если такая транзакция фиксируется, но фиксация позже отменяется (из-за

отказа системы), то ничего страшного не происходит, потому что приложение, которое обрабатывает эти записи, просто обработает упомянутую запись снова — оно является “перезапускаемым”.

Если взглянуть на эти категории приложений, вы заметите, что все три представляют фоновые, неинтерактивные приложения. Они не взаимодействуют с человеком напрямую. В любом приложении, которое взаимодействует с человеком и сообщает ему, что фиксация завершена, должна применяться синхронная фиксация. Асинхронные фиксации не подходят для онлайн-приложений, ориентированных на пользователя. Они пригодны только для пакетных приложений, которые способны автоматически перезапускаться после сбоя. Интерактивные приложения не перезапускаются автоматически в случае отказа — транзакцию должен повторить человек. Следовательно, у вас есть еще один признак, который указывает на то, должна ли эта возможность приниматься во внимание: каким приложением вы располагаете — пакетным или интерактивным? Если оно не пакетное, то необходимо использовать синхронные фиксации.

Таким образом, за рамками трех названных категорий пакетных приложений средство `COMMIT WRITE NOWAIT`, скорее всего, применяться не должно. Если вы все же используете его, то должны задаться вопросом: что случится, если приложение сообщит о *выполненной фиксации*, но позже фиксация отменяется? Вам необходимо иметь возможность ответить на этот вопрос и прийти к заключению, что такое допустимо. Если же вы не можете ответить на этот вопрос, или же потеря зафиксированного изменения может привести к серьезным последствиям, то от средства асинхронной фиксации следует отказаться.

Операторы `COMMIT` в нераспределенном блоке PL/SQL

С момента появления в версии Oracle 6 инструмента PL/SQL в нем автоматически применялась асинхронная фиксация. Такой подход работал, потому что инструмент PL/SQL целиком подобен пакетной программе — конечный пользователь не знает результатов процедуры до тех пор, пока она не будет полностью завершена. По той же причине асинхронная фиксация используется только в нераспределенном блоке кода PL/SQL; если в нем вовлечено более одной базы данных, тогда есть две сущности, две базы данных, полагающиеся на постоянство фиксации. В этом случае должны применяться синхронные протоколы или же изменение может быть зафиксировано в одной базе данных, но не в другой.

На заметку! Разумеется, конвейерные функции PL/SQL отличаются от “нормальных” функций PL/SQL. В нормальных функциях PL/SQL результаты не будут известны до конца вызова хранимой процедуры. Вообще говоря, конвейерные функции способны возвращать данные клиенту задолго до своего завершения (они возвращают клиенту “порции” данных, по одной за раз). Но поскольку конвейерные функции вызываются из операторов `SELECT`, и в любом случае не будут фиксироваться, в этом обсуждении они не участвуют.

По этой причине средство PL/SQL было разработано для использования асинхронной фиксации, позволяя оператору `COMMIT` в PL/SQL не ожидать завершения физического ввода-вывода (избегая ожидания синхронизации журнального файла). Это вовсе не означает, что нельзя полагаться на то, что процедура PL/SQL, которая фиксирует и возвращает управление приложению, не будет постоянной в отноше-

нии своих изменений. Напротив, PL/SQL будет ожидать записи на диск своей информации повторения перед возвратом управления клиентскому приложению, но будет ждать только однажды, прямо перед возвратом.

На заметку! В приведенном ниже примере демонстрируется плохая практика, которую я называю “замедленной обработкой” или “обработкой строки за строкой”, т.к. в реляционной базе данных они практически являются синонимами. Пример предназначен только для иллюстрации того, как PL/SQL обрабатывает оператор COMMIT.

Для начала создадим таблицу T:

```
EODA@ORA12CR1> create table t
2 as
3 select *
4   from all_objects
5  where 1=0
6 /
```

Table created.

Таблица создана.

Рассмотрим следующую процедуру PL/SQL:

```
EODA@ORA12CR1> create or replace procedure p
2 as
3 begin
4   for x in ( select * from all_objects )
5   loop
6       insert into t values X;
7       commit;
8   end loop;
9 end;
10 /
```

Procedure created.

Процедура создана.

Этот код PL/SQL читает по одной записи за раз из ALL_OBJECTS, вставляет запись в таблицу T и фиксирует каждую запись после вставки. Логически он эквивалентен такому коду:

```
EODA@ORA12CR1> create or replace procedure p
2 as
3 begin
4   for x in ( select * from all_objects )
5   loop
6       insert into t values X;
7       commit write NOWAIT;
8   end loop;
9
10  -- сделать здесь внутренний вызов, чтобы удостовериться
11  -- в записи данных повторения процессом LGWR
12 end;
13 /
```

Procedure created.

Процедура создана.

Таким образом, фиксации, выполняемые в процедуре, производятся с конструкцией `WRITE NOWAIT` и перед тем, как блок кода PL/SQL возвращает управление клиентскому приложению. Инструмент PL/SQL гарантирует, что сгенерированная им последняя порция данных повторения будет благополучно записана на диск, что обеспечивает постоянство блока кода PL/SQL и внесенных им изменений.

В главе 11 вы увидите характерные результаты этой возможности PL/SQL при измерении производительности индексов по реверсированным ключам. Если хотите посмотреть, каким образом PL/SQL работает в описанной выше манере, загляните туда, чтобы ознакомиться с показателями производительности индексов по реверсированным ключам.

Ограничения целостности и транзакции

Интересно отметить точные обстоятельства, когда проверяются ограничения целостности. По умолчанию ограничения целостности проверяются после обработки всего оператора SQL. Существуют также откладываемые ограничения, которые допускают отсрочивание проверки ограничений целостности до момента, когда либо приложение запросит их проверку, выдав `SET CONSTRAINTS ALL IMMEDIATE`, либо встретится оператор `COMMIT`.

Ограничения **IMMEDIATE**

Для первой части этого обсуждения мы предположим, что ограничения находятся в режиме `IMMEDIATE`, что является нормой. В таком случае ограничения целостности проверяются немедленно после обработки всего оператора SQL. Обратите внимание на применение термина “оператор SQL”, а не просто “оператор”. Если в хранимой процедуре PL/SQL имеется много операторов SQL, то проверка ограничений целостности будет следовать сразу же за выполнением каждого из них, а не после того, как завершится вся процедура.

Итак, почему ограничения проверяются *после* выполнения оператора SQL? Почему не *во время*? Дело в том, что для отдельных операторов вполне естественно делать отдельные строки таблицы на мгновение несогласованными. Просмотр частичной работы оператора может привести к отклонению этих результатов базой данных Oracle, даже если конечный исход всей работы будет корректным. Например, предположим, что есть следующая таблица:

```
EODA@ORA12CR1> create table t ( x int unique );
Table created.
Таблица создана.

EODA@ORA12CR1> insert into t values ( 1 );
1 row created.
1 строка создана.

EODA@ORA12CR1> insert into t values ( 2 );
1 row created.
1 строка создана.

EODA@ORA12CR1> commit;
Commit complete.
Фиксация завершена.
```


И мы хотим выполнить многострочный оператор UPDATE:

```

EODA@ORA12CR1> update t set x=x-1;
2 rows updated.
2 строки обновлено.

```

Если бы база данных Oracle проверяла ограничение после обновления каждой строки, то в любой день существует вероятность 50/50 того, что UPDATE потерпит неудачу. Доступ к строкам в таблице T производится в *некотором* порядке, и если Oracle обновит строку X=1 первой, то получится мгновенное дублированное значение для X, и оператор UPDATE будет отклонен. Поскольку Oracle терпеливо ожидает окончания работы оператора, данный оператор завершается успешно, т.к. к этому времени дубликатов уже нет.

Ограничения DEFERRABLE и каскадные обновления

Начиная с версии Oracle 8.0, мы также располагаем возможностью *отложить* проверку ограничений, что довольно полезно для разнообразных операций. На ум сразу же приходит требование каскадным образом распространить обновление первичного ключа на дочерние ключи. Многие люди заявят, что необходимость в этом никогда не должна возникать, ведь первичные ключи неизменны (я тоже принадлежу к их числу), но многие другие настаивают на своем желании иметь каскадное обновление. Отложенные ограничения делают это возможным.

На заметку! Выполнение каскадных обновлений для модификации первичного ключа считается чрезвычайно плохой практикой. Это нарушает сам замысел первичного ключа. Если вы должны предпринять такое действие один раз, чтобы исправить некорректную информацию — это одно, но если вы постоянно поступаете подобным образом в своем приложении, то стоит заново обдумать проектное решение — в качестве ключа были выбраны неправильные атрибуты!

В ранних выпусках Oracle каскадное обновление можно было выполнять, но это требовало огромного объема работы и обладало определенными пределами. Благодаря отложенным ограничениям, задача становится почти тривиальной. Код может выглядеть примерно так:

```

EODA@ORA12CR1> create table parent
2 ( pk int primary key )
3 /
Table created.
Таблица создана.

EODA@ORA12CR1> create table child
2 ( fk constraint child_fk_parent
3     references parent(pk)
4     deferrable
5     initially immediate
6 )
7 /
Table created.
Таблица создана.

```

```

EODA@ORA12CR1> insert into parent values ( 1 );
1 row created.
1 строка создана.

EODA@ORA12CR1> insert into child values ( 1 );
1 row created.
1 строка создана.

```

У нас есть родительская таблица PARENT и дочерняя таблица CHILD. Таблица CHILD ссылается на таблицу PARENT, и ограничение, используемое для внедрения этого правила, называется CHILD_FK_PARENT (дочерний внешний ключ к родительской таблице). Ограничение создано как DEFERRABLE, но установлено в INITIALLY IMMEDIATE. Это значит, что проверку ограничения можно отложить до выдачи COMMIT или до какого-то другого момента. Однако по умолчанию оно будет проверяться на уровне оператора. Это наиболее распространенный случай применения отложенных ограничений. Большинство существующих приложений не проверяют условие нарушения ограничений при выполнении оператора COMMIT, и лучше это не менять. Согласно определению, таблица CHILD ведет себя таким же образом, как всегда должны вести себя таблицы, но при этом предоставляет нам возможность явно изменить ее поведение. Теперь давайте попробуем запустить некоторый код DML в отношении таблиц и посмотрим, что произойдет:

```

EODA@ORA12CR1> update parent set pk = 2;
update parent set pk = 2
*
ERROR at line 1:
ORA-02292: integrity constraint (EODA.CHILD_FK_PARENT) violated -
child record found
ОШИБКА в строке 1:
ORA-02292: нарушение ограничения целостности (EODA.CHILD_FK_PARENT) -
найдена дочерняя запись

```

Поскольку ограничение находится в режиме IMMEDIATE, оператор UPDATE завершился неудачей. Изменим режим и попробуем снова:

```

EODA@ORA12CR1> set constraint child_fk_parent deferred;
Constraint set.
Ограничение установлено.

EODA@ORA12CR1> update parent set pk = 2;
1 row updated.
1 строка обновлена.

```

Теперь оператор выполняется успешно. В целях иллюстрации я покажу, как явно проверить отложенное ограничение перед фиксацией, чтобы увидеть, согласуются ли модификации с бизнес-правилами (другими словами, проверить, не нарушены ли в текущий момент ограничения). Это неплохо сделать перед фиксацией или передачей управления какой-то другой части программы (которая может не ожидать наличия отложенных ограничений):

```

EODA@ORA12CR1> set constraint child_fk_parent immediate;
set constraint child_fk_parent immediate
*
ERROR at line 1:

```

```
ORA-02291: integrity constraint (EODA.CHILD_FK_PARENT) violated -
parent key not found
```

ОШИБКА в строке 1:

```
ORA-02291: нарушение ограничения целостности (EODA.CHILD_FK_PARENT) -
не найден родительский ключ
```

Как и ожидалось, оператор завершается неудачей и немедленно возвращает ошибку, поскольку мы знали, что ограничение было нарушено. Откат оператора UPDATE для таблицы PARENT не производился (т.к. это нарушило бы атомарность на уровне оператора); он все еще ожидает выполнения. Также обратите внимание, что наша транзакция по-прежнему работает с отложенным ограничением CHILD_FK_PARENT, поскольку команда SET CONSTRAINT завершилась неудачей. Теперь продолжим каскадированием оператора UPDATE в отношении таблицы CHILD:

```
EODA@ORA12CR1> update child set fk = 2;
```

```
1 row updated.
```

1 строка обновлена.

```
EODA@ORA12CR1> set constraint child_fk_parent immediate;
```

```
Constraint set.
```

Ограничение установлено.

```
EODA@ORA12CR1> commit;
```

```
Commit complete.
```

Фиксация завершена.

Именно так это работает. Чтобы отложить ограничение, вы должны создать его так, как продемонстрировано — удалить и создать ограничение заново, чтобы изменить его с неотложенного на отложенное. Это может привести к уверенности в том, что все ограничения следует создавать как DEFERRABLE и INITIALLY IMMEDIATE просто на тот случай, если в какой-то момент их понадобится отложить. Вообще говоря, это *не* так. Разрешать ограничениям быть отложенными необходимо лишь при наличии действительной потребности в этом. Создавая отложенные ограничения, вы вносите в физическую реализацию (в структуру ваших данных) расхождение, которые могут быть неочевидными. Например, если вы создаете отложенное ограничение UNIQUE или PRIMARY KEY, то индекс, который Oracle создает для поддержки реализации этого ограничения, будет неуникальным. Обычно вы ожидаете, что для применения ограничения уникальности предусмотрен уникальный индекс, но поскольку вы указали, что ограничение может быть временно проигнорировано, оно не может использовать уникальный индекс. Будут наблюдаться и другие тонкие изменения, например, в ограничениях NOT NULL. В главе 11 показано, что индекс на столбце NOT NULL может применяться во многих ситуациях, когда это не может делать похожий индекс на столбце NULL. Если вы позволите ограничению NOT NULL быть отложенным, то оптимизатор начнет трактовать данный столбец, как будто он поддерживает значения NULL, потому что он на самом деле *поддерживает* значения NULL во время транзакции. Например, предположим, что есть таблица со следующими столбцами и данными:

```
EODA@ORA12CR1> create table t
```

```
2 ( x int constraint x_not_null not null deferrable,
```

```
3 y int constraint y_not_null not null,
```

```
4 z varchar2(30)
```

```
5 );
```

Table created.

Таблица создана.

```
EODA@ORA12CR1> insert into t(x,y,z)
  2 select rownum, rownum, rpad('x',30,'x')
  3   from all_users;
```

45 rows created.

45 строк создано.

```
EODA@ORA12CR1> exec dbms_stats.gather_table_stats( user, 'T' );
```

PL/SQL procedure successfully completed.

Процедура PL/SQL успешно завершена.

В этом примере столбец X создается так, что при выдаче COMMIT значение X не должно быть NULL. Однако во время транзакции в X допускаются значения NULL, т.к. ограничение является отложенным. С другой стороны, столбец Y всегда NOT NULL. Давайте проиндексируем столбец Y:

```
EODA@ORA12CR1> create index t_idx on t(y);
```

Index created.

Индекс создан.

После этого запустим запрос, который использует этот индекс на Y, но только если Y является NOT NULL, как в следующем запросе:

```
EODA@ORA12CR1> explain plan for select count(*) from t;
```

Explained.

Объяснено.

```
EODA@ORA12CR1> select * from table(dbms_xplan.
display(null,null,'BASIC'));
```

```
-----
| Id | Operation          | Name |
-----
|  0 | SELECT STATEMENT   |      |
|  1 |  SORT AGGREGATE    |      |
|  2 |    INDEX FULL SCAN | T_IDX |
-----
```

Вы увидите, что для подсчета строк оптимизатор предпочел задействовать небольшой индекс на Y, а не выполнять полное сканирование всей таблицы T. Но давайте удалим этот индекс и взамен проиндексируем столбец X:

```
EODA@ORA12CR1> drop index t_idx;
```

Index dropped.

Индекс удален.

```
EODA@ORA12CR1> create index t_idx on t(x);
```

Index created.

Индекс создан.

Если теперь снова запустить запрос для подсчета строк, то обнаружится, что база данных не применяет этот индекс — и в действительности не может это делать:

```
EODA@ORA12CR1> explain plan for select count(*) from t;
```

Explained.

Объяснено.

```
EODA@ORA12CR1> select * from table(dbms_xplan.display(null,null,'BASIC'));
```

| Id | Operation | Name |
|----|-------------------|------|
| 0 | SELECT STATEMENT | |
| 1 | SORT AGGREGATE | |
| 2 | TABLE ACCESS FULL | T |

Запрос полностью сканирует таблицу. Полное сканирование требуется для подсчета строк. Это вызвано тем фактом, что в индексе со структурой В-дерева (B*Tree) базы данных Oracle записи индексного ключа, которые полностью равны NULL, не создаются. То есть индекс не будет содержать записи для таких строк в таблице T, для которых все столбцы в индексе равны NULL. Поскольку столбцу X разрешено иметь значение NULL *временно*, оптимизатор обязан предположить, что X может содержать NULL, и, следовательно, не должен попадать в индекс по X. Поэтому количество, полученное на основе индекса, может (ошибочно) отличаться от количества, полученного из таблицы.

Можно заметить, что если на столбец X помещено неотложенное ограничение, то это ограничение удаляется; т.е. столбец X фактически так же хорош, как столбец Y, *если* ограничение NOT NULL не является отложенным:

```
EODA@ORA12CR1> alter table t drop constraint x_not_null;
```

Table altered.

Таблица изменена.

```
EODA@ORA12CR1> alter table t modify x constraint x_not_null not null;
```

Table altered.

Таблица изменена.

```
EODA@ORA12CR1> explain plan for select count(*) from t;
```

Explained.

Объяснено.

```
EODA@ORA12CR1> select * from table(dbms_xplan.display(null,null,'BASIC'));
```

| Id | Operation | Name |
|----|------------------|-------|
| 0 | SELECT STATEMENT | |
| 1 | SORT AGGREGATE | |
| 2 | INDEX FULL SCAN | T_IDX |

Таким образом, суть в том, что отложенные ограничения должны использоваться только тогда, когда в этом имеется четкая потребность. Они вводят тонкие побочные эффекты, которые могут привести к расхождениям в физической реализации (неуникальные индексы вместо уникальных) или в планах выполнения запросов, как только что было продемонстрировано.

Плохие привычки в отношении транзакций

У многих разработчиков есть плохие привычки, когда речь заходит о транзакциях. Я часто это замечаю у разработчиков, имеющих дело с базами данных, которые “поддерживают”, но не “стимулируют” применение транзакций. Например, в

Informix (по умолчанию), Sybase и SQL Server транзакция должна начинаться явно посредством `BEGIN`; иначе каждый индивидуальный оператор сам по себе будет представлять транзакцию. В манере, подобной способу, которым Oracle помещает изолированные операторы внутрь оболочки `SAVEPOINT`, эти базы данных помещают каждый оператор в пару `BEGIN WORK` и `COMMIT` или `ROLLBACK`. Причина в том, что в упомянутых базах данных блокировки являются ценным ресурсом и процессы чтения блокируют процессы записи, а также наоборот. В попытках увеличить степень параллелизма эти базы данных способствуют тому, что транзакции будут делаться насколько возможно короткими — иногда за счет утери целостности данных.

В Oracle принят противоположный подход. Транзакции всегда неявны, и не существует способа включить “автоматическую фиксацию”, если только приложение не реализует ее (детали приведены в разделе “Использование автоматической фиксации” далее в главе). В Oracle каждая транзакция должна быть зафиксирована, когда это необходимо, но не раньше. Транзакции должны быть настолько крупными, насколько это нужно. Аспекты наподобие блокировок, блокирования и тому подобного не должны рассматриваться в качестве движущих сил, диктующих размер транзакции; целостность данных — вот *главный фактор*, который влияет на размер транзакции. Блокировки не являются дефицитным ресурсом, и отсутствуют проблемы конкуренции между параллельными процессами чтения и записи данных. Это позволяет иметь надежные транзакции в базе данных. Такие транзакции не обязаны быть короткими по продолжительности — они могут быть настолько длинными, насколько это необходимо (но не больше). Транзакции не задумывались как удобство для компьютера и его программного обеспечения; они предназначены для защиты ваших данных.

Фиксация в цикле

Столкнувшись с задачей обновления многих строк, большинство программистов пытаются найти какой-то процедурный способ решить ее в цикле, чтобы можно было фиксировать каждую из большого количества строк. Я слышал два (ошибочных!) довода в пользу такого способа:

- быстрее и эффективнее часто фиксировать множество мелких транзакций, чем обрабатывать и фиксировать одну крупную транзакцию;
- нет достаточного пространства для отката.

Оба довода ошибочны. Более того, слишком частая фиксация подвергает базу данных опасности остаться в “неопределенном” состоянии, если обновление откажет где-то на полпути. Реализовать процесс, который будет гладко перезапускаться в случае отказа, требует написания сложной логики. Намного лучше выполнять фиксацию с частотой, продиктованной нуждами бизнес-процесса, и соответствующим образом настроить размеры сегментов отката.

Давайте рассмотрим эти проблемы более подробно.

Влияние на производительность

Частая фиксация обычно получается не быстрее — почти всегда быстрее делать работу в единственном операторе SQL. В качестве небольшого примера предположим, что есть таблица `T` с большим количеством строк, и необходимо обновить зна-

чение столбца для каждой строки в этой таблице. Мы будем устанавливать такую таблицу следующим образом (запускайте эти четыре оператора перед каждым из трех описанных ниже случаев):

```
EODA@ORA12CR1> drop table t;
Table dropped.
Таблица удалена.

EODA@ORA12CR1> create table t as select * from all_objects;
Table created.
Таблица создана.

EODA@ORA12CR1> exec dbms_stats.gather_table_stats( user, 'T' );
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.

EODA@ORA12CR1> variable n number
```

Обновление может быть просто реализовано с помощью единственного оператора UPDATE:

```
EODA@ORA12CR1> exec :n := dbms_utility.get_cpu_time;
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.

EODA@ORA12CR1> update t set object_name = lower(object_name);
72614 rows updated.
72614 строк обновлено.

EODA@ORA12CR1> exec dbms_output.put_line((dbms_utility.get_cpu_time-:n)||
' cpu hsecs...' );
49 cpu hsecs...
```

Многие люди по разным причинам вынуждены делать это замедленно, строка за строкой, фиксируя каждые N записей:

```
EODA@ORA12CR1> exec :n := dbms_utility.get_cpu_time;
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.

EODA@ORA12CR1> begin
2   for x in ( select rowid rid, object_name, rownum r
3               from t )
4   loop
5       update t
6           set object_name = lower(x.object_name)
7           where rowid = x.rid;
8       if ( mod(x.r,100) = 0 ) then
9           commit;
10      end if;
11  end loop;
12  commit;
13 end;
14 /
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.

EODA@ORA12CR1> exec dbms_output.put_line((dbms_utility.get_cpu_time-:n)||
' cpu hsecs...' );
275 cpu hsecs...
```

В этом простом примере видно, что многократная фиксация в цикле во много раз *медленнее*. Если вы можете сделать что-то в *единственном* операторе SQL, то так и поступайте, поскольку это почти наверняка окажется быстрее. Даже если мы “оптимизируем” процедурный код, применив групповую обработку обновлений (как показано ниже), то все действительно начнет выполняться быстрее, но по-прежнему намного медленнее, чем могло бы быть:

```
EODA@ORA12CR1> exec :n := dbms_utility.get_cpu_time;
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.
```

```
EODA@ORA12CR1> declare
2     type ridArray is table of rowid;
3     type vcArray is table of t.object_name%type;
4
5     l_rids ridArray;
6     l_names vcArray;
7
8     cursor c is select rowid, object_name from t;
9 begin
10    open c;
11    loop
12        fetch c bulk collect into l_rids, l_names LIMIT 100;
13        forall i in 1 .. l_rids.count
14            update t
15                set object_name = lower(l_names(i))
16                where rowid = l_rids(i);
17        commit;
18        exit when c%notfound;
19    end loop;
20    close c;
21 end;
22 /
```

```
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.
```

```
EODA@ORA12CR1> exec dbms_output.put_line((dbms_utility.get_cpu_time-n)||
' cpu hsecs... ');
67 cpu hsecs...
```

```
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.
```

Вдобавок вы должны заметить, что код становится все более и более сложным. Двигаясь от абсолютной простоты единственного оператора UPDATE к процедурному коду и затем к еще более сложному процедурному коду, мы явно избрали неверное направление! Более того (да, тут есть на что пожаловаться), предыдущий процедурный код еще не готов. Он не учитывает ситуацию, когда произойдет отказ (не *если* произойдет, а *когда*). Что случится, если код будет выполнен наполовину и затем система аварийно завершится? Каким образом вы перезапустите процедурный код с фиксацией? Вам придется предусмотреть дополнительный код, чтобы выяснить, с какого места продолжить обработку. При наличии единственного оператора UPDATE мы просто перезапускаем его. Мы знаем, что он либо полностью выполнится успешно, либо потерпит неудачу; здесь не придется беспокоиться о том, как поступать

в отношении частично проделанной работы. Мы еще вернемся к этому вопросу в разделе “Перезапускаемые процессы требуют сложной логики” далее в главе.

Теперь, просто чтобы поддержать противоположную точку зрения в этом обсуждении, вспомните главу 7, в которой шла речь о концепции согласованности записи и о том, как можно было бы обеспечить автоматический перезапуск, к примеру, оператора UPDATE. В случае, когда предшествующий оператор UPDATE должен был выполняться в отношении подмножества строк (он имел конструкцию WHERE, и другие пользователи модифицировали столбцы, которые присутствовали в конструкции WHERE этого оператора UPDATE), тогда имело смысл либо использовать последовательность меньших транзакций вместо одной большой транзакции, либо блокировать таблицу перед выполнением массового обновления. Цель заключалась в сокращении возможности перезапуска.

Если мы собирались обновить значительную часть строк таблицы, то это приводило к необходимости применения команды LOCK TABLE. Однако согласно моему опыту подобного рода массовые обновления или массовые удаления (только операторы этих типов в действительности подлежат перезапуску) делаются в изоляции. Такое крупное, однократное пакетное обновление или удаление старых данных обычно не выполняется в период высокой активности. На самом деле это вообще не должно касаться очистки данных, т.к. для нахождения информации, подлежащей очистке, вы обычно будете использовать какое-нибудь поле с датой, а другие приложения обычно эти данные не модифицируют.

Ошибка “snapshot too old”

Давайте теперь рассмотрим вторую причину возникновения у разработчиков соблазна фиксировать обновления в процедурном цикле, которая является результатом (ошибочных) попыток экономно расходовать “ограниченный ресурс” (сегменты отката). Это проблема конфигурации; вы *должны* обеспечить наличие достаточного пространства для сегментов отмены, чтобы корректно устанавливать размеры своих транзакций. Фиксация в цикле, помимо того, что работает медленнее, также часто вызывает ошибку ORA-01555. Взглянем на все это более детально.

После изучения глав 1 и 7 вы уже знаете, что многоверсионная модель Oracle применяет сегменты отмены для реконструирования блоков в том виде, в каком они были на момент запуска вашего оператора или транзакции (в зависимости от режима изоляции). Если необходимая информация отмены больше не существует, вы получите сообщение об ошибке ORA-01555: snapshot too old (ORA-01555: устаревший снимок) и запрос выполняться не будет. Таким образом, в случае модификации таблицы, которая также читается (как в предыдущем примере), происходит генерация данных отмены, необходимых для запроса. Оператор UPDATE генерирует информацию отмены, которую ваш запрос, вероятно, будет использовать, чтобы получить согласованное по чтению представление данных, подлежащих обновлению. Выполняя фиксацию, вы позволяете системе повторно использовать пространство сегментов отмены, которое вы только что заполнили. И если она задействует это пространство, уничтожив старые данные отмены, которые впоследствии понадобятся вашему запросу, то возникнет серьезная проблема. Оператор SELECT завершится неудачей, а UPDATE остановится на полпути. Вы получаете частично завершенную логическую транзакцию и, скорее всего, не располагаете подходящим способом ее перезапуска (подробности будут приведены ниже).

Давайте посмотрим на эту концепцию в действии с помощью демонстрации. В небольшой тестовой базе данных я подготовил такую таблицу:

```
EODA@ORA12CR1> create table t as select * from all_objects;
Table created.
Таблица создана.
EODA@ORA12CR1> create index t_idx on t(object_name);
Index created.
Индекс создан.
EODA@ORA12CR1> exec dbms_stats.gather_table_stats( user, 'T', cascade=>true );
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.
```

Затем я создал очень маленькое табличное пространство отмены и заставил систему применять его. Обратите внимание, что за счет установки параметра AUTOEXTEND в OFF я ограничил размер всего пространства отмены в этой системе до 10 Мбайт или меньше:

```
EODA@ORA12CR1> create undo tablespace undo_small
2   datafile '/u01/dbfile/ORA12CR1/undo_small.dbf'
3   size 10m reuse
4   autoextend off
5   /
Tablespace created.
Табличное пространство создано.
EODA@ORA12CR1> alter system set undo_tablespace = undo_small;
System altered.
Система изменена.
```

Теперь, имея одно маленькое табличное пространство отмены, я запустил следующий блок кода для обновления:

```
EODA@ORA12CR1> begin
2   for x in ( select /*+ INDEX(t t_idx) */ rowid rid, object_name, rownum r
3             from t
4             where object_name > ' ' )
5   loop
6       update t
7           set object_name = lower(x.object_name)
8           where rowid = x.rid;
9       if ( mod(x.r,100) = 0 ) then
10          commit;
11      end if;
12  end loop;
13  commit;
14  end;
15  /
```

begin

*

ERROR at line 1:

ORA-01555: snapshot too old: rollback segment number with name "" too small

ORA-06512: at line 2

ОШИБКА в строке 1:

ORA-01555: устаревший снимок: номер сегмента отката с именем "" слишком мал

ORA-06512: в строке 2

Возникла ошибка. Я должен обратить внимание на добавление к запросу индексной подсказки и конструкции WHERE, чтобы обеспечить чтение таблицы в случайном порядке (вместе они заставляют оптимизатор по стоимости читать таблицу, “отсортированную” по ключу индекса). Когда мы обрабатываем таблицу через индекс, то стремимся читать блок для одиночной строки, и следующая необходимая строка будет находиться в другом блоке. В конечном итоге мы обработаем все строки из блока 1, просто не одновременно. Блок 1 может содержать, скажем, данные всех строк, у которых значение OBJECT_NAME начинается с букв A, M, N, Q и Z. Поэтому мы будем попадать в блок 1 много раз, т.к. читаем данные, отсортированные по OBJECT_NAME, и предположительно немало значений OBJECT_NAME будут начинаться с букв от A до M. Поскольку мы часто производим фиксацию и повторно используем пространство отката, то со временем повторно наведем блок, где просто больше невозможно выполнить откат к точке во времени, когда начался наш запрос, и в данный момент получим ошибку.

Это был очень надуманный пример, предназначенный только для того, чтобы наглядно продемонстрировать происходящее. Оператор UPDATE генерировал информацию отката. У меня было очень маленькое табличное пространство отката (всего 10 Мбайт). Я многократно переписывал сегменты отката, т.к. они применяются в циклической манере. При каждой фиксации я позволял Oracle перезаписывать сгенерированные данные отката. В конце концов, мне понадобилась порция данных, сгенерированная ранее, но она больше не существовала и возникла ошибка ORA-01555.

Вы вправе указать, что в этом случае, если бы не выполнялась фиксация в строке 10, то была бы получена следующая ошибка:

```
begin
*
ERROR at line 1:
ORA-30036: unable to extend segment by 8 in undo tablespace 'UNDO_SMALL'
ORA-06512: at line 6
ОШИБКА в строке 1:
ORA-30036: не удается расширить сегмент на 8 в табличном пространстве
отмены UNDO_SMALL
ORA-06512: в строке 6
```

Ниже описаны главные отличия между этими двумя ошибками.

- Пример с ORA-01555 оставляет мое обновление в полностью неопределенном состоянии. Часть работы сделана, часть — нет.
- Я абсолютно ничего не могу сделать, чтобы избежать ошибки ORA-01555, если стану выполнять фиксацию в цикле FOR курсора.
- Ошибки ORA-30036 можно избежать, выделив соответствующие ресурсы в системе. Данная ошибка устраняется корректным определением размера сегмента отката; первая ошибка — нет. Более того, даже если я не смогу избежать этой ошибки, то по крайней мере обновление откатывается и база данных остается в известном, согласованном состоянии — я не оставляю ее на полпути какого-то крупного обновления.

Суть здесь в том, что вы не можете “сэкономить” пространство отката за счет частой фиксации — эти данные отката вам нужны.

Я работал в однопользовательской системе, когда получил ошибку ORA-01555. Чтобы привести к этой ошибке, достаточно было всего лишь одного сеанса, и даже в реальности часто оказывается, что какой-то отдельный сеанс вызывает собственные ошибки ORA-01555. Разработчики и администраторы баз данных должны работать совместно, чтобы адекватно определить размер этих сегментов для задач, решаемых в системе. Здесь не может быть кратковременных изменений. Тщательно анализируя свою систему, вы должны выяснить, каковы самые большие транзакции, и в соответствии с ними установить размеры сегментов отката. Динамическое представление производительности V\$UNDOSTAT может быть очень полезным для отслеживания объема генерируемых данных отката и продолжительности наиболее длительно выполняющихся запросов. Многие люди склонны воспринимать вещи вроде временной области, области отмены и области повторения как накладные расходы, т.е. то, подо что выделять пространство минимально возможного объема. Это напоминает проблему, с которой компьютерная индустрия столкнулась 1 января 2000 года — она была вызвана стремлением сэкономить 2 байта в поле даты. Упомянутые компоненты базы данных не являются накладными расходами, а ключевыми компонентами системы. Их размеры должны определяться соответствующим образом (не слишком большими, но и не слишком малыми).

На заметку! Поэкспериментировав со слишком маленькими сегментами отмены, не забудьте после выполнения этих примеров вернуть табличное пространство отмены в обычное состояние, иначе вы будете сталкиваться с ошибками ORA-30036 до конца книги!

Перезапускаемые процессы требуют сложной логики

Самая серьезная проблема подхода с фиксацией перед завершением логической транзакции связана с тем фактом, что он часто оставляет базу данных в неизвестном состоянии, если оператор UPDATE терпит отказ где-то на середине пути. Если вы не подготовились к этому заблаговременно, то будет очень трудно перезапустить отказавший процесс, позволив ему восстановиться в месте, где он был прерван. Например, предположим, что мы применяем к столбцу не функцию LOWER(), как в предыдущем примере, а какую-то другую функцию столбца:

```
last_ddl_time = last_ddl_time + 1;
```

Если мы остановим цикл UPDATE на полпути, то каким образом его перезапустить? Мы не можем просто выполнить его повторно, потому что тогда к некоторым датам будет добавлено 2, а к остальным — 1. Если снова произойдет отказ, то при следующем запуске к одним датам добавится 3, к другим — 2, к остальным — 1 и т.д. Мы нуждаемся в более сложной логике — некотором способе “секционирования” данных. Например, можно было бы обрабатывать данные со значением OBJECT_NAME, начинающимся с A, затем — с B и т.д.:

```
EODA@ORA12CR1> create table to_do
2 as
3 select distinct substr( object_name, 1,1 ) first_char
4   from T
5   /
```

Table created.

Таблица создана.

```

EODA@ORA12CR1> begin
2      for x in ( select * from to_do )
3      loop
4          update t set last_ddl_time = last_ddl_time+1
5              where object_name like x.first_char || '%';
6
7          dbms_output.put_line( sql%rowcount || ' rows updated' );
8          delete from to_do where first_char = x.first_char;
9
10         commit;
11     end loop;
12 end;
13 /
238 rows updated.
238 строк обновлено.
5730 rows updated.
5730 строк обновлено.
1428 rows updated.
1428 строк обновлено.
...
262 rows updated.
262 строки обновлено.
1687 rows updated.
1687 строк обновлено.
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.

```

Теперь можно перезапускать процесс, если он завершится сбоем, т.к. мы не будем обрабатывать данные OBJECT_NAME, которые уже были обработаны успешно. Тем не менее, этому подходу присуща одна проблема: если нет атрибута, который позволит равномерно распределить данные, то мы получим очень широкое распределение строк. Второму оператору UPDATE придется выполнять намного больше работы, чем всем остальным вместе взятым. Вдобавок если другие сеансы получают доступ к этой же таблице и модифицируют данные, они могут обновить также и поле OBJECT_NAME. Предположим, что какой-то другой сеанс обновляет объект по имени Z, изменяя его на A, *после* того, как мы уже обработали все объекты с именами A. Мы потеряем эту запись. Более того, это очень неэффективный процесс по сравнению с UPDATE T SET LAST_DDL_TIME = LAST_DDL_TIME+1. Возможно, мы воспользуемся индексом для чтения каждой строки таблицы или выполним полное сканирование *n* раз — и то, и другое нежелательно. О таком подходе можно сказать слишком много плохого.

На заметку! В главе 14 мы рассмотрим функциональное средство, которое доступно, начиная с версии Oracle 11g Release 2 — пакет DBMS_PARALLEL_EXECUTE. Там мы вернемся к этому подходу с перезапуском, а также будем иметь дело с шаблонами неунифицированных обновлений.

Лучшим подходом является тот, который я отстаивал в начале главы 1: делайте все просто. Если что-то можно реализовать на SQL — так и поступайте. Если это сделать на SQL невозможно, реализуйте его на PL/SQL. Старайтесь использовать

минимально возможный объем кода. Выделяйте достаточные ресурсы. Всегда думайте о том, что произойдет в случае ошибки. Мне очень часто приходилось видеть людей, кодирующих циклы обновления, которые великолепно работали с тестовыми данными, но отказывали на полпути, когда применялись к реальным данным. Затем они действительно попадали в тупик, поскольку не имели ни малейшего понятия о том, где цикл останавливал обработку. Намного проще корректно установить размер области отката, чем создавать перезапускаемую программу. При наличии действительно больших таблиц, подлежащих обновлению, вы должны использовать секции (более подробно обсуждаются в главе 10), которые можно обновлять индивидуально. Для выполнения обновлений вы можете даже применять параллельный DML или пакет DBMS_PARALLEL_EXECUTE в Oracle 11g Release 2 и последующих версиях.

Использование автоматической фиксации

Заключительные слова по поводу плохих привычек относительно транзакций касаются использования популярных API-интерфейсов ODBC и JDBC. Эти API-интерфейсы выполняют автоматическую фиксацию по умолчанию. Рассмотрим следующие операторы, которые переводят сумму \$1000 с расчетного счета на сберегательный счет:

```
update accounts set balance = balance - 1000 where account_id = 123;
update accounts set balance = balance + 1000 where account_id = 456;
```

Если при отправке этих запросов программа применяет API-интерфейс ODBC или JDBC, то он (молча) внедряет оператор фиксации после *каждого* оператора UPDATE. Представьте последствия этого, если в системе произойдет отказ после первого UPDATE, но перед вторым. Сумма \$1000 будет потеряна!

Отчасти я могу понять, почему ODBC так поступает. Дело в том, что ODBC проектировали разработчики SQL Server, а эта база данных требует использования очень коротких транзакций из-за своей модели параллелизма (операции записи блокируют операции чтения, операции чтения блокируют операции записи, а блокировки являются дефицитным ресурсом). Чего я не могу понять, так это почему подобное попало в JDBC — в API-интерфейс, предполагающий поддержку уровня “предприятия”. Я убежден в том, что следующая строка кода после открытия подключения в JDBC всегда должна выглядеть так:

```
Connection conn = DriverManager.getConnection
    ("jdbc:oracle:oci:@database", "scott", "tiger");
conn.setAutoCommit (false);
```

Это возвращает контроль над транзакциями обратно вам, разработчику, которому он и должен принадлежать. Затем можно безопасно кодировать транзакцию перевода денег между счетами и фиксировать ее после того, как оба оператора завершатся успешно. В данном случае недостаток знаний API-интерфейса может оказаться губительным. Мне приходилось встречать немало разработчиков, не подозревающих об этой “функции” автоматической фиксации, которые получали огромные проблемы со своими приложениями, когда возникала ошибка.

Распределенные транзакции

Одной из действительно замечательных характеристик СУБД Oracle является ее способность прозрачно обрабатывать распределенные транзакции. Внутри одной транзакции я могу обновлять данные во многих базах данных. Когда я выполняю фиксацию этой транзакции, то либо фиксируются все обновления во всех экземплярах, либо не фиксируется ни одно из них (все они откатываются). Для достижения такой цели не нужен никакой дополнительный код; я просто записываю `commit`.

Ключом к распределенным транзакциям в Oracle является *связь баз данных* (*database link*), которая представляет собой объект базы данных, описывающий способ входа в другой экземпляр базы из вашего экземпляра. Однако цель этого раздела не в том, чтобы раскрыть синтаксис команд связи баз данных (они полностью документированы в руководстве по языку SQL (*Oracle Database SQL Language Reference*)), а скорее в том, чтобы уведомить о самом его существовании. Как только вы настроите связь баз данных, доступ к удаленным объектам становится очень простым:

```
select * from T@another_database;
```

Этот запрос выберет данные из таблицы `T` в экземпляре базы данных, определенном связью `ANOTHER_DATABASE`. Обычно вы будете “скрывать” тот факт, что `T` — удаленная таблица, создавая ее представление, или синоним. Например, можно выполнить следующую команду и затем обращаться к `T`, как если бы она была локальной таблицей:

```
create synonym T for T@another_database;
```

Теперь, когда есть настроенная связь баз данных и возможность чтения некоторых таблиц, эти таблицы можно также модифицировать (конечно, при наличии подходящих привилегий). Выполнение распределенной транзакции ничем не отличается от запуска локальной транзакции. Вот все, что для этого нужно:

```
update local_table set x = 5;
update remote_table@another_database set y = 10;
commit;
```

И ничего больше. СУБД Oracle выполнит фиксацию либо в обеих базах данных, либо ни в одной из них. Для этого Oracle применяет протокол двухфазной фиксации (*Two Phase Commit — 2PC*) — распределенный протокол, который позволяет команде модификации, затрагивающей несколько разрозненных баз данных, фиксироваться атомарным образом. Перед фиксацией протокол 2PC пытается насколько возможно перекрыть все пути для распределенного отказа. При двухфазной фиксации между многими базами данных одна из баз — обычно та, в которую клиент вошел первоначально — станет координатором распределенной транзакции. Этот узел будет опрашивать другие узлы на предмет готовности к фиксации. По существу данный узел будет обращаться к другим узлам, предлагая им подготовиться к фиксации. Каждый из других узлов сообщит свое “состояние готовности” в форме “да” или “нет”. Если любой из узлов ответит “нет”, будет произведен откат всей транзакции. Если все узлы ответили “да”, то узел-координатор рассылает сообщение с командой фиксации на каждом из узлов.

Это сужает окно, в котором могла бы возникнуть серьезная ошибка. Перед “голосованием” по протоколу 2PC любая распределенная ошибка приведет к выпол-

нению отката на всех узлах. Не будет никаких сомнений относительно исхода распределенной транзакции. После выдачи команды на фиксацию или откат снова нет никаких сомнений об исходе распределенной транзакции. Только в течение очень короткого временного окна, когда координатор собирает ответы от узлов, исход может быть неоднозначным из-за отказа.

Предположим, например, что в транзакции участвуют три узла, причем узел 1 выступает в роли координатора. Узел 1 просит узел 2 подготовиться к фиксации, и узел 2 делает это. Затем узел 1 предлагает узлу 3 подготовиться к фиксации, и он также делает это. В данный момент времени узел 1 является единственным, которому известен исход транзакции, и теперь он отвечает за рассылку сведений об этом исходе другим узлам. Если ошибка произойдет прямо сейчас (отказ сети, сбой электропитания на узле 1 или что-нибудь еще), то узлы 2 и 3 останутся в подвешенном состоянии. Они получают то, что называется *сомнительной распределенной транзакцией*. Протокол 2PC пытается минимизировать возможность возникновения ошибок, но не способен исключить ее полностью. Узлы 2 и 3 должны удерживать транзакцию открытой, ожидая от узла 1 уведомления об исходе транзакции.

Если вы вспомните обсуждение архитектуры в главе 5, то там было сказано, что для решения такой проблемы существует процесс RECO. Это также тот случай, когда в игру вступают команды COMMIT и ROLLBACK с опцией FORCE. Если причиной проблемы был отказ сети между узлами 1, 2 и 3, то администраторы баз данных на узлах 2 и 3 должны обратиться к администратору базы данных на узле 1, узнать у него исход транзакции и соответствующим образом применить фиксацию или откат вручную.

Существует несколько (небольшое количество) ограничений относительно того, что можно делать в распределенной транзакции, и все они оправданы (во всяком случае, для меня они выглядят оправданными). Ниже перечислены самые важные из них.

- Вы не можете выдать команду COMMIT через связь баз данных. То есть команда наподобие COMMIT@удаленный_узел не разрешена. Фиксацию допускается делать только на узле, который инициировал транзакцию.
- Вы не можете выполнять операторы DDL через связь баз данных. Это прямое следствие предыдущего ограничения. Операторы DDL осуществляют фиксацию. Поскольку фиксацию нельзя производить ни из каких других узлов, кроме того, который инициировал транзакцию, то не допускается и запуск операторов DDL через связь баз данных.
- Нельзя выдавать команду SAVEPOINT через связь баз данных. Короче говоря, через связь баз данных не разрешено выполнять операторы управления транзакциями. Все управление транзакциями наследуется от сеанса, в котором первоначально открывалась связь баз данных. Другие средства управления транзакцией в распределенных экземплярах транзакции не доступны.

Отсутствие возможности управлять транзакциями через связь баз данных вполне оправдано, т.к. инициирующий узел является единственным, который имеет список всех, кто вовлечен в транзакцию. Если в нашей конфигурации из трех узлов узел 2 попытается зафиксировать транзакцию, то у него не будет способа узнать о том, что в транзакции участвует узел 3. В Oracle только узел 1 может выдать команду фиксации. В этой точке узлу 1 разрешено делегировать ответственность за управление распределенной транзакцией другому узлу.

Мы можем повлиять на то, какой узел будет действительным фиксирующим узлом, устанавливая для узла параметр `COMMIT_POINT_STRENGTH` (приоритет завершения транзакции). Этот параметр ассоциирует относительный уровень важности с сервером в распределенной транзакции. Чем более важен сервер (чем больше доступных данных должно быть на нем), тем более вероятно, что он будет координировать распределенную транзакцию. Поступать подобным образом может потребоваться в случае, если необходимо выполнить распределенную транзакцию между производственной и тестовой машинами. Поскольку координатор транзакции *никогда* не сомневается в исходе транзакции, лучше, если распределенную транзакцию будет координировать производственная машина. Вас не должно особо беспокоить то, что на тестовой машине есть открытые транзакции и заблокированные ресурсы. Вас определенно должна волновать ситуация, если подобное произойдет на производственной машине.

В невозможности выполнения операторов DDL через связь баз данных на самом деле нет ничего плохого. Во-первых, операторы DDL применяются редко. Вы запускаете их один раз при установке или во время модернизации. Производственные системы не выполняют операторы DDL (во всяком случае, *не должны*). Во-вторых, существует метод запуска операторов DDL через связь баз данных до известной степени, предусматривающий использование средства организации очереди заданий `DBMS_JOB` или, в Oracle 10g и последующих версиях, пакета планировщика `DBMS_SCHEDULER`. Вместо попытки выдавать операторы DDL через связь, вы применяете эту связь для планирования удаленного задания, чтобы оно выполнилось, как только будет произведена фиксация. В таком случае задание запустится на удаленной машине, оно не будет частью распределенной транзакции и сможет выполнить операторы DDL. Фактически это метод, посредством которого службы репликации Oracle (Oracle Replication Services) выполняют распределенные команды DDL для проведения репликации схемы.

Автономные транзакции

Автономные транзакции позволяют создать “транзакцию внутри транзакции”, которая зафиксирует или произведет откат изменений независимо от родительской транзакции. Они дают возможность приостановить транзакцию, выполняемую в текущий момент, запустить новую транзакцию, сделать некоторую работу и зафиксировать или откатить ее — и все это не затрагивая состояние текущей выполняемой транзакции. Автономные транзакции предоставляют новый метод управления транзакциями в PL/SQL и могут использоваться в следующих конструкциях:

- анонимные блоки верхнего уровня;
- локальные (процедура в процедуре), автономные или пакетные функции и процедуры;
- методы объектных типов;
- триггеры баз данных.

Прежде чем посмотреть, как работают автономные транзакции, я хотел бы подчеркнуть, что такой тип транзакций является мощным и, следовательно, опасным инструментом, если применять его неправильно. Реальная потребность в автоном-

ных транзакциях возникает очень редко. Я с подозрением отношусь к любому коду, в котором они используются — такой код заслуживает дополнительного изучения. Применяя автономные транзакции, чрезвычайно легко неумышленно внести в систему проблемы с логической целостностью данных. В последующих разделах мы обсудим, когда их безопасно использовать, после того, как ознакомимся с их работой.

Как работают автономные транзакции

Лучший способ продемонстрировать действие и последствия автономных транзакций — рассмотреть пример. Создадим простую таблицу для хранения сообщений:

```
EODA@ORA12CR1> create table t ( msg varchar2(25) );
Table created.
Таблица создана.
```

Далее создадим две процедуры, каждая из которых просто вставляет свое имя в таблицу сообщений и производит фиксацию. Однако одна из этих процедур будет нормальной, а другая — закодированной в виде автономной транзакции. Мы будем применять эти объекты, чтобы показать, какая работа сохраняется (фиксируется) в базе данных при различных обстоятельствах.

Для начала вот процедура AUTONOMOUS_INSERT:

```
EODA@ORA12CR1> create or replace procedure Autonomous_Insert
2 as
3     pragma autonomous_transaction;
4 begin
5     insert into t values ( 'Autonomous Insert' );
6     commit;
7 end;
8 /
Procedure created.
Процедура создана.
```

Обратите внимание на использование прагмы AUTONOMOUS_TRANSACTION. Эта директива сообщает базе данных, что процедура должна выполняться как новая автономная транзакция, не зависящая от ее родительской транзакции.

На заметку! Прагма (pragma) — это просто директива компилятора, т.е. метод инструктирования компилятора о том, что он должен выполнить некоторую опцию компиляции. Доступны и другие прагмы. Их список можно найти в предметном указателе руководства по языку PL/SQL (*Oracle Database PL/SQL Language Reference*).

А вот “нормальная” процедура NONAUTONOMOUS_INSERT:

```
EODA@ORA12CR1> create or replace procedure NonAutonomous_Insert
2 as
3 begin
4     insert into t values ( 'NonAutonomous Insert' );
5     commit;
6 end;
7 /
Procedure created.
Процедура создана.
```

Теперь давайте понаблюдаем за поведением *неавтономной* транзакции в анонимном блоке кода PL/SQL:

```
EODA@ORA12CR1> begin
2      insert into t values ( 'Anonymous Block' );
3      NonAutonomous_Insert;
4      rollback;
5 end;
6 /
```

PL/SQL procedure successfully completed.

Процедура PL/SQL успешно завершена.

```
EODA@ORA12CR1> select * from t;
```

MSG

```
-----
Anonymous Block
NonAutonomous Insert
```

Как видите, работа, выполненная анонимным блоком — его оператор INSERT — была *зафиксирована* процедурой NONAUTONOMOUS_INSERT. Зафиксированы обе строки данных, так что команде ROLLBACK откатывать нечего. Сравните это с поведением процедуры автономной транзакции:

```
EODA@ORA12CR1> delete from t;
```

2 rows deleted.

2 строки удалено.

```
EODA@ORA12CR1> commit;
```

Commit complete.

Фиксация завершена.

```
EODA@ORA12CR1> begin
2      insert into t values ( 'Anonymous Block' );
3      Autonomous_Insert;
4      rollback;
5 end;
6 /
```

PL/SQL procedure successfully completed.

Процедура PL/SQL успешно завершена.

```
EODA@ORA12CR1> select * from t;
```

MSG

```
-----
Autonomous Insert
```

Здесь сохраняется только работа, сделанная и зафиксированная в автономной транзакции. Оператор INSERT, выполненный в анонимном блоке, откатывается оператором ROLLBACK в строке 4. Оператор COMMIT процедуры автономной транзакции не оказывает никакого воздействия на родительскую транзакцию, которая стартовала в анонимном блоке. В действительности здесь продемонстрирована сущность автономных транзакций и то, что они делают.

Подводя итоги, можно сказать, что если выполняется оператор COMMIT внутри “нормальной” процедуры, то он касается не только ее собственной работы, но также всей незавершенной работы, произведенной в данном сеансе. Тем не менее, оператор COMMIT, выполненный в процедуре с автономной транзакцией, распространяется только на работу этой процедуры.

Когда использовать автономные транзакции

Внутренне автономные транзакции в Oracle поддерживаются довольно давно. Мы видим их в форме рекурсивного SQL. Например, рекурсивная транзакция может выполняться при выборе значения из последовательности, чтобы немедленно инкрементировать значение последовательности в таблице SYS.SEQ\$. Обновление таблицы SYS.SEQ\$ для поддержки вашей последовательности сразу же фиксируется и становится видимым другим транзакциям, хотя ваша транзакция пока еще не зафиксирована. Вдобавок, если вы осуществите откат своей транзакции, то результат инкрементирования значения последовательности останется в силе; он не откатывается вместе с транзакцией, т.к. уже зафиксирован. Управление пространством, аудит и другие внутренние операции выполняются в аналогичной рекурсивной манере.

Теперь это средство открыто для применения всеми. Однако я обнаружил, что оправданное использование автономных транзакций в реальности *чрезвычайно ограничено*. Временами я встречаю их применение в качестве обходного пути для решения таких задач, как ограничение мутирующей таблицы в триггере. Однако это почти всегда ведет к проблемам с целостностью данных, поскольку причиной мутирующей таблицы является попытка читать таблицу в то время, когда инициирован ее триггер. Используя автономную транзакцию, вы можете запросить таблицу, но не увидите внесенных вами изменений (для чего в первую очередь и предназначено ограничение мутирующей таблицы; таблица находится в процессе модификации, поэтому результаты запроса были бы несогласованными). Любые решения, принятые на основе запроса из триггера, оказались бы сомнительными — в этот момент времени вы читаете “старые” данные.

Потенциально допустимое применение автономной транзакции связано со специальным аудитом, но я делаю акцент на словах “потенциально допустимое”. Существуют более эффективные способы для аудита информации в базе данных, чем написание специального триггера. Например, можно использовать пакет DBMS_FGA или просто саму команду AUDIT.

Разработчиками приложений часто задают мне вопрос: как можно протоколировать в журнале информацию об ошибках в процедурах PL/SQL, чтобы она сохранялась, даже если работа, выполненная этими процедурами, откатывается? Ранее было показано, что операторы PL/SQL являются *атомарными* — они либо полностью успешны, либо полностью неудачны. Если сведения об ошибке протоколируются в процедурах PL/SQL, то по умолчанию журнальная информация об ошибке откатывается, когда Oracle производит откат оператора. Автономные транзакции позволяют изменить это поведение, сохраняя информацию об ошибке в журнале даже в случае отката остальной частичной работы процедуры.

Давайте начнем с подготовки простой таблицы для протоколирования ошибок; мы будем записывать в нее отметку времени возникновения ошибки, сообщение об ошибке и стек ошибок PL/SQL (для идентификации источника ошибки):

```

ODA@ORA12CR1> create table error_log
2 ( ts timestamp,
3   err1 clob,
4   err2 clob )
5 /
Table created.
Таблица создана.

```

Теперь нам нужна процедура PL/SQL для протоколирования ошибок в этой таблице. Воспользуемся следующей небольшой процедурой:

```

EODA@ORA12CR1> create or replace
  2 procedure log_error
  3 ( p_err1 in varchar2, p_err2 in varchar2 )
  4 as
  5     pragma autonomous_transaction;
  6 begin
  7     insert into error_log( ts, err1, err2 )
  8     values ( systimestamp, p_err1, p_err2 );
  9     commit;
 10 end;
 11 /
Procedure created.
Процедура создана.

```

Вся “магия” этой процедуры сосредоточена в строке 5, где применяется директива `pragma autonomous_transaction`, которая информирует PL/SQL о том, что данная процедура должна запустить новую транзакцию, сделать в ней определенную работу и зафиксировать ее, не влияя на любую другую транзакцию, выполняющуюся в текущий момент. Оператор `COMMIT` в строке 9 может затрагивать только SQL-код, выполняемый процедурой `LOG_ERROR`.

Теперь давайте проведем тестирование. Ради интереса создадим пару процедур, которые будут вызывать друг друга:

```

EODA@ORA12CR1> create table t ( x int check (x>0) );
Table created.
Таблица создана.

EODA@ORA12CR1> create or replace procedure p1( p_n in number )
  2 as
  3 begin
  4     -- некоторый код
  5     insert into t (x) values ( p_n );
  6 end;
  7 /
Procedure created.
Процедура создана.

EODA@ORA12CR1> create or replace procedure p2( p_n in number )
  2 as
  3 begin
  4     -- код
  5     -- код
  6     p1(p_n);
  7 end;
  8 /
Procedure created.
Процедура создана.

```

Теперь вызовем эти процедуры из анонимного блока:

```

EODA@ORA12CR1> begin
2      p2( 1 );
3      p2( 2 );
4      p2( -1 );
5  exception
6      when others
7      then
8          log_error( sqlerrm, dbms_utility.format_error_backtrace );
9      RAISE;
10 end;
11 /
begin
*
ERROR at line 1:
ORA-02290: check constraint (EODA.SYS_C0061527) violated
ORA-06512: at line 9
ОШИБКА в строке 1:
ORA-02290: нарушено проверочное ограничение целостности (EODA.SYS_C0061527)
ORA-06512: в строке 9

```

Мы видим, что код завершился неудачей (эту ошибку необходимо возратить, для чего предусмотрен оператор RAISE в строке 9). Можно удостовериться в том, что база данных Oracle отменила сделанную работу (мы знаем, что первые два вызова процедуры P2 завершились успешно; значения 1 и 2 благополучно вставлены в таблицу T):

```

EODA@ORA12CR1> select * from t;
no rows selected
нет выбранных строк

```

Но можно также убедиться, что информация в журнале ошибок сохранена и в действительности зафиксирована:

```

EODA@ORA12CR1> rollback;
Rollback complete.
Откат завершен.

EODA@ORA12CR1> select * from error_log;
TS
-----
ERR1
-----
ERR2
-----
09-MAY-14 05.11.03.818918 PM
ORA-02290: check constraint (EODA.SYS_C00204351) violated
ORA-06512: at "EODA.P1", line 5
ORA-06512: at "EODA.P2", line 6
ORA-06512: at line 4
ORA-02290: нарушено проверочное ограничение целостности (EODA.SYS_C00204351)
ORA-06512: в EODA.P1, строка 5
ORA-06512: в EODA.P2, строка 6
ORA-06512: в строке 4

```

Согласно моему опыту, единственным по-настоящему допустимым использованием автономных транзакций является протоколирование ошибок или информационных сообщений в манере, позволяющей их фиксировать независимо от родительской транзакции.

Резюме

В настоящей главе были рассмотрены многие аспекты управления транзакциями в Oracle. Транзакции — одно из главных средств, отличающих базу данных от файловой системы. Понимание их работы и правильное их применение необходимо для корректной реализации приложений в любой базе данных. Очень важно знать, что в Oracle любой оператор является атомарным (включая все его побочные эффекты), и что атомарность распространяется на хранимые процедуры. Мы видели, как помещение в блок PL/SQL обработчика исключений WHEN OTHERS может радикально повлиять на то, какие изменения произойдут в базе данных. Для разработчиков базы данных четкое понимание работы транзакций имеет важнейшее значение.

Мы взглянули на отчасти сложные взаимодействия между ограничениями целостности (уникальные ключи, проверочные ограничения и т.п.) и транзакциями в Oracle. Мы обсудили, каким образом Oracle обычно обрабатывает ограничения целостности немедленно после выполнения операторов, а затем показали, как при желании можно отложить проверку ограничений до конца транзакции. Это средство является ключевым в реализации сложных многотабличных обновлений, когда модифицируемые таблицы зависят друг от друга — примером может служить каскадное обновление.

Мы проанализировали ряд вредных привычек, связанных с транзакциями, которые вырабатываются у разработчиков из-за эксплуатации баз данных, “поддерживающих”, а не “стимулирующих” использование транзакций. Было описано главное правило транзакций: они должны быть насколько возможно короткими, но не короче, чем необходимо. *Размером транзакций управляет целостность данных* — вот ключевая концепция, которую вы должны вынести из этой главы. Единственный фактор, влияющий на размер транзакций — это бизнес-правила, которым подчиняется ваша система. Не пространство отката, не блокировки, а именно бизнес-правила.

Мы раскрыли тему распределенных транзакций и их отличия от транзакций в единственной базе данных. Мы исследовали ограничения, накладываемые на распределенные транзакции, и объяснили причины их существования. Прежде чем приступить к построению распределенной системы, следует четко уяснить эти ограничения. То, что работает в одиночном экземпляре, может не работать в распределенной базе данных.

Глава была завершена рассмотрением автономных транзакций и описанием того, что они собой представляют, а также, что более важно — когда их стоит применять, а когда нет. Я хотел бы еще раз подчеркнуть, что оправданное использование автономных транзакций в реальном мире встречается исключительно редко. Если вы обнаружите, что применяете их постоянно, стоит серьезно задуматься над причинами.

Повтор и отмена

В этой главе описаны две из наиболее важных частей в базе данных Oracle: повторение (redo) и отмена (undo). *Redo* — это информация, которую Oracle записывает в оперативные (и архивные) журнальные файлы повторения (redo-журналы) на случай “повторного воспроизведения” транзакции после сбоя. *Undo* — это информация, которую Oracle сохраняет в сегментах отмены с целью аннулирования, или отката, транзакции.

Здесь мы обсудим такие темы, как генерация redo- и undo-журналов (rollback-журналов), а также их участие в транзакциях, восстановлении и тому подобных действиях. Мы начнем с высокоуровневого обзора того, что собой представляют redo- и undo-журналы, и каким образом они работают вместе. Затем мы углубимся в каждую из этих тем, раскрывая их во всех подробностях, и рассмотрим то, что необходимо знать о них разработчику.

Эта глава ориентирована в первую очередь на точку зрения разработчика — мы не будем рассматривать проблемы, за выявление и решение которых должен отвечать исключительно администратор базы данных. Например, мы не показываем, как найти оптимальные значения для параметров `RECOVERY_PARALLELISM` или `FAST_START_MTTR_TARGET`. Тем не менее, журналы redo и undo — это тот предмет, который связывает роли разработчика и администратора базы данных. Они оба должны обладать фундаментальным пониманием назначения журналов redo и undo, особенностей их функционирования и способов избежать потенциальных сложностей, связанных с их применением. Знание журналов undo и redo также поможет администраторам баз данных и разработчикам лучше понять работу базы данных в целом.

В настоящей главе я представлю псевдокод для этих механизмов Oracle и концептуальные пояснения того, что в действительности происходит. Я не буду приводить все внутренние детали о том, какими байтами данных обновляются те или иные файлы. То, что на самом деле происходит — немного сложнее, но наличие хорошего понимания рабочего потока весьма ценно и поможет в оценке последствий ваших действий.

На заметку! Я то и дело получаю вопросы о тех или иных битах и байтах в журналах повторения и отмены. Похоже, люди хотят иметь очень детальную спецификацию того, что в точности в них находится. Я никогда не отвечаю на такие вопросы. Вместо этого я направляю их внимание на предназначение повторения и отмены, на концепции, лежащие в их основе. Я концентрируюсь на использовании повторения и отмены — не на битах и байтах. Я и сам не разбираю подробно файлы журналов повторения или сегменты отмены.

На заметку! Я пользуюсь предлагаемыми инструментами, такими как средством Log Miner для чтения redo-информации и хронологией ретроспективных транзакций для чтения undo-информации, но эти инструменты представляют информацию в формате, читабельном для человека. Поэтому в настоящей главе мы не станем углубляться в их внутренний состав, а вместо этого постараемся заложить прочный фундамент знаний.

Что собой представляет redo

Журнальные файлы повторения (redo-журналы) критически важны для базы данных Oracle. Это журналы транзакций, выполненных в базе данных. СУБД Oracle поддерживает два типа redo-журналов: *оперативные* и *архивные*. Они применяются в целях восстановления: их главное назначение — использоваться в случае выхода из строя экземпляра или носителя.

Если на сервере базы данных пропадает электропитание, что приводит к отказу работы экземпляра, то Oracle будет применять оперативные журналы redo для восстановления системы в точности до точки фиксации, которая непосредственно предшествовала отключению электропитания. Если происходит поломка дискового привода (выход из строя носителя), то Oracle будет использовать архивные и оперативные журналы redo для восстановления из резервной копии данных, которые находились на этом приводе, на подходящий момент времени. Более того, если вы “нечаянно” очистите таблицу либо удалите какую-то важную информацию и зафиксируете эту операцию, то сможете восстановить эти данные из резервной копии в состоянии на момент, непосредственно предшествовавший “несчастному случаю”, с применением оперативных и архивных журнальных файлов redo.

Архивные журнальные файлы redo — это просто копии старых, заполненных оперативных журнальных файлов redo. По мере того, как система наполняет журнальные файлы, процесс ARCn делает копию оперативного журнального файла redo в другом месте, а также дополнительно помещает несколько других копий в локальные и удаленные местоположения. Эти архивные файлы журналов redo используются при выполнении восстановления носителя, когда отказ вызван неполадками дискового устройства или другим физическим дефектом. СУБД Oracle может взять эти архивные файлы журналов redo и применить их к резервным копиям файлов данных, чтобы наверстать остальную часть базы данных. Они являются хронологией транзакций базы данных.

На заметку! С выходом версии Oracle 10g стала доступной технология ретроспективы (flash-back). Она позволяет выполнять ретроспективные запросы (запросы данных на определенный момент времени в прошлом), восстанавливать удаленные таблицы базы, помещать таблицы обратно в состояние, которое они имели некоторое время назад, и т.д. В результате количество случаев, когда может понадобиться традиционное восстановление с использованием резервных копий и архивных журналов redo, существенно снизилось. Тем не менее, возможность выполнить восстановление является наиболее важной работой администратора базы данных. Восстановление базы данных — это то, в чем администратору ошибаться нельзя.

Каждая база данных Oracle имеет, по крайней мере, две группы оперативных журналов redo, в каждой из которых присутствует минимум один член (журнальный файл). Запись в эти оперативные группы журналов производится в циклической манере. СУБД Oracle выполняет запись в файлы журналов группы 1, а когда заканчиваются файлы в группе 1, переключается на файлы журналов группы 2 и начинает записывать в них. После заполнения журнальных файлов группы 2 она снова переключается на файлы группы 1 (предполагается, что есть только две группы файлов журналов redo; если их, скажем, три, естественно, Oracle продолжит обработкой третьей группы).

Журналы redo, или журналы транзакций, представляют собой одно из важнейших средств, которые делают базу данных именно базой данных. Вероятно, они являются наиболее важной структурой восстановления, хотя без других частей, таких как сегменты отмены, восстановление распределенных транзакций и т.п., ничего работать не будет. Они являются главным компонентом, отличающим базы данных от файловой системы. Оперативные журналы redo позволяют эффективно восстанавливать базу данных после перебоев в электропитании, которые могут произойти, когда Oracle находится в процессе записи. Архивные журналы redo позволяют проводить восстановление после отказа носителей, когда, например, жесткий диск выходит из строя или вследствие человеческой ошибки утрачиваются данные. Без журналов redo база данных не смогла бы предложить более высокую защиту, чем файловая система.

Что собой представляет undo

Концептуально undo является противоположностью redo. Информация undo регистрируется базой данных во время модификации данных и позволяет вернуть данные в состояние, которое они имели до проведения модификации. Это может делаться для поддержки многоверсионности, как было показано в главе 7, в случае отказа выполняемой транзакции или оператора по любой причине либо при запросе посредством оператора ROLLBACK. В то время как redo применяется для повторения транзакции в случае отказа, чтобы восстановить транзакцию, информация undo ориентирована на аннулирование результатов выполнения оператора или набора операторов. Информация undo, в отличие от redo, сохраняется внутри базы данных в специальных сегментах, называемых сегментами отмены.

На заметку! Сегмент отката (rollback segment) и сегмент отмены (undo segment) считаются терминами-синонимами. Используя ручное управление пространством отмены, администратор базы данных создает сегменты отката. Применяя автоматическое управление пространством отмены, система автоматически создает и уничтожает сегменты отмены по мере надобности. В рамках нашего обсуждения эти термины следует трактовать как идентичные.

Существует распространенное заблуждение, что информация undo используется для *физического* восстановления базы данных в состояние, которое она имела перед выполнением оператора или транзакции, хотя это не так. База данных *логически* восстанавливается в предыдущее состояние — любые изменения логически отменяются, — но структуры данных, сами блоки базы данных, после отката вполне могут

отличаться. Причина в том, что в любой многопользовательской системе происходят десятки, сотни или даже тысячи параллельных транзакций. Одной из главных функций базы данных является посредничество в параллельном доступе к своим данным. Блоки, модифицированные нашей транзакцией, также в общем случае модифицируются и многими другими транзакциями. Следовательно, мы не можем просто поместить блок обратно точно таким же, каким он был на момент запуска транзакции — это могло бы привести к отмене чьей-нибудь работы!

Например, предположим, что наша транзакция выполняет оператор INSERT, который вызывает выделение нового экстенда (т.е. таблица растет). Этот оператор INSERT должен привести к выделению нового блока, его форматированию перед применением и помещению в него некоторых данных. В тот же самый момент какая-то другая транзакция может вставлять свои данные в этот блок. Если мы произведем откат нашей транзакции, очевидно, что мы не сможем отменить форматирование и выделение нового блока. Таким образом, когда СУБД Oracle осуществляет откат, в действительности она выполняет логический эквивалент противоположности того, что было сделано первоначально. Для каждого оператора INSERT база данных Oracle выполнит оператор DELETE. Для каждого оператора DELETE база данных Oracle выполнит INSERT. Для каждого оператора UPDATE база данных Oracle выполнит “анти-UPDATE”, или оператор UPDATE, который приведет строку в состояние, предшествующее модификации.

На заметку! Генерация undo не касается операций прямого режима, которые обладают способностью обходить генерацию undo в таблице. Мы вскоре обсудим такие операции более подробно.

Как увидеть все это в действии? Возможно, простейший путь предусматривает следование перечисленным ниже шагам.

1. Создать пустую таблицу.
2. Выполнить полное сканирование таблицы и понаблюдать за объемом операций ввода-вывода, необходимых для ее чтения.
3. Наполнить таблицу множеством строк (без фиксации).
4. Произвести откат этой работы, отменив ее.
5. Выполнить полное сканирование таблицы во второй раз и посмотреть, каким будет объем операций ввода-вывода.

Итак, давайте создадим пустую таблицу:

```

EODA@ORA12CR1> create table t
2 as
3 select *
4   from all_objects
5  where 1=0;
Table created.
Таблица создана.

```

А теперь мы будем запускать запросы к ней с включенным режимом AUTOTRACE в SQL*Plus, чтобы измерить ввод-вывод.

На заметку! В рассматриваемом примере каждый раз мы будем выполнять полное сканирование таблицы дважды. Целью является измерение количества операций ввода-вывода только при втором проходе. Это позволяет избежать учета дополнительных операций ввода-вывода, выполненных оптимизатором во время разбора и оптимизации.

Изначально запрос *вообще не* требует операций ввода-вывода для полного сканирования таблицы:

```
EODA@ORA12CR1> select * from t;
no rows selected
нет выбранных строк

EODA@ORA12CR1> set autotrace traceonly statistics
EODA@ORA12CR1> select * from t;
no rows selected
нет выбранных строк

Statistics
-----
      0 recursive calls
      0 db block gets
      0 consistent gets
      0 physical reads

EODA@ORA12CR1> set autotrace off
```

Поначалу может удивить, особенно пользователей версий, предшествовавших Oracle 11g Release 2, что сообщается о *нуле* операций ввода-вывода в таблице. Причина кроется в новом средстве Oracle 11g Release 2 — отложенном создании сегментов.

На заметку! Средство отложенного создания сегментов доступно только в редакции Enterprise системы Oracle. По умолчанию оно включено в Oracle 11g Release 2 и последующих версиях. Это стандартное поведение можно переопределить при создании таблицы. За дополнительными сведениями обращайтесь в главу 10.

Если вы запустите этот пример в более старых выпусках, то, скорее всего, будете наблюдать три или около того операции ввода-вывода. Вскоре мы это обсудим, а пока давайте продолжать работу с примером. Далее мы добавим в таблицу большой объем данных. Мы заставим ее “вырасти”, после чего выполним откат:

```
EODA@ORA12CR1> insert into t select * from all_objects;
18371 rows created.
18371 строк создано.

EODA@ORA12CR1> rollback;
Rollback complete.
Откат завершен.
```

Запросив таблицу снова, мы обнаружим, что для чтения таблицы на этот раз требуется заметно больше операций ввода-вывода:

```
EODA@ORA12CR1> select * from t;
no rows selected
нет выбранных строк
```

```
EODA@ORA12CR1> set autotrace traceonly statistics
EODA@ORA12CR1> select * from t;
no rows selected
нет выбранных строк
```

```
Statistics
```

```
-----
0 recursive calls
0 db block gets
317 consistent gets
0 physical reads
```

```
EODA@ORA12CR1> set autotrace off
```

Блоки, которые были добавлены нашим оператором INSERT ниже маркера максимального уровня заполнения (high-water mark — HWM) таблицы, остались на месте — форматированные, но пустые. Наше полное сканирование должно было прочитать их, чтобы выяснить, содержат ли они какие-то строки. Кроме того, при первом выполнении запроса мы наблюдали *ноль* операций ввода-вывода. Это произошло из-за стандартного режима создания таблиц в Oracle Database 11g Release 2, предусматривающего использование отложенного создания сегментов. Когда выдается команда CREATE TABLE, никакого пространства в хранилище, ни одного экстен-та не выделяется. Создание сегмента отсрочивается до появления первого оператора INSERT, и когда мы производим откат, этот сегмент сохраняется. В этом легко убедиться с помощью простого примера. Здесь отложенное создание сегментов запрашивается явно, хотя оно по умолчанию включено в версии Oracle 11g Release 2:

```
EODA@ORA12CR1> drop table t purge;
```

```
Table dropped.
```

```
Таблица удалена.
```

```
EODA@ORA12CR1> create table t ( x int )
```

```
2 segment creation deferred;
```

```
Table created.
```

```
Таблица создана.
```

```
EODA@ORA12CR1> select extent_id, bytes, blocks
```

```
2 from user_extents
3 where segment_name = 'T'
4 order by extent_id;
```

```
no rows selected
```

```
нет выбранных строк
```

```
EODA@ORA12CR1> insert into t(x) values (1);
```

```
1 row created.
```

```
1 строка создана.
```

```
EODA@ORA12CR1> rollback;
```

```
Rollback complete.
```

```
Откат выполнен.
```

```
EODA@ORA12CR1> select extent_id, bytes, blocks
```

```
2 from user_extents
3 where segment_name = 'T'
4 order by extent_id;
```

```
EXTENT_ID      BYTES      BLOCKS
-----
0          65536          8
```

Как видите, после того, как таблица была первоначально создана, никакое пространство в хранилище не выделялось — таблица не задействовала ни одного экстенста. После выполнения оператора `INSERT`, за которым немедленно следует `ROLLBACK`, мы видим, что оператор `INSERT` вызвал выделение пространства в хранилище, но `ROLLBACK` не “освободил” его.

Вместе эти два факта — то, что сегмент в действительности был создан оператором `INSERT`, но его создание не было “отменено” оператором `ROLLBACK`, и то, что новые форматированные блоки, созданные `INSERT`, были просканированы во второй раз — говорят о том, что откат является операцией логического “возврата базы данных в предыдущее состояние”. База данных не будет точно такой, как была раньше — она будет лишь логически эквивалентной.

Совместная работа redo и undo

Теперь давайте посмотрим, как redo и undo работают вместе в различных сценариях. Мы обсудим, например, что происходит во время обработки оператора `INSERT` в отношении генерации журналов redo и undo, и каким образом Oracle применяет эту информацию в случае отказов в разнообразные моменты времени.

Интересно отметить, что информация undo, сохраненная в табличном пространстве отмены или сегментах отмены, также защищается информацией redo. Другими словами, данные undo трактуются точно так же, как данные таблиц или индексов — изменения в undo генерируют информацию redo, которая записывается в журнал (в журнальный буфер, а затем в файл журнала повторения). Причина проявится очень скоро, когда речь пойдет о том, что происходит при аварийном отказе системы. Данные undo добавляются в сегмент отмены и кешируются в буферном кеше подобно любой другой порции данных.

Пример сценария `INSERT-UPDATE-DELETE-COMMIT`

В этом примере мы предположим, что была создана следующая таблица с индексом:

```
create table t(x int, y int);
create index ti on t(x);
```

Мы посмотрим, что может произойти при выполнении следующего набора операторов:

```
insert into t (x,y) values (1,1);
update t set x = x+1 where x = 1;
delete from t where x = 2;
```

Мы проследим эту транзакцию различными путями и найдем ответы на перечисленные ниже вопросы.

- Что случится, если система откажет в разных точках обработки этих операторов?
- Что случится, если переполнится буферный кеш?
- Что случится, если выполнить `ROLLBACK` в любой точке?
- Что случится, если все пройдет успешно и будет выполнен `COMMIT`?

Оператор *INSERT*

Первоначальный оператор `INSERT INTO T` сгенерирует и redo-, и undo-информацию. Сгенерированной информации undo будет достаточно, чтобы отменить результат работы `INSERT`. Информации redo, сгенерированной `INSERT INTO T`, окажется достаточно, чтобы повторить вставку еще раз.

После того, как вставка произведена, мы получаем сценарий, представленный на рис. 9.1.

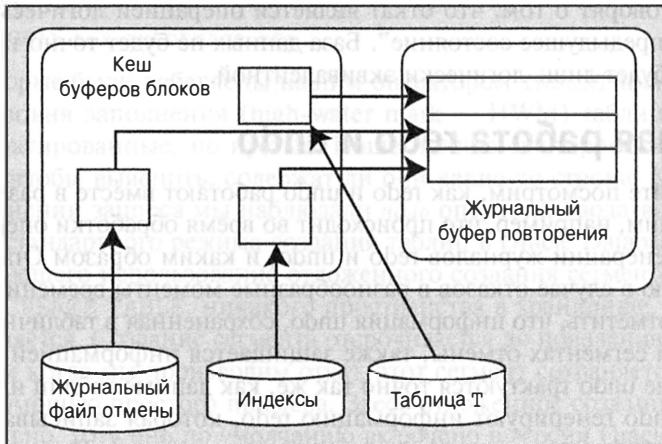


Рис. 9.1. Состояние системы после выполнения оператора `INSERT`

Здесь есть кешированные, модифицированные блоки отмены, индексные блоки и блоки данных таблиц. Каждый из этих блоков защищен записями в журнальном буфере повторения.

Гипотетический сценарий: аварийный отказ системы происходит прямо сейчас

При таком сценарии аварийный отказ системы случается перед выдачей `COMMIT` или перед тем, как записи повторения будут сохранены на диске (детальные сведения о механизмах, вызывающих запись данных redo на диск, приведены в главе 4). Все в порядке. Область SGA очищена, но мы не нуждаемся ни в чем, что находилось в SGA. При перезапуске транзакции все будет выглядеть так, будто она никогда ранее не выполнялась. Ни один из блоков с изменениями не сбрасывается на диск, и никакая информация redo также не сбрасывается на диск. Для восстановления после отказа экземпляра эта информация undo или redo не нужна.

Гипотетический сценарий: буферный кеш переполняется прямо сейчас

Ситуация такова, что процесс записи DBWn должен освободить место, и наши модифицированные блоки должны быть вытолкнуты из кеша. В этом случае DBWn начнет с сообщения процессу LGWR о необходимости выталкивания записей redo, которые защищают эти блоки базы данных. Прежде чем DBWn сможет записать любые измененные блоки на диск, процесс LGWR должен вытолкнуть (на диск) информацию redo, связанную с этими блоками. Это имеет смысл: если мы должны вытолкнуть модифицированные блоки для таблицы T (но не блоки undo, связанные с модификациями) без выталкивания записей redo, ассоциированных с блоками undo,

а система терпит аварийный отказ, то мы получим модифицированный блок таблицы Т без связанной с ним информации undo. Перед записью этих блоков мы должны вытолкнуть буферы журналов redo, чтобы можно было повторить все изменения, необходимые для возврата области SGA в состояние, в котором она находится прямо сейчас, так что может быть выполнен откат.

Второй сценарий демонстрирует определенную предусмотрительность, которая была проявлена во всем этом. Набор условий, описываемый в виде “если мы вытолкнем блоки таблицы Т и не вытолкнем информацию redo для блоков undo и произойдет аварийный отказ системы”, начинает становиться сложным. По мере добавления пользователей, объектов, параллельной обработки и тому подобного он станет еще сложнее.

В этот момент мы имеем дело с ситуацией, показанной на рис. 9.1. Мы сгенерировали ряд модифицированных блоков таблицы и индексов. Они имеют связанные блоки сегментов undo, и все три типа блоков располагают информацией redo, сгенерированной для их защиты. Вспомните из обсуждения журнального буфера повторения в главе 4, что он сбрасывается, *по меньшей мере*, каждые три секунды, когда заполнен на треть или содержит 1 Мбайт буферизованных данных, либо когда встречается оператор COMMIT или ROLLBACK. Весьма вероятно, что в некоторой точке во время обработки журнальный буфер повторения будет очищен. В этом случае картина выглядит так, как представлено на рис. 9.2.

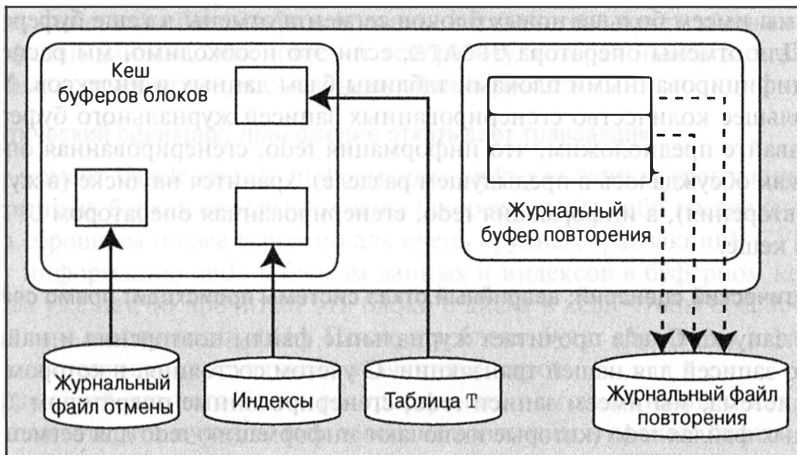


Рис. 9.2. Состояние системы после сброса на диск журнального буфера повторения

То есть мы будем иметь модифицированные блоки, представляющие незафиксированные данные, в кеше буферов и информацию redo для этих незафиксированных изменений на диске. Это совершенно нормальный сценарий, который встречается довольно часто.

Оператор UPDATE

Оператор UPDATE будет инициировать в основном ту же самую работу, что и оператор INSERT. На этот раз объем информации undo будет больше; у нас есть ряд образов “до того”, предназначенных для сохранения в качестве результата UPDATE.

Теперь картина имеет вид, приведенный на рис. 9.3 (темный прямоугольник в обозначении журнального файла повторения представляет информацию redo, сгенерированную оператором INSERT, а информация redo для оператора UPDATE по-прежнему находится в области SGA, и она пока еще не была записана на диск).

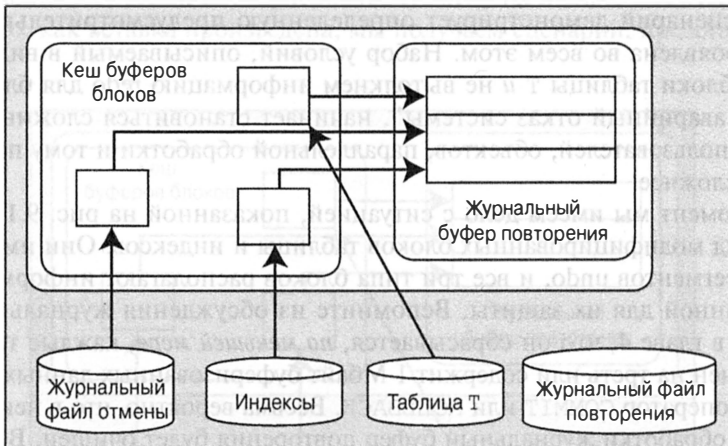


Рис. 9.3. Состояние системы после выполнения оператора UPDATE

Здесь мы имеем больше новых блоков сегмента отмены в кеше буферов блоков данных. Для отмены оператора UPDATE, если это необходимо, мы располагаем в кеше модифицированными блоками таблицы базы данных и индексов. Мы также имеем большее количество сгенерированных записей журнального буфера повторения. Давайте предположим, что информация redo, сгенерированная оператором INSERT (как обсуждалось в предыдущем разделе), хранится на диске (в журнальном файле повторения), а информация redo, сгенерированная оператором UPDATE, находится в кеше.

Гипотетический сценарий: аварийный отказ системы происходит прямо сейчас

После запуска Oracle прочитает журнальные файлы повторения и найдет в них несколько записей для нашей транзакции. С учетом состояния, в котором была оставлена система, мы имеем записи redo, сгенерированные оператором INSERT, в журнальных файлах redo (которые включают информацию redo для сегментов undo, ассоциированных с оператором INSERT). Тем не менее, информация redo для оператора UPDATE присутствовала только в журнальном буфере и никогда не сбрасывалась на диск (вдобавок она была очищена, когда произошел аварийный отказ системы). Ситуация вполне нормальна, транзакция не фиксировалась и файлы данных на диске отражают состояние системы перед запуском UPDATE.

Однако информация redo для оператора INSERT была записана в журнальный файл. Следовательно, Oracle выполнит “накат” (roll forward) оператора INSERT в базе данных. Мы попадаем в состояние, во многом похожее на то, что показано на рис. 9.1, с модифицированными блоками undo (информацией о том, как произвести откат INSERT), модифицированными табличными блоками (сразу после выполнения INSERT) и модифицированными индексными блоками (также сразу после INSERT).

База данных Oracle обнаружит, что наша транзакция не была зафиксирована, и осуществит ее откат, т.к. система выполняет восстановление после сбоя и, конечно же, наш сеанс больше не подключен.

Чтобы произвести откат незафиксированного оператора INSERT, база данных Oracle воспользуется информацией undo, накат которой был только что выполнен (из информации redo, которая теперь в буферном кеше), и применит ее к блокам данных и индексов, возвращая им вид, который они имели до выдачи оператора INSERT. Теперь все становится таким, каким было ранее. Блоки на диске могут отражать или не отражать INSERT (это зависит от того, были ли вытолкнуты наши блоки перед аварийным отказом). Если блоки на диске отражают INSERT, значит, работа оператора INSERT будет отменена, когда блоки сбросятся из буферного кеша. Если же они не отражают отмененный INSERT, то так и быть — они в любом случае будут перезаписаны позднее.

На заметку! Подробное обсуждение создания контрольных точек и ситуации, когда модифицированные (грязные) буферы записываются из буферного кеша на диск, приведено в главе 3.

Описанный сценарий раскрывает элементарные детали восстановления после аварии. Система выполняет его как двухшаговый процесс. Сначала она производит накат базы данных, приводя систему точно в точку отказа, а затем продолжает откатывать все, что еще не было зафиксировано. Это действие будет заново синхронизировать файлы данных. Оно повторит всю работу, которая находилась в процессе выполнения, и отменит все, что не было зафиксировано.

Гипотетический сценарий: приложение откатывает транзакцию

В этот момент Oracle найдет информацию undo для данной транзакции — либо в кешированных блоках сегмента отмены (скорее всего), либо на диске, если они были туда сброшены (более вероятно для очень крупных транзакций). СУБД Oracle применит информацию undo к блокам данных и индексов в буферном кеше, либо если их там уже нет, то прочитает эти блоки с диска в кеш, чтобы обеспечить применение к ним информации undo. Позже эти блоки будут сброшены в файлы данных с восстановленными исходными значениями строк.

Этот сценарий намного более распространен, чем аварийный отказ системы. Полезно отметить, что журналы redo никак не вовлечены в процесс отката. Журналы redo читаются только во время восстановления и архивации. Это ключевая концепция настройки: в журналы redo производится запись. СУБД Oracle не читает их во время нормальной обработки. До тех пор, пока у вас есть достаточное количество устройств, чтобы при чтении файла процессом ARCn процесс LGWR осуществлял запись на другое устройство, конкуренция за журналы redo отсутствует. Многие другие базы данных трактуют журнальные файлы как “журналы транзакций”. Они не имеют такого разделения между redo и undo. Для таких систем действие по откату может оказаться гибельным — процесс отката должен читать журналы, в которые процесс записи в журнал пытается выполнить запись. Они вводят конкуренцию в часть системы, которая может, как минимум, привести к ее останову. Цель Oracle — сделать так, чтобы журналы redo записывались последовательно, и никто никогда не читал их во время записи.

Оператор **DELETE**

В результате выполнения оператора **DELETE** снова генерируется информация **undo**, блоки модифицируются, а информация **redo** отправляется в буфер журнала повторения. Это не особо отличается от описанного ранее. В действительности все настолько похоже на случай с **UPDATE**, что мы можем сразу переходить к рассмотрению **COMMIT**.

Оператор **COMMIT**

Мы взглянули на разные сценарии отказов и различные пути и теперь, наконец, сделаем это в отношении оператора **COMMIT**. Здесь Oracle будет сбрасывать буферы журнала повторения на диск, и картина выглядит так, как представлено на рис. 9.4.

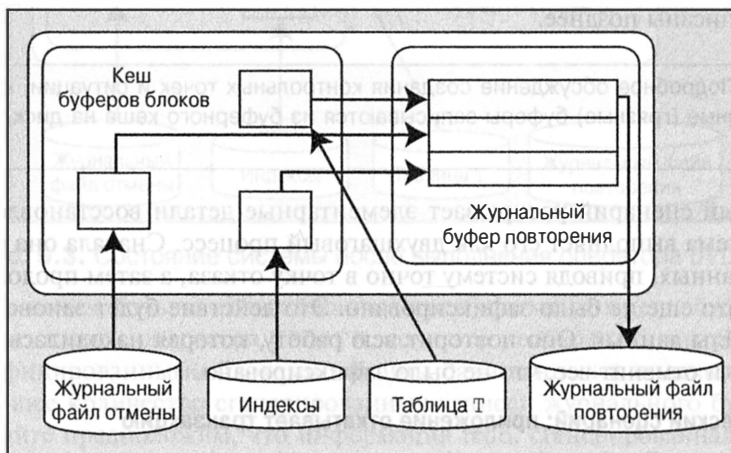


Рис. 9.4. Состояние системы после выполнения оператора **COMMIT**

Модифицированные блоки располагаются в буферном кеше; возможно, некоторые из них были сброшены на диск. *Вся информация redo*, необходимая для повторения транзакции, благополучно находится на диске и изменения теперь постоянны. Если бы мы читали данные прямо из файлов данных, то вероятно видели бы блоки в состоянии, в котором они пребывали *перед* транзакцией, поскольку процесс **DBWn**, скорее всего, еще их не записал. Это нормально — журнальные файлы повторения могут использоваться для приведения этих блоков в актуальное состояние в случае отказа. Информация **undo** будет храниться до тех пор, пока сегменты отмены не начнут перезаписываться, снова утилизировав эти блоки. Эта информация **undo** будет применяться Oracle, чтобы обеспечить согласованное чтение затронутых объектов для любых сеансов, которые в них нуждаются.

Обработка **COMMIT** и **ROLLBACK**

Важно понимать, каким образом файлы журналов **redo** могут воздействовать на разработчиков. Мы посмотрим, как разные способы написания кода влияют на использование журналов **redo**. Ранее в настоящей главе вы уже видели механизм повторения, а теперь мы рассмотрим некоторые специфические проблемы. Вы можете

обнаружить многие из этих сценариев, но они будут исправлены администратором базы данных, т.к. затрагивают экземпляр базы данных в целом. Мы начнем с того, что происходит во время выполнения оператора COMMIT, а затем перейдем к часто задаваемым вопросам и проблемам, сопровождающим оперативные журналы redo.

Что делает оператор COMMIT?

Как разработчик, вы должны хорошо понимать, что в точности делается во время выполнения оператора COMMIT. В этом разделе мы исследуем то, что происходит, когда оператор COMMIT обрабатывается в Oracle. Обычно COMMIT — очень быстрая операция, не зависящая от размера транзакции. Может показаться, что чем крупнее транзакция (т.е. чем больший объем данных она затрагивает), тем дольше будет выполняться оператор COMMIT. Это не так. Время ответа COMMIT обычно “однородно” и не зависит от размера транзакции. Причина в том, что на самом деле COMMIT не приходится выполнять слишком много работы, но то, что он делает, является жизненно важным.

Одна из причин важности понимания и принятия этого факта состоит в том, что транзакции могут быть настолько велики, насколько это необходимо. Как было показано в предыдущей главе, многие разработчики искусственно ограничивают размеры своих транзакций, производя фиксацию небольших групп строк вместо того, чтобы выполнять фиксацию по окончании логической единицы работы. Они поступают так из-за ошибочной уверенности в том, что тем самым сберегают дефицитные системные ресурсы, когда на деле только увеличивают их расход. Если оператор COMMIT для одной строки занимает X единиц времени и COMMIT для 1000 строк требует тех же X единиц, то выполнение работы в стиле с 1000 однострочных операторов COMMIT займет $1000 \cdot X$ дополнительных единиц времени. Производя фиксацию только тогда, когда она должна делаться (при завершении логической единицы работы), вы не только увеличиваете производительность, но и сокращаете конкуренцию за разделяемые ресурсы (журнальные файлы, разнообразные внутренние защелки и т.п.). Простой пример демонстрирует, что многократная фиксация неминуемо занимает больше времени. Мы будем применять Java-программу, хотя похожих результатов можно ожидать практически от любого клиента — за исключением в этом случае PL/SQL (причины объясняются после примера). Для начала понадобится таблица для вставки данных:

```
SCOTT@ORA12CR1> create table test
2 ( id          number,
3   code        varchar2(20),
4   descr       varchar2(20),
5   insert_user varchar2(30),
6   insert_date date
7 )
8 /
Table created.
Таблица создана.
```

Наша Java-программа (сохраненная в файле по имени `perfctest.java`) будет принимать два входных параметра: количество строк для вставки (`iters`) и количество фиксируемых строк (`commitCnt`). Она начинается с подключения к базе дан-

ных, отключения режима автоматической фиксации (что должно делаться в любом коде Java) и двухкратного вызова метода `doInserts()`:

- первый раз, чтобы “разогреть” процедуру (обеспечить загрузку всех нужных классов);
- второй раз, при включенном средстве трассировки SQL, с указанием количества строк для вставки наряду с количеством, которые должны фиксироваться за раз (т.е. фиксировать каждые N строк).

Затем программа закрывает соединение и завершается. Метод `main` выглядит следующим образом:

```
import java.sql.*;

public class perfctest
{
    public static void main (String arr[]) throws Exception
    {
        DriverManager.registerDriver(new oracle.jdbc.OracleDriver());
        Connection con = DriverManager.getConnection
            ("jdbc:oracle:thin:@csxdev:1521:ORA12CR1", "scott", "tiger");
        Integer iters = new Integer(arr[0]);
        Integer commitCnt = new Integer(arr[1]);

        con.setAutoCommit(false);
        doInserts( con, 1, 1 );

        Statement stmt = con.createStatement ();
        stmt.execute( "begin dbms_monitor.session_trace_enable(waits=>true); end;" );
        doInserts( con, iters.intValue(), commitCnt.intValue() );
        con.close();
    }
}
```

На заметку! Учетная запись SCOTT или любая другая учетная запись, которую вы используете для тестирования этого кода, должна иметь привилегию EXECUTE, выданную для пакета DBMS_MONITOR.

Метод `doInserts()` довольно прост. Он начинается с подготовки (разбора) оператора INSERT, чтобы мы могли привязывать и выполнять его снова и снова:

```
static void doInserts(Connection con, int count, int commitCount )
throws Exception
{
    PreparedStatement ps =
        con.prepareStatement
            ("insert into test " +
             "(id, code, descr, insert_user, insert_date)"
             + " values (?, ?, ?, user, sysdate)");
```

Затем организуется цикл по количеству строк для вставки, на каждом шаге которого производится привязка и выполнение оператора INSERT. Дополнительно в цикле проверяется счетчик строк, чтобы выяснить, выполнять ли COMMIT внутри цикла:

```

int rowcnt = 0;
int committed = 0;
for (int i = 0; i < count; i++)
{
    ps.setInt(1,i);
    ps.setString(2,"PS - code" + i);
    ps.setString(3,"PS - desc" + i);
    ps.executeUpdate();
    rowcnt++;
    if ( rowcnt == commitCount )
    {
        con.commit();
        rowcnt = 0;
        committed++;
    }
}
con.commit();
System.out.println
("pstatement rows/commitcnt = " + count + " / " + committed );
}
}

```

Теперь мы будем многократно запускать этот код с разными входными параметрами и просматривать результирующие файлы TKPROF. Мы будем производить 100 000 вставок строк, фиксируя по 1 строке за раз, затем по 10 и т.д. Результаты, которые попадут в файлы TKPROF, представлены в табл. 9.1.

Таблица 9.1. Результаты вставки 100 000 строк

| Количество вставленных строк | Фиксация через каждые N строк, $N=$ | Время центрального процессора для оператора INSERT (в секундах) | Время ожидания синхронизации журнальных файлов (в секундах) |
|------------------------------|---------------------------------------|---|---|
| 100 000 | 1 | 2,30 | 31,17 |
| 100 000 | 10 | 2,16 | 3,48 |
| 100 000 | 100 | 2,20 | 0,62 |
| 100 000 | 1 000 | 2,02 | 0,08 |
| 100 000 | 10 000 | 1,46 | 0,02 |
| 100 000 | 100 000 | 2,01 | 0,00 |

Как видите, чем чаще выполняется фиксация, тем дольше приходится ожидать (ваши показатели будут варьироваться). К тому же промежуток времени ожидания более или менее прямо пропорционален количеству фиксаций. Не забывайте, что это только однопользовательский сценарий; при множестве пользователей, делающих ту же самую работу, слишком частые фиксации быстро увеличат эти цифры.

Мы слышим ту же историю снова и снова с похожими ситуациями. Например, мы видели, что отказ от переменных привязки и частое выполнение полных разборов значительно снижает степень параллелизма из-за конкуренции за содержимое библиотечного кеша и чрезмерной утилизации центрального процессора. Даже

если мы перейдем на применение переменных привязки, то слишком часто производимый частичный разбор, вызванный закрытием курсоров, несмотря на то, что мы собираемся вскоре использовать их повторно, ведет к значительным накладным расходам. Мы должны выполнять операции только тогда, когда это необходимо — COMMIT является просто одной из таких операций. Лучше устанавливать размеры транзакций на основе бизнес-потребностей, а не на ошибочном стремлении снизить потребление ресурсов в базе данных.

Существуют два фактора, которые приводят к увеличению накладных расходов операции COMMIT в этом примере.

- Мы очевидным образом увеличиваем количество круговых обменов с базой данных. Фиксация каждой записи вызывает генерирование намного большего трафика в обоих направлениях. Я даже не измерял то, что добавилось бы к общему времени выполнения.
- Каждый раз, выполняя фиксацию, мы должны ждать записи на диск информации redo. Это приводит к ожиданию. В данном случае ожидание называется “синхронизацией журнальных файлов”.

Таким образом, мы производим фиксацию после каждого оператора INSERT, каждый раз ждем короткое время, и если ожидание случается на протяжении короткого времени, но часто, то все эти периоды суммируются. Почти 30 секунд времени выполнения было потрачено на ожидание фиксации, когда оператор COMMIT выдавался 100 000 раз — другими словами, на ожидание записи на диск информации redo процессом LGWR. Разительный контраст с этим представляет случай с однократной фиксацией — нам не приходится долго ждать (на самом деле вообще неизмеримо малое время). Это доказывает, что COMMIT — быстрая операция; мы предполагаем, что время ее выполнения должно быть более или менее однородным, а не функцией от объема сделанной работы.

Почему же время ответа COMMIT почти однородно независимо от размера транзакции? Причина в том, что перед выдачей COMMIT в базе данных уже было сделано действительно много работы. Данные в базе уже модифицированы, так что 99,9% работы выполнено. Например, операции вроде перечисленных ниже уже произошли.

- Блоки undo сгенерированы в области SGA.
- Модифицированные блоки данных сгенерированы в области SGA.
- Буферизованная информация redo для предыдущих двух позиций сгенерирована в области SGA.
- В зависимости от размера предыдущих трех позиций и объема потраченного времени некоторые комбинации существующих данных уже могут быть сброшены на диск.
- Все блокировки получены.

Когда мы запускаем оператор COMMIT, выполняются следующие действия.

- Для транзакции генерируется системный номер изменения (system change number — SCN). На тот случай, если вы не знакомы с этим понятием, то знайте, что SCN представляет собой простой механизм синхронизации, который Oracle применяет для гарантирования порядка выполнения транзакций и пре-

доставления возможности восстановления после аварийного отказа. Он также используется для обеспечения согласованности чтения и сохранения контрольных точек в базе данных. Воспринимайте SCN как счетчик; всякий раз, когда кто-то выполняет COMMIT, значение SCN увеличивается на единицу.

- Процесс LGWR записывает все оставшиеся буферизованные записи журнала redo на диск и также сохраняет SCN в оперативных журнальных файлах повторения. Данный шаг в действительности и есть COMMIT. Если этот шаг выполнен, то фиксация произошла. Наша запись о транзакции “удаляется” из представления V\$TRANSACTION, и это показывает, что фиксация выполнена.
- Все блокировки, записанные в представление V\$LOCK нашим сеансом, освобождаются, и все, кто стоит в очереди и ожидает снятия блокировок, пробуждаются и могут продолжать свою работу.
- Некоторые блоки, модифицированные нашей транзакцией, будут посещены и “очищены” в быстром режиме, если они все еще находятся в буферном кеше. Очистка блока имеет отношение к информации, связанной с блокировками, которую мы сохраняем в заголовке блока базы данных. В основном мы очищаем в блоке информацию о транзакции, так что это не придется делать следующему клиенту, посетившему блок. Мы делаем это таким способом, который не требует генерации журнальной информации повторения, что позже экономит значительный объем работы (подробные сведения ищите в разделе “Очистка блоков” далее в главе).

Как видите, для обработки оператора COMMIT нужно сделать совсем немного. Самой длинной операцией является, и всегда будет таковой, деятельность, выполняемая процессом LGWR, т.к. она связана с физическим дисковым вводом-выводом. Объем затрат времени, потраченного здесь процессом LGWR, будет значительно сокращен из-за того факта, что содержимое буфера журнала redo уже сбрасывается на диск на регулярной основе. Процесс LGWR не будет буферизовать всю проделываемую вами работу на протяжении ее выполнения. Вместо этого он инкрементным образом выталкивает содержимое буфера журнала redo в фоновом режиме. Это позволяет оператору COMMIT не ждать очень долго, пока вся информация redo будет сброшена на диск за один раз.

Таким образом, даже при наличии длительной транзакции большая часть сгенерированного ею буферизованного журнала redo будет сброшена на диск еще до самой фиксации. Обратной стороной является тот факт, что когда производится COMMIT, нам обычно приходится ждать, пока *вся* буферизованная информация redo, которая еще не была записана, сохранится на диске. То есть обращение к LGWR по умолчанию является *синхронным*. Хотя процесс LGWR может применять асинхронный ввод-вывод для параллельной записи в журнальные файлы, наша транзакция, как правило, будет ожидать, пока LGWR не завершит все операции записи и не получит подтверждение, что данные существуют на диске, прежде чем возвращать управление.

На заметку! В главе 8 было указано, что в Oracle 11g Release 1 и последующих версиях доступно асинхронное ожидание. Тем не менее, как там упоминалось, в общих ситуациях этот стиль фиксации имеет ограниченное использование. Фиксации в любом приложении, ориентированном на конечного пользователя, должны быть синхронными.

Ранее уже говорилось о том, что мы применяли программу на Java, а не на PL/SQL, по определенной причине, и причина эта связана с оптимизацией времени фиксации PL/SQL, как обсуждалось в главе 8. Было указано, что обращение к процессу LGWR по умолчанию является синхронным, и приходится ожидать, пока он не завершит свою запись. Это справедливо в Oracle 12c Release 1 и предшествующих версиях для всех языков программирования *кроме PL/SQL*.

Механизм PL/SQL, выяснив, что клиент не узнает, произойдет ли COMMIT в процедуре PL/SQL до ее окончания, делает асинхронную фиксацию. Он не ждет завершения записи процессом LGWR, а осуществляет немедленный возврат из вызова COMMIT. Однако когда процедура PL/SQL завершена и происходит возврат управления из базы данных клиенту, процедура PL/SQL будет ожидать, пока процесс LGWR выполнит все оставшиеся фиксации COMMIT.

Таким образом, если вы 100 раз произвели фиксацию в коде PL/SQL и затем возвращаете управление клиенту, то, скорее всего, обнаружите, что благодаря этой оптимизации ожидание процесса LGWR происходит только один раз, а не 100. Подразумевает ли это, что частая фиксация в PL/SQL — хорошая или допустимая идея? Вовсе нет! Это просто *не настолько плохая идея*, как в других языках. Руководящее правило заключается в том, что фиксацию нужно выполнять по завершении логической единицы работы, но не раньше.

На заметку! Такая оптимизация времени фиксации в PL/SQL может быть приостановлена, когда выполняются распределенные транзакции либо имеется конфигурация Data Guard в режиме максимальной доступности. Поскольку есть два участника, инструмент PL/SQL должен ждать полного окончания фиксации, прежде чем продолжить работу. Также в Oracle 11g Release 1 и последующих версиях указанную оптимизацию можно приостановить прямым вызовом COMMIT WORK WRITE WAIT в PL/SQL.

Для демонстрации того, что COMMIT является операцией с “однородным временем ответа”, мы будем генерировать разные по объему информации redo и времени операторы INSERT и COMMIT. По мере выполнения этих операторов INSERT и COMMIT мы будем измерять объем информации redo, которая генерируется нашим сеансом, с использованием небольшой служебной функции:

```
EODA@ORA12CR1> create or replace function get_stat_val( p_name in
varchar2 ) return number
2 as
3     l_val number;
4 begin
5     select b.value
6     into l_val
7     from v$statname a, v$mystat b
8     where a.statistic# = b.statistic#
9           and a.name = p_name;
10
11     return l_val;
12 end;
13 /
Function created.
Функция создана.
```

На заметку! Владельцу предыдущей функции понадобится напрямую выдать привилегию SELECT на представлениях V\$STATNAME и V\$MYSTAT.

Удалите таблицу T (если она существует) и создайте пустую таблицу T с той же структурой, что и у BIG_TABLE:

```
EODA@ORA12CR1> drop table t purge;
EODA@ORA12CR1> create table t
2 as
3 select *
4 from big_table
5 where 1=0;
Table created.
Таблица создана.
```

На заметку! Указания о том, как создавать и наполнять таблицу BIG_TABLE, которая применяется во многих примерах, приведены в разделе “Настройка среды” в начале этой книги.

Мы будем измерять время центрального процессора (CPU) и общее (Elapsed) время, используемое для фиксации транзакции, с помощью процедур GET_CPU_TIME и GET_TIME из пакета DBMS_UTILITY. Действительный блок PL/SQL, применяемый для генерации рабочей нагрузки и отчета, выглядит следующим образом:

```
EODA@ORA12CR1> declare
2   l_redo number;
3   l_cpu number;
4   l_ela number;
5   begin
6     dbms_output.put_line
7       ( '-' || '      Rows' || '      Redo' ||
8         '      CPU' || ' Elapsed' );
9     for i in 1 .. 6
10    loop
11      l_redo := get_stat_val( 'redo size' );
12      insert into t select * from big_table where rownum <= power(10,i);
13      l_cpu := dbms_utility.get_cpu_time;
14      l_ela := dbms_utility.get_time;
15      commit work write wait;
16      dbms_output.put_line
17        ( '-' ||
18          to_char( power( 10, i ), '9,999,999' ) ||
19          to_char( (get_stat_val('redo size')-l_redo), '999,999,999' )
20          ||
21          to_char( (dbms_utility.get_cpu_time-l_cpu), '999,999' ) ||
22          to_char( (dbms_utility.get_time-l_ela), '999,999' ) );
23    end loop;
24  end;
```

| | Rows | Redo | CPU | Elapsed |
|---|-----------|-------------|-----|---------|
| - | 10 | 7,072 | 0 | 1 |
| - | 100 | 10,248 | 0 | 0 |
| - | 1,000 | 114,080 | 0 | 0 |
| - | 10,000 | 1,146,484 | 0 | 2 |
| - | 100,000 | 11,368,512 | 0 | 2 |
| - | 1,000,000 | 113,800,488 | 1 | 2 |

PL/SQL procedure successfully completed.

Процедура PL/SQL успешно завершена.

* Этот тест был выполнен на однопользовательской машине с журнальным буфером 1,7 Мбайт и тремя оперативными журнальными файлами повторения размером по 500 Мбайт. Значения времени указаны в сотых долях секунды.

Как видите, хотя мы генерируем варьирующиеся объемы информации redo, от 7072 байтов до 113 Мбайт, разница во времени выполнения COMMIT не поддается измерению таймером с разрешающей способностью в одну сотую секунды. По мере обработки и генерации журнала redo, процесс LGWR постоянно сбрасывает буферизованную информацию redo на диск в фоновом режиме. Таким образом, когда мы сгенерировали 113 Мбайт журнальной информации redo, процесс LGWR был занят выталкиванием на диск каждого 1 Мбайт данных или около того. Когда дело дошло до COMMIT, оставалось выполнить не особо много — не намного больше, чем при создании десяти строк данных. Вы должны ожидать увидеть похожие (но не точно такие же) результаты независимо от объема сгенерированной информации redo.

Что делает оператор ROLLBACK?

Заменив COMMIT оператором ROLLBACK, мы можем ожидать совершенно другого результата. Время на откат определенно является функцией от объема модифицированных данных. Я изменил сценарий, разработанный в предыдущем разделе, чтобы он вместо COMMIT выполнял ROLLBACK, и показатели времени стали существенно отличаться. Взгляните на результаты:

```

EODA@ORA12CR1> declare
2     l_redo number;
3     l_cpu  number;
4     l_ela  number;
5 begin
6     dbms_output.put_line
7     ( '-' || '      Rows' || '      Redo' ||
8       '      CPU' || ' Elapsed' );
9     for i in 1 .. 6
10    loop
11        l_redo := get_stat_val( 'redo size' );
12        insert into t select * from big_table where rownum <= power(10,i);
13        l_cpu  := dbms_utility.get_cpu_time;
14        l_ela  := dbms_utility.get_time;
15        --commit work write wait;
16        rollback;
17        dbms_output.put_line
18        ( '-' ||
19          to_char( power( 10, i ), '9,999,999' ) ||

```

```

20      to_char( (get_stat_val('redo size')-l_redo), '999,999,999' ) ||
21      to_char( (dbms_utility.get_cpu_time-l_cpu), '999,999' ) ||
22      to_char( (dbms_utility.get_time-l_ela), '999,999' ) );
23  end loop;
24 end;
25 /

```

| | Rows | Redo | CPU | Elapsed |
|---|-----------|-------------|-----|---------|
| - | 10 | 7,180 | 0 | 0 |
| - | 100 | 10,872 | 0 | 0 |
| - | 1,000 | 121,880 | 0 | 0 |
| - | 10,000 | 1,224,864 | 0 | 0 |
| - | 100,000 | 12,148,416 | 2 | 4 |
| - | 1,000,000 | 121,733,580 | 25 | 36 |

PL/SQL procedure successfully completed.

Процедура PL/SQL успешно завершена.

Такой разницы в значениях показателей CPU и Elapsed следовало ожидать, потому что оператор ROLLBACK должен отменить сделанную нами работу. Подобно COMMIT, здесь также должна быть выполнена последовательность операций. До того как приступить к ROLLBACK, база данных уже выполнила много работы. Ниже перечислено, что должно было произойти.

- Записи сегмента undo сгенерированы в области SGA.
- Модифицированные блоки данных сгенерированы в области SGA.
- Буферизованная информация redo для предыдущих двух позиций сгенерирована в области SGA.
- В зависимости от размера предыдущих трех позиций и объема потраченного времени, некоторые комбинации существующих данных уже могут быть сброшены на диск.
- Все блокировки получены.

Когда мы запускаем оператор ROLLBACK, выполняются следующие действия.

- Отменяются все произведенные изменения. Это достигается чтением данных из сегмента отмены и, по сути, реверсированием операции с последующей пометкой записи undo как примененной. Если была вставлена строка, то ROLLBACK удалит ее. Если строка была обновлена, откат обратит обновление. Если строка была удалена, откат вставит ее обратно.
- Все блокировки, удерживаемые сеансом, освобождаются, и все, кто стоит в очереди и ожидает снятия блокировок, продолжают свою работу.

С другой стороны, оператор COMMIT просто выталкивает все данные, оставшиеся в буферах журнала redo. По сравнению с ROLLBACK он выполняет очень мало работы. Смысл здесь в том, что обычно вы не хотите выполнять откат, если только не обязаны. Это дорогая в плане ресурсов операция, потому что вы потратили много времени на проведение работ, и теперь еще придется уделить массу времени на отмену работ. Не делайте работу, если вы не уверены в том, что впоследствии ее зафиксируете. Звучит как прописная истина: конечно же, не имеет смысла выполнять все работы, если не планируется зафиксировать их посредством COMMIT. Однако

я много раз наблюдал, как разработчик использует “реальную” таблицу в качестве временной, наполняет ее данными, генерирует отчет по ним и затем производит откат, чтобы избавиться от временных данных. Позже мы поговорим о подлинных временных таблицах и о том, как избежать этой проблемы.

Исследование redo

Разработчику часто важно иметь возможность измерить объем генерируемой его операторами информации redo. Чем больше такой информации вы генерируете, тем больше времени могут потребовать ваши операции, и тем медленнее будет работать вся система. Вы повлияете не только на *свой* сеанс, а на абсолютно *все* сеансы. Управление пространством redo представляет собой точку сериализации внутри базы данных. В любом экземпляре Oracle есть только один процесс LGWR, и в конечном итоге все транзакции попадают в LGWR, запрашивая у него обработку их информации redo и выполнение фиксации. Чем больше работы этому процессу приходится делать, тем медленнее функционирует система. За счет оценки объема информации redo, генерируемой операцией, и сравнения нескольких подходов к решению задачи можно найти наилучший способ реализации.

На заметку! Начиная с версии Oracle 12c (в многопроцессорных системах), автоматически запускаются рабочие процессы записи в журнал (LG00), направленные на улучшение производительности записи в журнальный файл повторения.

Измерение redo

Как было показано ранее в этой главе, узнать объем сгенерированной информации redo очень просто. Я использовал встроенное в SQL*Plus средство AUTOTRACE. Однако AUTOTRACE работает только с простыми операциями DML — оно не может, к примеру, применяться для выяснения, какая хранимая процедура была вызвана. Я также использовал представленную ранее служебную функцию GET_STAT_VAL, чтобы извлечь значение “redo size” (размер redo) из таблиц V\$. Она же применяется и в следующем упражнении.

Давайте взглянем на отличие в журналах redo, сгенерированных обычными повседневными операторами INSERT и прямыми (direct-path) операторами INSERT, которые используются при загрузке больших объемов данных в базу. В этом простом примере применяется средство AUTOTRACE и ранее созданные таблицы T и BIG_TABLE. Сначала мы загрузим таблицу с использованием обычного оператора INSERT:

```
EODA@ORA12CR1> set autotrace traceonly statistics;
EODA@ORA12CR1> truncate table t;
Table truncated.
Таблица усечена.

EODA@ORA12CR1> insert into t
2  select * from big_table;
1000000 rows created.
1000000 строк создано.
```

Statistics

```

-----
      90 recursive calls
    123808 db block gets
     39407 consistent gets
    13847 physical reads
113875056 redo size
      1177 bytes sent via SQL*Net to client
     1354 bytes received via SQL*Net from client
         4 SQL*Net roundtrips to/from client
         2 sorts (memory)
         0 sorts (disk)
    1000000 rows processed

```

Несложно заметить, что оператор INSERT генерирует около 113 Мбайт информации redo; мы ожидали это, учитывая предыдущий пример на PL/SQL.

На заметку! Пример, рассматриваемый в этом разделе, запускался в базе данных, работающей в режиме NOARCHIVELOG. Если база данных функционирует в режиме ARCHIVELOG, то таблица должна быть создана или установлена как NOLOGGING, чтобы можно было наблюдать это разительное отличие. В разделе “Установка NOLOGGING в SQL” мы обсудим атрибут NOLOGGING более подробно. В реальной системе обязательно согласовывайте все операции, не записывающие в журналы, с администратором базы данных.

В случае применения прямой загрузки в базу данных, работающую в режиме NOARCHIVELOG, мы получим следующие результаты:

```

EODA@ORA12CR1> truncate table t;
Table truncated.
Таблица усечена.

EODA@ORA12CR1> insert /*+ APPEND */ into t
  2 select * from big_table;
1000000 rows created.
1000000 строк создано.

```

Statistics

```

-----
     551 recursive calls
    16645 db block gets
    15242 consistent gets
    13873 physical reads
220504 redo size
     1160 bytes sent via SQL*Net to client
    1368 bytes received via SQL*Net from client
         4 SQL*Net roundtrips to/from client
        86 sorts (memory)
         0 sorts (disk)
    1000000 rows processed

```

```
EODA@ORA12CR1> set autotrace off
```

Этот оператор INSERT генерирует только около 240 Кбайт информации redo — *Кбайт, а не Мбайт*. Обрисованный метод использования представления V\$MYSTAT в целом удобен для оценки побочных эффектов различных опций. Сценарий GET_STAT_VAL удобен для небольших тестов с одной или двумя операциями.

Можно ли отключить генерацию журналов redo?

Этот вопрос задают часто. Простой и краткий ответ на него — нет, поскольку ведение журналов redo критически важно для базы данных; это не накладные расходы и не пустая трата ресурсов. Оно действительно необходимо вне зависимости от того, верите вы в это или нет. Это жизненный факт, а также способ работы базы данных. Если вы отключите генерацию журналов redo, то любой временный отказ дисковых устройств, отключение электропитания или аварийное завершение программного обеспечения приведет к необратимой порче базы данных без возможности ее восстановления. Тем не менее, существует несколько операций, которые в ряде случаев могут быть выполнены без генерации журналов redo.

На заметку! В Oracle9i Release 2 администратор базы данных может перевести базу данных в режим `FORCE LOGGING`. В этом случае *все* операции осуществляют запись в журнал. Чтобы посмотреть, включен ли этот режим, можно воспользоваться запросом `SELECT FORCE_LOGGING FROM V$DATABASE`. Эту возможность поддерживает конфигурация Data Guard — средство восстановления в аварийных ситуациях Oracle, которое опирается на информацию redo при ведении резервной копии базы данных.

Установка `NOLOGGING` в SQL

Некоторые операторы и операции SQL поддерживают конструкцию `NOLOGGING`. Это вовсе не означает, что все операции в отношении объекта будут выполняться без генерации журнала redo — просто определенные очень *специфические* операции будут генерировать *значительно меньше* информации redo, чем обычно. Обратите внимание — “значительно меньше информации redo”, а не “вообще никакой информации redo”. Все операции генерируют некоторую информацию redo — все операции со словарем данных будут заносить сведения в журнал независимо от режима ведения журналов. Однако объем сгенерированной информации redo может быть существенно уменьшен. Для этого примера применения конструкции `NOLOGGING` были выполнены следующие команды в базе, запущенной в режиме `ARCHIVELOG`:

```
EODA@ORA12CR1> select log_mode from v$database;
LOG_MODE
-----
ARCHIVELOG
EODA@ORA12CR1> drop table t purge;
Table dropped.
Таблица удалена.
EODA@ORA12CR1> variable redo number
EODA@ORA12CR1> exec :redo := get_stat_val( 'redo size' );
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.
EODA@ORA12CR1> create table t
2 as
3 select * from all_objects;
Table created.
```

```

EODA@ORA12CR1> exec dbms_output.put_line( (get_stat_val('redo size')-:redo)
❖|| ' bytes of redo generated...' );
4487796 bytes of redo generated...
PL/SQL procedure successfully completed.
Сгенерировано 4487796 байтов информации redo...
Процедура PL/SQL успешно завершена.

```

Этот оператор CREATE TABLE сгенерировал около 4 Мбайт информации redo (ваши результаты могут варьироваться в зависимости от того, сколько строк вставлено в таблицу T). Мы удалим ее и создадим заново, но уже в режиме NOLOGGING:

```

EODA@ORA12CR1> drop table t;
Table dropped.
Таблица удалена.

EODA@ORA12CR1> variable redo number
EODA@ORA12CR1> exec :redo := get_stat_val( 'redo size' );
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.

EODA@ORA12CR1> create table t
  2 NOLOGGING
  3 as
  4 select * from all_objects;
Table created.
Таблица создана.

EODA@ORA12CR1> exec dbms_output.put_line( (get_stat_val('redo size')-:redo)
❖|| ' bytes of redo generated...' );
90108 bytes of redo generated...
PL/SQL procedure successfully completed.
Сгенерировано 90108 байтов информации redo...
Процедура PL/SQL успешно завершена.

```

На этот раз генерируется всего 90 Кбайт информации redo. Как видите, разница огромная — 4 Мбайт против 90 Кбайт. Данные объемом 4 Мбайт, записанные в первом примере, представляют собой копию действительных данных самой таблицы; они были занесены в журнал redo, когда таблица создавалась без конструкции NOLOGGING.

Если вы протестируете это в базе данных, работающей в режиме NOARCHIVELOG, то не увидите никаких отличий. В таком режиме базы данных оператор CREATE TABLE не будет записывать сведения в журнал за исключением модификаций словаря данных. Из этого факта также вытекает важный совет: тестируйте свою систему в том режиме, в котором она будет эксплуатироваться в производственной среде, поскольку поведение может отличаться. Производственная система будет функционировать в режиме ARCHIVELOG; если вы выполняете множество операций, которые генерируют информацию redo в этом режиме, но не в режиме NOARCHIVELOG, то лучше обнаружить это на этапе тестирования, а не по отзывам пользователей!

Разумеется, теперь очевидно, что в режиме NOARCHIVELOG вы будете делать все возможное, правильно? В действительности ответ снова отрицателен. Этот режим должен использоваться очень осторожно и только после обсуждения проблем с лицом, отвечающим за резервное копирование и восстановление. Предположим, что вы создали показанную таблицу, и она стала частью вашего приложения (т.е. вы применяли оператор CREATE TABLE AS SELECT NOLOGGING в сценарии модерни-

зации). Ваши пользователи модифицируют эту таблицу на протяжении дня. Ночью диск, содержащий таблицу, выходит из строя. “Ничего страшного”, — говорит администратор базы данных, — “мы работаем в режиме ARCHIVELOG и можем выполнить восстановление носителя”. Однако проблема в том, что поскольку изначально созданная таблица не была занесена в журнал, восстановить ее из архивного журнала redo невозможно. Таблица не подлежит восстановлению, и это приводит нас к наиболее важному выводу относительно операций NOLOGGING: они должны быть согласованы с администратором базы данных и системой в целом. Если вы их используете, а другие не осведомлены об этом факте, то вы можете подвергнуть риску способность администратора базы данных провести полное восстановление вашей базы после отказа носителя. Операции NOLOGGING следует применять рассудительно и осторожно.

Ниже перечислены важные замечания об операциях NOLOGGING.

- Некоторый объем информации redo будет генерироваться в любом случае. Эта информация redo предназначена для защиты словаря данных. Избежать ее генерации нельзя. Объем этой информации может быть значительно меньше, чем ранее, но он будет ненулевым.
- Конструкция NOLOGGING не предотвращает генерирования информации redo всеми последующими операциями. В предыдущем примере я не создавал таблицу, которая никогда не записывается в журнал. Не протоколировалась в журнале только одна операция создания таблицы. Все последующие “нормальные” операции, такие как INSERT, UPDATE и DELETE, будут записываться в журнал. Другие специальные операции, подобные загрузке в прямом режиме с использованием SQL*Loader либо вставке в прямом режиме с применением синтаксиса INSERT /*+ APPEND */ , протоколироваться в журнале не будут (если только вы снова не включите для этой таблицы полное ведение журналов). Однако в общем случае операции, выполняемые вашим приложением в отношении этой таблицы, будут записываться в журнал.
- После выполнения операций NOLOGGING в базе данных, работающей в режиме ARCHIVELOG, вы должны как можно скорее создать новую опорную резервную копию затронутых ими файлов данных, чтобы избежать утери последующих изменений этих объектов из-за отказа носителя. Поскольку данные, созданные операцией NOLOGGING, не находятся в журнальных файлах redo и пока еще отсутствуют в резервных копиях, нет никакого способа их восстановления!

Установка NOLOGGING для индекса

Существуют два метода использования опции NOLOGGING. Вы уже видели один из них — встраивание ключевого слова NOLOGGING в команду SQL. Другой метод, предполагающий установку атрибута NOLOGGING на сегменте (индекса или таблицы), позволяет некоторым операциям неявно выполняться в режиме NOLOGGING. Например, можно изменить таблицу или индекс так, чтобы они были NOLOGGING по умолчанию. Для индекса это означает, что последующие перестройки этого индекса не будут протоколироваться в журнале (индекс не будет генерировать информацию redo; другие индексы и сама таблица могут это делать, но такой индекс — нет). Мы можем наблюдать это с применением только что созданной таблицы T:

```

EODA@ORA12CR1> select log_mode from v$database;
LOG_MODE
-----
ARCHIVELOG
EODA@ORA12CR1> create index t_idx on t(object_name);
Index created.
Индекс создан.
EODA@ORA12CR1> variable redo number
EODA@ORA12CR1> exec :redo := get_stat_val( 'redo size' );
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.
EODA@ORA12CR1> alter index t_idx rebuild;
Index altered.
Индекс изменен.
EODA@ORA12CR1> exec dbms_output.put_line( (get_stat_val('redo size')-:redo)
|| ' bytes of redo generated...');
672264 bytes of redo generated...
PL/SQL procedure successfully completed.
Сгенерировано 672264 байтов информации redo...
Процедура PL/SQL успешно завершена.

```

На заметку! И снова этот пример выполнялся в базе данных, функционирующей в режиме ARCHIVELOG. Вы не увидите разницы в размере информации redo в базе данных, функционирующей в режиме NOARCHIVELOG, т.к. операции CREATE и REBUILD на индексе в этом режиме в журнале не протоколируются.

Когда индекс находится в режиме LOGGING (по умолчанию), его перестройка генерирует около 600 Кбайт информации redo. Тем не менее, индекс можно изменить:

```

EODA@ORA12CR1> alter index t_idx nologging;
Index altered.
Индекс изменен.
EODA@ORA12CR1> exec :redo := get_stat_val( 'redo size' );
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.
EODA@ORA12CR1> alter index t_idx rebuild;
Index altered.
Индекс изменен.
EODA@ORA12CR1> exec dbms_output.put_line( (get_stat_val('redo size')-:redo)
|| ' bytes of redo generated...');
39352 bytes of redo generated...
PL/SQL procedure successfully completed.
Сгенерировано 39352 байтов информации redo...
Процедура PL/SQL успешно завершена.

```

Теперь он генерирует всего 39 Кбайт информации redo. Но теперь этот индекс “не защищен”. Если файлы данных, в которых он хранится, повреждаются и должны быть восстановлены из резервной копии, то данные такого индекса теряются.

Понимание этого факта критически важно. Прямо сейчас индекс является невозстановимым — нам нужна резервная копия. В качестве альтернативы администратор базы данных мог бы просто создать индекс повторно, как мы можем пересоздавать его непосредственно по данным таблицы.

Заключительные соображения по поводу режима NOLOGGING

Ниже перечислены операции, которые могут выполняться в режиме NOLOGGING.

- Создание и изменение (перестройка) индексов.
- Массовые вставки в таблицу с использованием оператора INSERT в прямом режиме, такого как INSERT, доступный через подсказку /*+ APPEND */, или прямые загрузки в SQL*Loader. Табличные данные не будут генерировать информацию redo, но это будут делать все операции изменения индексов (индексы такой не протоколируемой в журнале таблицы будут генерировать информацию redo).
- Операции над данными LOB (обновления больших объектов не должны записываться в журнал).
- Создание таблиц посредством CREATE TABLE AS SELECT.
- Разнообразные операции ALTER TABLE, такие как MOVE и SPLIT.

Правильно применяемая с базой данных, работающей в режиме ARCHIVELOG, конструкция NOLOGGING может ускорить многие операции за счет значительного сокращения объема генерируемой информации redo. Предположим, что есть таблица, которую необходимо переместить из одного табличного пространства в другое. Выполнение этой операции может быть запланировано на момент непосредственно перед резервным копированием — вы должны изменить таблицу так, чтобы она стала NOLOGGING, переместить ее, перестроить индексы (также без записи в журнал) и затем снова изменить таблицу, вернув ее в режим ведения журнала. Теперь операция, которая требовала для своего выполнения X часов, возможно, выполнится за X/2 часов (хотя я не обещаю пятидесятипроцентное сокращение времени выполнения). Корректное использование этого средства предусматривает вовлечение администратора базы данных или сотрудника, отвечающего за резервное копирование и восстановление базы данных либо ведение любой запасной базы данных. Если упомянутый сотрудник не будет осведомлен о том, что вы применяете это средство, и происходит отказ носителя, то вы можете потерять данные или может быть нарушена целостность резервной базы данных. К этому следует относиться очень серьезно.

Почему не удается разместить новый журнал?

Этот вопрос мне задают постоянно. Вы получаете сообщение, предупреждающее об этом (его можно найти в журнале alert.log на сервере):

```
Thread 1 cannot allocate new log, sequence 1466
Checkpoint not complete
Current log# 3 seq# 1465 mem# 0: /.../...redo03.log
```

Потоку 1 не удается разместить новый журнал, последовательность 1466

Контрольная точка не завершена

Текущий журнал # 3 последовательность # 1465 память # 0: /.../...redo03.log

Вместо Checkpoint not complete (Контрольная точка не завершена) сообщение может содержать предупреждение Archival required (Требуется архивация), но результат будет почти тем же. Это действительно то, на что должен взглянуть администратор базы данных. Такое сообщение будет заноситься в журнал alert.log на сервере всякий раз, когда база данных пытается повторно использовать оперативный журнал redo и обнаруживает, что не может этого сделать. Подобная ситуация возникает, когда процесс DBWn еще не завершил сохранение контрольной точки данных, защищенных журналом redo, или процесс ARCn не закончил копирование файла журнала redo в архив. С точки зрения конечного пользователя в этот момент база данных *останавливается*. Она “замораживается”. Процессу DBWn или ARCn будет предоставлен приоритет для сбрасывания блоков данных на диск. После завершения установки контрольной точки или архивирования все возвращается к норме. Причина того, что база данных приостанавливает пользовательскую активность, связана с тем, что просто нет места для записи изменений, производимых пользователем. СУБД Oracle пытается повторно использовать оперативный журнальный файл redo, но из-за того, что либо этот файл потребовался для восстановления базы данных в случае отказа (Checkpoint not complete), либо архиватор еще не завершил его копирование (Archival required), Oracle приходится ждать (и конечным пользователям также), пока не появится возможность снова безопасно работать с этим файлом.

Если вы видите, что ваши сеансы тратят много времени на ожидание событий “log file switch” (“переключение журнальных файлов”), “log buffer space” (“буферное пространство журнала”) или “log file switch checkpoint or archival incomplete” (“не завершено сохранение контрольной точки или не закончена архивация”), то, скорее всего, вы столкнулись с описанной ситуацией. Вы заметите это по замедлению операций модификации базы, если размеры журнальных файлов установлены некорректно либо процессы DBWn и ARCn нуждаются в дополнительной настройке администратором базы данных или системным администратором. Я часто наблюдал эту проблему в “стартовых” базах данных, которые не были настроены. Обычно в “стартовой” базе данных размеры журналов redo установлены слишком малыми, чтобы можно было выполнить любой существенный объем работы (включая начальное построение собственного словаря базы данных). Начав загрузку информации в базу данных, вы замечаете, что первые 1000 строк проходят быстро, а затем все происходит скачкообразно: 1000 строк загружаются быстро, затем происходит задержка, потом опять загрузка строк идет быстро, снова возникает задержка и т.д. Это признаки того, что вы попали в рассматриваемую ситуацию.

Существует несколько действий, которые вы можете предпринять, чтобы решить эту проблему.

- Ускорить процесс DBWn. Попросите администратора базы данных настроить процесс DBWn путем включения асинхронного ввода-вывода, применения подчиненных процессов ввода-вывода DBWn либо использования нескольких процессов DBWn. Если при анализе ввода-вывода в системе обнаруживается, что один диск или набор дисков является “горячим”, то данные понадобятся разнести по разным дискам. Та же общая рекомендация касается и процесса ARCn. Преимущество такого решения состоит в том, что вы получаете здесь “кое-что просто так” — увеличение производительности без действительных изменений логики/структур/кода. У этого подхода нет отрицательных сторон.

- Добавить дополнительные журнальные файлы redo. В некоторых случаях это отсрочит появление предупреждения Checkpoint not complete, а со временем отсрочит его так надолго, что, возможно, оно никогда не появится (поскольку вы предоставили процессу DBWn достаточное пространство для создания контрольных точек). То же самое применимо и к предупреждению Archival required. Преимуществом такого подхода является устранение “пауз” в работе системы. Недостаток связан с потреблением большего дискового пространства, но преимущество значительно перевешивает этот недостаток.
- Создать журнальные файлы заново с большими размерами. Это увеличит период времени между моментом наполнения оперативного журнала redo и моментом, когда возникнет необходимость в его повторном использовании. То же самое касается предупреждения Archival required, если утилизация журнального файла redo носит “пульсирующий” характер. Если у вас есть период массивированной генерации данных redo (ночные загрузки, пакетные процессы), за которыми следуют периоды относительного затишья, то наличие журнальных файлов большего размера предоставит процессу ARCn достаточное время, чтобы наверстать невыполненную работу в периоды затишья. Преимущества и недостатки здесь идентичны подходу с добавлением дополнительных журнальных файлов. Кроме того, данный прием может перенести момент сохранения контрольной точки на более позднее время, т.к. это действие случается при каждом переключении журналов (как минимум), а переключения журналов будут происходить реже.
- Обеспечить более частое и непрерывное сохранение контрольных точек. Применяйте кеш буферов блоков меньшего размера (не совсем желательно) или различные настройки параметров, такие как FAST_START_MTTR_TARGET, LOG_CHECKPOINT_INTERVAL и LOG_CHECKPOINT_TIMEOUT. Это заставит процесс DBWn выталкивать грязные блоки более часто. Преимущество такого подхода связано с сокращением времени восстановления после отказа. В оперативных журналах redo всегда будет оставаться меньше работы, подлежащей применению. Недостаток заключается в том, что блоки могут записываться на диск более часто в случае их частой модификации. Буферный кеш будет не настолько эффективным, как мог бы быть, и это может помешать работе механизма очистки блоков, который обсуждается в следующем разделе.

Подход будет выбираться в зависимости от конкретных обстоятельств. Это то, что должно быть установлено на уровне базы данных, принимая во внимание целый экземпляр.

Очистка блоков

В настоящем разделе мы обсудим *очистку блоков* (block cleanouts), или удаление информации, связанной с блокировками, из блоков базы данных, которые мы модифицировали. Эта концепция будет важна для понимания печально известной ошибки ORA-01555: snapshot too old (ORA-1555: устаревший снимок), о которой пойдет речь в следующем разделе.

В главе 6 говорилось о блокировках данных и об управлении ними. Было показано, что на самом деле они являются атрибутами данных, хранящимися в заголов-

ке блока. Побочный эффект от этого заключается в том, что при доступе к блоку в следующий раз может потребоваться его очистить — другими словами, удалить информацию о транзакции. Такое действие генерирует информацию redo и приводит к тому, что блок становится *грязным*, если он еще не был таковым, а это означает, что простой оператор SELECT *может генерировать информацию redo* и вызвать запись множества блоков на диск в следующей контрольной точке. Однако в большинстве нормальных случаев подобное не происходит. Если в вашей системе преобладают транзакции от небольшого до среднего размера (OLTP) или вы располагаете хранилищем данных, которое производит прямые загрузки либо использует пакет DBMS_STATS для анализа таблиц после операций загрузки, то вы обнаружите, что блоки обычно очищаются автоматически. Как было указано в разделе “Что делает оператор COMMIT?” ранее в главе, одним из шагов, выполняемых во время обработки оператора COMMIT, является повторное посещение определенных блоков, если они по-прежнему находятся в области SGA и доступны (никто другой не модифицирует их), и проведение их очистки. Это действие называется *очисткой при фиксации*, и его сущность заключается в очистке модифицированных блоков от информации, касающейся транзакции. В оптимальном случае наш оператор COMMIT может очистить блоки, так что последующему оператору SELECT (чтению) это делать не придется. Только оператор UPDATE для этого блока будет действительно очищать остаточную информацию о транзакции, и поскольку UPDATE уже генерирует информацию redo, очистка остается незаметной.

Мы можем принудительно *запретить* очистку и таким образом понаблюдать за ее побочными эффектами, благодаря пониманию того, как она работает. В список фиксации, ассоциированный с транзакцией, Oracle будет записывать перечни блоков, которые мы модифицировали. Каждый из этих перечней имеет длину в 20 блоков, и Oracle выделит столько таких списков, сколько будет необходимо. Когда сумма модифицированных блоков превышает 10% от размера кеша буферов блоков, Oracle прекращает выделение новых списков. Например, если буферный кеш настроен на кеширование 3000 блоков, то Oracle будет поддерживать список длиной до 300 блоков (10% от 3000). Получив COMMIT, Oracle обработает каждый из этих списков из 20 указателей на блоки, и если блок еще доступен, то выполнит быструю очистку. Таким образом, пока количество модифицированных блоков не превышает 10% от числа блоков в кеше и наши блоки по-прежнему находятся в кеше и доступны, Oracle будет очищать их при поступлении COMMIT. В противном случае они пропускаются (т.е. не будут очищены).

Зная все это, мы можем создать искусственные условия, чтобы посмотреть, каким образом работает очистка. Я установил параметр DB_CACHE_SIZE в минимальное значение — 16 Мбайт, чего достаточно для того, чтобы вместить 2048 блоков по 8 Кбайт (размер блока у меня составляет 8 Кбайт). Далее я создам таблицу так, чтобы любая строка занимала строго один блок — не будет ни одного блока, вмещающего две строки. После этого я заполню эту таблицу 10 000 записями и выдам оператор COMMIT. Нам известно, что количество 10 000 блоков намного превосходит 10% от 2048, поэтому база данных не способна очистить эти грязные блоки при фиксации — большая их часть уже даже не будет находиться в буферном кеше. Я измеряю объем информации redo, сгенерированной до сих пор, запущу оператор SELECT, который посетит каждый блок, и затем измеряю объем информации redo, сгенерированной этим оператором SELECT.

На заметку! Чтобы рассматриваемый пример был воспроизводимым и предсказуемым, понадобится отключить автоматическое управление памятью SGA. Если оно включено, то есть шанс, что база данных увеличит размер буферного кеша, сводя на нет всю “математику”, задействованную в примере.

К удивлению многих, оператор SELECT генерирует информацию redo. Более того, он также делает “грязными” модифицированные блоки, заставляя процесс DBWn записывать их заново. Это происходит из-за очистки блоков. Далее я выдам оператор SELECT еще раз и увижу, что никакой информации redo на этот раз не генерируется. Ситуация вполне ожидаема, т.к. в этот момент все блоки “чистые”. Мы начнем с создания таблицы:

```
EODA@ORA12CR1> create table t
2  ( id number primary key,
3    x char(2000),
4    y char(2000),
5    z char(2000)
6  )
7  /
```

Table created.

Таблица создана.

```
EODA@ORA12CR1> exec dbms_stats.set_table_stats( user, 'T',
                                                    numrows=>10000, numblks=>10000 );
```

PL/SQL procedure successfully completed.

Процедура PL/SQL успешно завершена.

Для получения статистических сведений по таблице применяется пакет DBMS_STATS, что позволяет избежать любых побочных эффектов от полного разбора позднее (СУБД Oracle склонна сканировать объекты, не имеющие статистики, во время полного разбора, и этот побочный эффект повлияет на пример). Итак, имеется таблица с одной строкой на блок (в базе данных с размером блока в 8 Кбайт). Далее мы рассмотрим блок кода, который будет выполнен в отношении этой таблицы:

```
EODA@ORA12CR1> declare
2    l_rec t%rowtype;
3  begin
4    for i in 1 .. 10000
5    loop
6      select * into l_rec from t where id=i;
7    end loop;
8  end;
9  /
```

declare

*

ERROR at line 1:

ORA-01403: no data found

ORA-06512: at line 6

ОШИБКА в строке 1:

ORA-01403: данные не найдены

ORA-06512: в строке 6

Блок работать отказался, но это нормально — мы знали, что так будет, поскольку в таблице пока еще нет данных. Я запустил его просто для выполнения полного разбора SQL и PL/SQL, так что когда мы запустим его позднее, не придется беспокоиться по поводу нежелательного учета побочных эффектов от полного разбора. Теперь все готово для загрузки данных в таблицу и их фиксации:

```
EODA@ORA12CR1> insert into t
  2  select rownum, 'x', 'y', 'z'
  3  from all_objects
  4  where rownum <= 10000;
10000 rows created.
10000 строк создано.
```

```
EODA@ORA12CR1> commit;
Commit complete.
Фиксация завершена.
```

И, наконец, можно измерить объем информации redo, сгенерированной во время первого чтения данных:

```
EODA@ORA12CR1> variable redo number
EODA@ORA12CR1> exec :redo := get_stat_val( 'redo size' );
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.
```

```
EODA@ORA12CR1> declare
  2  l_rec t%rowtype;
  3  begin
  4  for i in 1 .. 10000
  5  loop
  6  select * into l_rec from t where id=i;
  7  end loop;
  8  end;
  9  /
```

```
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.
```

```
EODA@ORA12CR1> exec dbms_output.put_line( (get_stat_val('redo size')-:redo)
  || ' bytes of redo generated...');
802632 bytes of redo generated...
```

```
PL/SQL procedure successfully completed.
Сгенерировано 802632 байтов информации redo...
Процедура PL/SQL успешно завершена.
```

Таким образом, оператор SELECT сгенерировал около 802 Кбайт информации redo во время своей работы. Это представляет заголовки блоков, модифицированные во время чтения индекса первичного ключа и последующего чтения таблицы T. Процесс DBWn запишет эти модифицированные блоки обратно на диск в какой-то момент в будущем (на самом деле, поскольку таблица не помещается в кеш, мы знаем, что DBWn уже записал, по крайней мере, некоторые из них). Запустим запрос снова:

```
EODA@ORA12CR1> exec :redo := get_stat_val( 'redo size' );
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.
```



```

EODA@ORA12CR1> declare
2     l_rec t%rowtype;
3 begin
4     for i in 1 .. 10000
5     loop
6         select * into l_rec from t where id=i;
7     end loop;
8 end;
9 /

```

PL/SQL procedure successfully completed.

Процедура PL/SQL успешно завершена.

```

EODA@ORA12CR1> exec dbms_output.put_line( (get_stat_val('redo size')-:redo)
|| ' bytes of redo generated...');

```

0 bytes of redo generated...

PL/SQL procedure successfully completed.

Сгенерировано 0 байтов информации redo...

Процедура PL/SQL успешно завершена.

Легко заметить, что информация redo не генерировалась — все блоки чисты.

Если перезапустить предыдущий пример с буферным кешем, настроенным так, чтобы удерживать немного больше 100 000 блоков, то обнаружится, что любой оператор SELECT будет генерировать малый или нулевой объем информации redo — не приходится очищать грязные блоки во время выполнения любого из операторов SELECT. Это связано с тем, что свыше 10 000 модифицированных нами блоков (помните, что индекс также был модифицирован) спокойно умещаются в 10% буферного кеша, и мы являемся единственным пользователем. Никто другой не затрагивает наши данные и никто другой не вызывает их сбрасывание на диск или обращение к этим блокам. В реальной системе будет нормой то, что, по крайней мере, некоторые из блоков иногда очищаться не будут.

Такое поведение оказывает наибольшее влияние после выполнения крупного оператора INSERT (как только что демонстрировалось), UPDATE или DELETE — когда затрагивается много блоков в базе данных (все, что превышает 10% от размера кеша). Вы заметите, что первый запрос, затрагивающий блок после этого, сгенерирует мало информации redo и сделает блок грязным, возможно, вызывая его перезапись, если процесс DBWn уже вытолкнул его на диск или экземпляр был остановлен, что привело к полной очистке буферного кеша. В такой ситуации сделать можно не особенно многое. Это нормально и ожидаемо. Если Oracle не выполнит эту отложенную очистку блока, то оператору COMMIT может потребоваться относительно много времени для обработки транзакции. Оператору COMMIT придется повторно посетить каждый блок, возможно, снова читая некоторые из них с диска (они уже могут быть вытолкнуты).

Если вы не осведомлены об очистке блоков и особенностях ее работы, то для вас это будет одной из тех мистических вещей, которые, как кажется, происходят безо всякой причины. Например, предположим, что вы выполняете обновление большого объема данных и затем фиксацию. После этого вы запускаете запрос по этим данным, чтобы проверить результаты. Ситуация выглядит так, что запрос генерирует громадный объем операций записи и информации redo. Подобное кажется невозможным, если не знать об очистке блоков; так было со мной, когда я впервые с этим столкнулся. Вы зовете кого-то, чтобы вместе понаблюдать за этим поведени-

ем, но оно не воспроизводимо, т.к. при следующем запросе блоки уже “чистые”. Вы просто списываете это на одну из загадок базы данных — загадок, которые случаются, когда вы одни.

В системах OLTP вы, возможно, никогда не увидите очистку блоков, поскольку такие системы характеризуются небольшими и короткими транзакциями, которые затрагивают лишь несколько блоков. Проектные решения для таких систем построены так, что большинство транзакций являются очень краткими и изящными. Модифицируйте пару блоков — и все они будут очищены. В хранилище данных, где выполняются массовые обновления данных после загрузки, очистки блоков могут оказаться фактором, который должен быть учтен при проектировании. Некоторые операции будут создавать данные в “чистых” блоках. Например, данные, сгенерированные оператором `CREATE TABLE AS SELECT`, данные, загруженные в прямом режиме, данные, вставленные в прямом режиме (с использованием подсказки `/*+ APPEND */`), будут создавать чистые блоки. Оператор `UPDATE`, нормальный оператор `INSERT` или оператор `DELETE` могут создавать блоки, которые нуждаются в очистке при первом чтении. Это действительно может оказать влияние, если обработка состоит из следующих действий:

- массовая загрузка большого объема новых данных в хранилище;
- выполнение обновлений всех только что загруженных данных (производя блоки, которым необходима очистка);
- предоставление возможности запрашивать данные.

Вы должны понимать, что первый запрос, касающийся этих данных, потребует некоторой дополнительной обработки, если блоки нуждаются в очистке. Осознавая это, вы сами должны “затронуть” данные после обновления. Вы только что загрузили или модифицировали крупный объем данных — их необходимо, по меньшей мере, проанализировать. Возможно, понадобится построить какие-то отчеты, чтобы проверить загруженные данные. В результате блоки очистятся, и следующему запросу не придется делать это. Более того, поскольку произошла массовая загрузка большого объема данных, в любом случае нужно обновить статистику. Запуск утилиты `DBMS_STATS` для сбора статистики может также очистить все блоки, т.к. он просто применяет оператор `SQL` для запрашивания информации и в ходе его выполнения естественным образом осуществляется очистка.

Конкуренция за журнал

Конкуренция за журнал, как и сообщение `cannot allocate new log` (не удается разместить новый журнал), требует вмешательства администратора базы данных, обычно в сотрудничестве с системным администратором. Однако это может обнаружить и разработчик, если администратор базы данных поблизости не находится.

Если вы сталкиваетесь с конкуренцией за журнал, то можете наблюдать в отчете `Statspack` большое время ожидания при наступлении события “`log file sync`” (“синхронизация журнальных файлов”) и длительное время записи при событии “`log file parallel write`” (“параллельная запись в журнальные файлы”). Если вы видите подобное, то это может быть конкуренция за журналы redo; запись в них производится недостаточно быстро. Причин может быть много. Одной из причин, связанных с приложением (которую не может устранить администратор базы данных, но должен

устранить разработчик), является слишком частая фиксация, например, фиксация внутри цикла, в котором выполняется оператор `INSERT`. Как демонстрировалось в разделе “Что делает оператор `COMMIT`?”, слишком частая фиксация, кроме того, что сама по себе считается плохой практикой программирования, представляет собой верный путь к привнесению многочисленных ожиданий синхронизации журнальных файлов. Предполагая, что размеры всех транзакций установлены корректно (вы не выполняете фиксацию чаще, чем продиктовано бизнес-правилами), самыми распространенными причинами ожиданий синхронизации журнальных файлов будут перечисленные ниже.

- Размещение журналов redo на медленном устройстве. Производительность дисков явно недостаточна. Пора приобрести более быстрые диски.
- Размещение журналов redo на одном устройстве с другими файлами, к которым производится частый доступ. Журналы redo спроектированы так, что запись в них должна осуществляться последовательно, а размещаться они должны на выделенных устройствах. Если другие компоненты вашей системы — даже другие компоненты Oracle — пытаются выполнять чтение и запись на этом устройстве в то же самое время, что и процесс LGWR, то вы столкнетесь с определенной долей конкуренции. По возможности необходимо обеспечить процессу LGWR монопольный доступ к устройству.
- Монтирование устройств с журнальными файлами в буферизованном режиме. Здесь вы используете “готовую” файловую систему (не диски RAW). Операционная система буферизует данные, и база данных также буферизует их (посредством буфера журнала redo). Двойная буферизация замедляет работу. По возможности монтируйте устройства в “прямом” режиме. Способ зависит от операционной системы и устройства, но обычно это возможно.
- Размещение журналов redo на устройствах, функционирующих по медленной технологии, такой как RAID-5. Устройство RAID-5 великолепно работает при выполнении чтения, но в целом отвратительно, когда дело доходит до записи. Как мы видели ранее, рассматривая, что происходит во время обработки оператора `COMMIT`, приходится ждать, пока процесс LGWR обеспечит запись данных на диск. Применение любой технологии, которая замедляет этот процесс — плохая идея.

Если это вообще возможно, следует использовать минимум пять выделенных устройств для ведения журналов и оптимально шесть таких устройств, чтобы обеспечить также зеркальное отражение архивов. В наши дни с дисками объемом 200 Гбайт, 300 Гбайт, 1 Тбайт и более это становится сложнее, но если удастся отыскать четыре самых маленьких и быстрых диска и один или два больших, то вы сможете добиться оптимальной работы процессов LGWR и ARCn. При этом диски понадобятся разделить на три группы (рис. 9.5):

- группа 1 журналов redo — диски 1 и 3;
- группа 2 журналов redo — диски 2 и 4;
- архив — диск 5 и дополнительно диск 6 (большие диски).

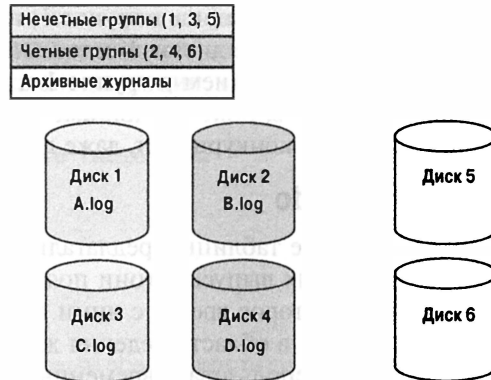


Рис. 9.5. Оптимальная конфигурация журналов redo

Вы должны поместить группу 1 журналов redo с членами А и В на диски 1 и 3, а группу 2 журналов redo с членами С и D — на диски 2 и 4. Если есть группы 3, 4 и т.д. журналов redo, то они должны соответствующим образом попасть на нечетные и четные группы дисков. Результатом будет то, что процесс LGWR, в то время когда база данных использует группу 1, будет параллельно выполнять запись на диски 1 и 3. Когда эта группа заполнится, процесс LGWR перейдет к дискам 2 и 4. Когда они заполнятся, LGWR возвратится к дискам 1 и 3. Тем временем процесс ARCn будет обрабатывать полные оперативные журналы redo и записывать их на диски 5 и 6 (большие диски). Совокупный эффект заключается в том, что процессы ARCn и LGWR никогда не будут читать диск, на который одновременно производится запись, или записывать на диск, который в данный момент читается, поэтому никакой конкуренции не возникает (рис. 9.6).

| |
|---------------------------|
| Нечетные группы (1, 3, 5) |
| Четные группы (2, 4, 6) |
| Архивные журналы |

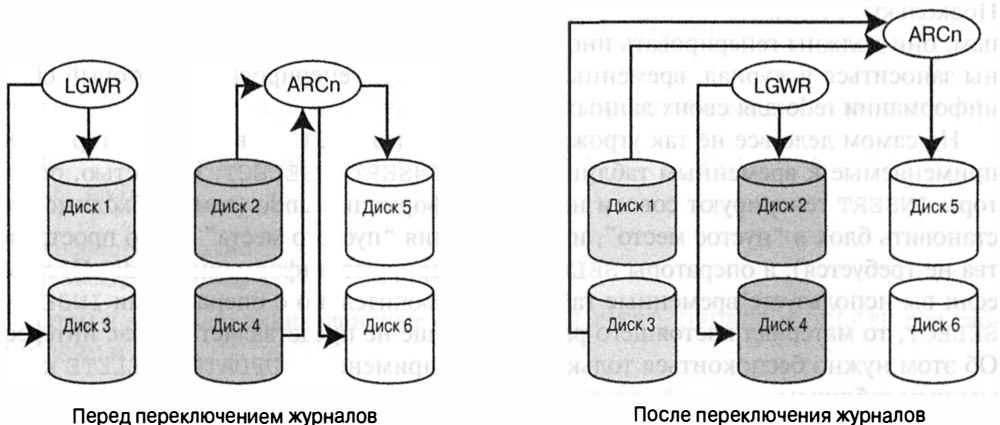


Рис. 9.6. Ведение журналов redo

Таким образом, когда LGWR выполняет запись в группе 1, процесс ARCn производит чтение в группе 2 и запись на архивные диски. Когда LGWR осуществляет запись в группе 2, процесс ARCn занимается чтением в группе 1 и записью на архивные диски. При такой организации процессы LGWR и ARCn имеют собственные выделенные устройства и никогда ни с кем не конкурируют, даже друг с другом.

Временные таблицы и redo/undo

К настоящему моменту временные таблицы предлагались в качестве функционального средства Oracle в нескольких выпусках (они появились в Oracle8i версии 8.1.5). Хотя они уже существуют некоторое время, с ними по-прежнему связана определенная неразбериха, в частности, в области ведения журналов. В главе 10 речь пойдет о том, как и почему можно использовать временные таблицы. Здесь же мы исследуем только один вопрос: что делают временные таблицы в плане записи в журнал изменений?

В Oracle 12c обработка информации undo для временных таблиц значительно улучшилась. По этой причине данная тема разбита на два раздела: “До версии Oracle 12c” и “Начиная с версии Oracle 12c”.

До версии Oracle 12c

Временные таблицы не генерируют информацию redo для своих блоков. Следовательно, операции над временными таблицами являются невозстановимыми. Когда вы модифицируете блок во временной таблице, никакие записи об этом изменении в файлы журналов redo не заносятся. Однако временные таблицы генерируют информацию undo, *и эта информация undo записывается в журнал*. Таким образом, временные таблицы все-таки генерируют некоторую информацию redo. На первый взгляд, в этом совершенно нет смысла: зачем им генерировать данные undo? Причина в том, что вы можете выполнить откат к точке сохранения внутри транзакции. Вы можете удалить 50 последних вставок во временную таблицу, оставив первые 50. Временные таблицы могут иметь ограничения и все остальное, что присуще нормальным таблицам. При вставке 500 строк одним оператором INSERT может произойти отказ на последней строке, что потребует отката этого оператора. Поскольку временные таблицы ведут себя в основном подобно нормальным таблицам, они должны генерировать информацию undo. А поскольку данные undo должны заноситься в журнал, временные таблицы будут генерировать некоторый объем информации redo для своих данных undo.

На самом деле все не так угрожающе, как выглядит. Основные операторы SQL, применяемые к временным таблицам — это INSERT и SELECT. К счастью, операторы INSERT генерируют совсем немного информации undo (вам необходимо восстановить блок в “пустое место”, и для хранения “пустого места” много пространства не требуется), а операторы SELECT не генерируют информацию undo. Поэтому, если вы используете временные таблицы исключительно с операторами INSERT и SELECT, то материал настоящего раздела вообще не представляет для вас интереса. Об этом нужно беспокоиться только в случае применения UPDATE и DELETE к временным таблицам.

Я подготовил небольшой тест для демонстрации объема информации redo, генерируемой при работе с временными таблицами, который, следовательно, является показателем объема данных undo, сгенерированного для временных таблиц, т.к. для

них в журнал записывается только информация undo. Для этой цели я создам идентично сконфигурированные постоянную и временную таблицы и затем выполню с каждой из них те же самые операции, каждый раз измеряя объем сгенерированной информации redo. Вот эти таблицы:

```
EODA@ORA11CR1> create table perm
  2  ( x char(2000) ,
  3    y char(2000) ,
  4    z char(2000) )
  5  /
```

Table created.

Таблица создана.

```
EODA@ORA11CR1> create global temporary table temp
  2  ( x char(2000) ,
  3    y char(2000) ,
  4    z char(2000) )
  5  on commit preserve rows
  6  /
```

Table created.

Таблица создана.

Я написал небольшую хранимую процедуру, которая позволит выполнять произвольный оператор SQL, и будет сообщать объем информации redo, сгенерированной этим оператором SQL. Я использую эту процедуру для выполнения операторов INSERT, UPDATE и DELETE в отношении как постоянной, так и временной таблицы:

```
EODA@ORA11CR1> create or replace procedure do_sql( p_sql in varchar2 )
  2  as
  3    l_start_redo    number;
  4    l_redo          number;
  5  begin
  6    l_start_redo := get_stat_val( 'redo size' );
  7
  8    execute immediate p_sql;
  9    commit;
 10
 11    l_redo := get_stat_val( 'redo size' ) - l_start_redo;
 12
 13    dbms_output.put_line
 14      ( to_char(l_redo, '99,999,999') || ' bytes of redo generated for "'
 15      ||
 16      substr( replace( p_sql, chr(10), ' '), 1, 25 ) || '"...' );
 17  end;
```

Procedure created.

Процедура создана.

Затем я запускаю эквивалентные операторы INSERT, UPDATE и DELETE для таблиц PERM и TEMP:

```
EODA@ORA11CR1> set serveroutput on format wrapped
EODA@ORA11CR1> begin
  2    do_sql( 'insert into perm
  3            select 1,1,1
```

```

4          from all_objects
5          where rownum <= 500' );
6
7      do_sql( 'insert into temp
8              select 1,1,1
9              from all_objects
10             where rownum <= 500' );
11      dbms_output.new_line;
12
13      do_sql( 'update perm set x = 2' );
14      do_sql( 'update temp set x = 2' );
15      dbms_output.new_line;
16
17      do_sql( 'delete from perm' );
18      do_sql( 'delete from temp' );
19  end;
20 /

```

3,313,088 bytes of redo generated for "insert into perm" ...
 72,584 bytes of redo generated for "insert into temp" ...
 3,313,088 байтов информации redo сгенерировано для "insert into perm" ...
 72,584 байтов информации redo сгенерировано для "insert into temp" ...
 3,268,384 bytes of redo generated for "update perm set x = 2" ...
 1,946,432 bytes of redo generated for "update temp set x = 2" ...
 3,268,384 байтов информации redo сгенерировано для "update perm set x = 2" ...
 1,946,432 байтов информации redo сгенерировано для "update temp set x = 2" ...
 3,245,112 bytes of redo generated for "delete from perm" ...
 3,224,460 bytes of redo generated for "delete from temp" ...
 3,245,112 байтов информации redo сгенерировано для "delete from perm" ...
 3,224,460 байтов информации redo сгенерировано для "delete from temp" ...

PL/SQL procedure successfully completed.

Процедура PL/SQL успешно завершена.

На основе приведенного отчета можно сделать следующие выводы.

- Оператор INSERT для “реальной” таблицы генерирует большой объем информации redo. Для временной таблицы информация redo практически не генерировалась. В этом есть смысл — для операторов INSERT генерируется очень мало данных undo и для временных таблиц в журнал записываются только данные undo.
- Оператор UPDATE для реальной таблицы генерирует примерно вдвое больше информации redo, чем для временной таблицы. Опять-таки, это имеет смысл. Должна быть сохранена примерно половина данных UPDATE — образ “до изменения”. Образ “после изменения” (информация redo) для временных таблиц сохраняться не должен.
- Оператор DELETE требует примерно того же объема информации redo. Это имеет смысл, потому что информация undo для DELETE имеет большой размер, но информации redo для модифицированных блоков очень мало. В результате выполнение DELETE в отношении временной таблицы происходит очень похоже на выполнение DELETE в постоянной таблице.

На заметку! Если вы видите, что временная таблица генерирует для оператора INSERT больше информации redo, чем постоянная таблица, это значит, что вы наблюдаете проблему в программном продукте, которая была устранена в Oracle 9.2.0.6 и в выпусках исправлений 10.1.0.4 и последующих версий.

Таким образом, можно сделать следующие обобщения относительно операций DML во временных таблицах.

- Оператор INSERT будет генерировать минимум информации undo/redo.
- Оператор UPDATE будет генерировать около половины объема информации redo по сравнению с постоянными таблицами.
- Оператор DELETE будет генерировать такой же объем информации redo, как и в постоянных таблицах.

С предпоследним оператором связаны значительные исключения. Например, если я выполняю обновление столбца, содержащего значение NULL, посредством 2000 байтов данных, то будет сгенерировано очень мало информации undo. Такой оператор UPDATE ведет себя подобно INSERT. С другой стороны, если я обновляю столбец с 2000 байтов значением NULL, то такой оператор UPDATE в плане генерации данных redo ведет себя как DELETE. В среднем можно ожидать, что оператор UPDATE в отношении временной таблицы будет генерировать около половины информации undo/redo по сравнению со случаем постоянной таблицы.

В дополнение вы должны принимать во внимание любые индексы на временных таблицах. Модификации индексов также сгенерируют информацию undo, которая, в свою очередь, генерирует информацию redo. Перезапустим приведенный выше пример со следующими двумя индексами:

```
EODA@ORA11CR1> create index perm_idx on perm(x);
Index created.
Индекс создан.

EODA@ORA11CR1> create index temp_idx on temp(x);
Index created.
Индекс создан.
```

Видно, что была сгенерирована информация redo (для краткости весь код из предыдущего примера здесь не повторяется):

```
...
19 end;
20 /

11,735,576 bytes of redo generated for "insert into perm      "...
3,351,864 bytes of redo generated for "insert into temp      "...
11,735,576 байтов информации redo сгенерировано для "insert into perm      "...
3,351,864 байтов информации redo сгенерировано для "insert into temp      "...

9,257,748 bytes of redo generated for "update perm set x = 2"...
5,465,868 bytes of redo generated for "update temp set x = 2"...
9,257,748 байтов информации redo сгенерировано для "update perm set x = 2"...
5,465,868 байтов информации redo сгенерировано для "update temp set x = 2"...

4,434,992 bytes of redo generated for "delete from perm"...
4,371,620 bytes of redo generated for "delete from temp"...
4,434,992 байтов информации redo сгенерировано для "delete from perm"...
4,371,620 байтов информации redo сгенерировано для "delete from temp"...
```


PL/SQL procedure successfully completed.

Процедура PL/SQL успешно завершена.

Числа подтверждают то, что мы видели ранее, но здесь видно, что индекс определенно увеличил объем сгенерированной информации redo. Оператор INSERT для глобальной временной таблицы, ранее почти не генерировавший информацию redo, теперь генерирует 3,3 Мбайт данных redo. Весь этот дополнительный объем redo связан с информацией undo, созданной для обслуживания индексов.

На заметку! Это несколько преувеличенный пример. Индекс был построен на столбце CHAR(2000); ключ индекса намного больше, чем можно встретить в реальности. Обычно не стоит ждать настолько значительного дополнительного объема redo.

В общем случае оценка объема генерируемой информации redo диктуется здравым смыслом. Если выполняемая вами операция вызывает создание данных undo, то определите, насколько легко или трудно будет обратить (отменить) результат этой операции. Если вы вставляете 2000 байтов, отменить это легко. Вы просто возвращаетесь к состоянию 0 байтов. Если же вы удаляете 2000 байтов, то обратной операцией будет вставка 2000 байтов. В этом случае объем информации redo является существенным.

Взяв на вооружение эти знания, вы будете избегать удаления во временных таблицах. Вы можете использовать TRUNCATE (разумеется, памятуя о том, что этот оператор DDL выполняет фиксацию текущей транзакции, а в Oracle9i и предшествующих версиях делает недействительными курсоры) либо просто позволить временным таблицам очищать себя автоматически после COMMIT или по завершении сеанса. Все эти методы не генерируют информацию undo и, следовательно, также информацию redo. Вы должны стараться избегать обновления временной таблицы, если только в этом нет настоятельной необходимости. Вы должны применять временные таблицы по большей части как место, предназначенное для вставки и выборки. При таком подходе вы обеспечите оптимальное использование их уникальной способности не генерировать информацию redo.

Начиная с версии Oracle 12c

В предыдущем разделе вы видели, что при выдаче операторов INSERT, UPDATE и DELETE во временной таблице информация undo для этих изменений записывается в табличное пространство отмены и, в свою очередь, будет генерировать данные redo. С выходом версии Oracle 12c вы можете инструктировать Oracle о сохранении информации undo для временной таблицы во временном табличном пространстве через параметр TEMP_UNDO_ENABLED. При модификации блоков во временном табличном пространстве информация redo не генерируется. Таким образом, когда параметр TEMP_UNDO_ENABLED установлен в TRUE, любая операция DML, выполняемая в отношении временной таблицы, будет генерировать очень мало или вообще никакой информации redo.

На заметку! По умолчанию параметр TEMP_UNDO_ENABLED установлен в FALSE. Поэтому, если только не сконфигурировано по-другому, временные таблицы в Oracle 12c будут генерировать такой же объем данных redo, как и в предшествующих выпусках.

Параметр `TEMP_UNDO_ENABLED` может быть установлен на уровне сеанса или системы. Вот пример его установки в `TRUE` на уровне сеанса:

```
EODA@ORA12CR1> alter session set temp_undo_enabled=true;
```

После включения для сеанса любые модификации данных во временной таблице, предпринятые в этом сеансе, будут иметь информацию `undo`, занесенную во временное табличное пространство. Любые модификации в постоянных таблицах по-прежнему будут иметь данные `undo`, сохраненные в табличном пространстве `undo`. Чтобы оценить влияние этого, я повторно выполняю тот же самый код (из раздела “До версии Oracle 12c”), который отображает объем информации `redo`, сгенерированной при выдаче операторов в отношении постоянной и временной таблиц; единственное отличие — параметр `TEMP_UNDO_ENABLED` теперь установлен в `TRUE`. Ниже показан вывод:

```
3,312,148 bytes of redo generated for "insert into perm      "...
      376 bytes of redo generated for "insert into temp      "...
3,312,148 байтов информации redo сгенерировано для "insert into perm  "...
      376 байтов информации redo сгенерировано для "insert into temp  "...

2,203,788 bytes of redo generated for "update perm set x = 2"...
      376 bytes of redo generated for "update temp set x = 2"...
2,203,788 байтов информации redo сгенерировано для "update perm set x = 2"...
      376 байтов информации redo сгенерировано для "update temp set x = 2"...

3,243,412 bytes of redo generated for "delete from perm"...
      376 bytes of redo generated for "delete from temp"...
3,243,412 байтов информации redo сгенерировано для "delete from perm"...
      376 байтов информации redo сгенерировано для "delete from temp"...
```

Результаты разительно отличаются: объем данных `redo`, сгенерированных операторами `INSERT`, `UPDATE` и `DELETE` во временной таблице, совершенно незначителен. Для сред, где выполняются крупные пакетные операции, которые запускают транзакции к временным таблицам, можно ожидать значительного снижения объема сгенерированной информации `redo`.

На заметку! Вас может интересовать, по какой причине было сгенерировано 376 байтов информации `redo` в предыдущем примере. Так как процессы потребляют пространство в базе данных, Oracle выполняет некоторые внутренние служебные действия. Получаемые в результате них изменения записываются в словарь данных, что, в свою очередь, приводит к генерации определенной информации `redo` и `undo`.

Начиная с версии Oracle 12c, в конфигурации Oracle Active Data Guard (Активный защитник данных Oracle) вы можете выдавать операторы DML напрямую во временную таблицу, которая существует в резервной базе данных. Чтобы просмотреть объем информации `redo`, сгенерированной для временной таблицы в резервной базе данных, можно запустить тот же самый код (из раздела “До версии Oracle 12c”) в отношении резервной базы данных. Единственное отличие от операторов, инициирующих транзакции в постоянных таблицах, должно исчезнуть (поскольку невозможно выполнить операцию DML над постоянной таблицей в резервной базе данных). Вывод показывает, что было сгенерировано 0 байтов информации `redo`:

```

0 bytes of redo generated for "insert into temp      "...
0 байтов информации redo сгенерировано для "insert into temp      "...
0 bytes of redo generated for "update temp set x = 2"...
0 байтов информации redo сгенерировано для "update temp set x = 2"...
0 bytes of redo generated for "delete from temp"...
0 байтов информации redo сгенерировано для "delete from temp"...

```

На заметку! В резервной базе данных нет необходимости устанавливать TEMP_UNDO_ENABLED. Это объясняется тем, что в резервной базе данных Oracle Active Data Guard генерация информации undo для временной таблицы всегда включена.

Глобальные временные таблицы часто применяются при построении отчетов, скажем, для генерации и сохранения промежуточных результатов запросов. Дополнение Oracle Active Data Guard обычно используется для разгрузки приложений, формирующих отчеты, в пользу резервной базы данных. Связывание временных таблиц с Oracle Active Data Guard позволяет получить в свое распоряжение более мощный инструмент для удовлетворения существующих требований к построению отчетов.

Исследование undo

Мы уже много говорили на тему сегментов undo. Мы видели, как они применяются во время восстановления, каким образом взаимодействуют с журналами redo и как используются для согласованного, неблокирующего чтения данных. В этом разделе мы рассмотрим наиболее часто возникающие проблемы, связанные с сегментами undo.

Основное внимание будет уделено ошибке ORA-01555: snapshot too old (ORA-1555: устаревший снимок), т.к. эта единственная проблема вызывает больше путаницы, чем любой другой вопрос из тематики, касающейся баз данных. Однако прежде чем приступить к этому, мы исследуем одну из двух проблем, имеющих отношение к undo: вопрос о том, какой тип операций DML генерирует максимальный и минимальный объем информации undo. (Вы уже можете быть в состоянии самостоятельно ответить на этот вопрос, основываясь на предшествующих примерах с временными таблицами.)

Что генерирует максимальный и минимальный объем информации undo?

Такой вопрос задается часто, но ответить на него легко. Присутствие индексов (или тот факт, что таблица является индекс-таблицей (index-organized table)), может радикальным образом повлиять на объем информации undo, поскольку индексы — это сложные структуры данных, которые могут генерировать крупные объемы undo.

Тем не менее, оператор INSERT в общем случае будет генерировать минимальный объем информации undo, поскольку все, что Oracle нужно для этого записать — идентификатор строки (rowid), подлежащей “удалению”. Оператор UPDATE занимает второе место (в большинстве случаев). Все, что необходимо сохранить для будущей отмены — это измененные байты. Чаще всего UPDATE изменяет какую-то небольшую часть всех данных строки. Таким образом, внутри undo должна запомин-

наться только эта небольшая часть строки. Многие из приведенных ранее примеров идут вразрез с этим эмпирическим правилом, но так происходит потому, что они обновляют крупные строки фиксированного размера и обновляют строки целиком. Однако для оператора UPDATE гораздо более типично обновлять лишь небольшую часть всей строки. Оператор DELETE в общем случае будет генерировать максимальный объем информации undo. Для DELETE в сегмент undo должен быть записан предыдущий образ целой строки. С точки зрения генерации redo этот факт был продемонстрирован в рассмотренном выше примере с временной таблицей. Оператор DELETE генерирует максимум информации redo, а поскольку единственным заносимым в журнал элементом операции DML на временной таблице являются данные undo, то мы на самом деле наблюдаем, что DELETE генерирует максимум информации undo. Оператор INSERT генерирует очень мало информации undo, предназначенной для записи в журнал. Оператор UPDATE генерирует информацию undo в объеме, равном размеру образа данных до изменения, а DELETE генерирует полный набор данных, записанных в сегмент undo.

Как упоминалось ранее, вы должны также принимать во внимание работу, выполняемую с индексом. Вы обнаружите, что обновление неиндексированного столбца не только выполняется намного быстрее, но и обычно генерирует значительно меньше информации undo, чем обновление индексированного столбца. Для примера мы создадим таблицу с двумя столбцами, содержащими идентичную информацию, а также индексом на одном из них:

```
EODA@ORA12CR1> create table t
  2  as
  3  select object_name unindexed,
  4         object_name indexed
  5  from all_objects
  6  /
Table created.
Таблица создана.

EODA@ORA12CR1> create index t_idx on t(indexed);
Index created.
Индекс создан.

EODA@ORA12CR1> exec dbms_stats.gather_table_stats(user, 'T');
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.
```

Теперь мы обновим эту таблицу, сначала изменив неиндексированный столбец, а затем индексированный. Нам понадобится новый запрос к представлению V\$ для измерения объема информации undo, сгенерированной в каждом случае. Этот запрос показан ниже. Он получает идентификатор текущего сеанса (SID) из V\$MYSTAT, применяет его для поиска нужной записи в представлении V\$SESSION и извлекает адрес транзакции (TADDR). Он использует TADDR для извлечения записи V\$TRANSACTION (если она есть) и выбирает столбец USED_UBLK — количество задействованных блоков undo. Так как мы в текущий момент не находимся внутри транзакции, от запроса можно ожидать возвращения 0 строк:

```
EODA@ORA12CR1> select used_ublk
  2  from v$transaction
  3  where addr = (select taddr
```

```

4          from v$session
5          where sid = (select sid
6                        from v$mystat
7                        where rownum = 1
8                        )
9          )
10 /
no rows selected
нет выбранных строк

```

Но после того как оператор UPDATE начнет транзакцию, запрос возвратит строку:

```

EODA@ORA12CR1> update t set unindexed = lower(unindexed);
72077 rows updated.
72077 строк обновлено.

```

```

EODA@ORA12CR1> select used_ublk
2   from v$transaction
3   where addr = (select taddr
4                  from v$session
5                  where sid = (select sid
6                                from v$mystat
7                                where rownum = 1
8                                )
9                  )
10 /
USED_UBLK
-----
      151

```

```

EODA@ORA12CR1> commit;
Commit complete.
Фиксация завершена.

```

Этот оператор UPDATE использовал 151 блок для сохранения своей информации undo. Фиксация должна освободить их, так что если мы снова запустим запрос к V\$TRANSACTION, он опять сообщит no rows selected (нет выбранных строк). Когда мы обновим те же самые данные, но на этот раз индексированные, то будем наблюдать следующий вывод:

```

EODA@ORA12CR1> update t set indexed = lower(indexed);
72077 rows updated.
72077 строк обновлено.

EODA@ORA12CR1> select used_ublk
2   from v$transaction
3   where addr = (select taddr
4                  from v$session
5                  where sid = (select sid
6                                from v$mystat
7                                where rownum = 1
8                                )
9                  )
10 /
USED_UBLK
-----
      854

```

Как видите, обновление индексированного столбца в этом примере генерирует в несколько раз больше информации undo. Это связано со сложностью самой структуры индексов и тем фактом, что мы обновляем каждую строку в таблице, перемещая значение каждого отдельного ключа индекса в этой структуре.

Ошибка ORA-01555: snapshot too old

В предыдущей главе мы провели краткие исследования ошибки ORA-01555 и взглянули на одну из причин ее появления: слишком частая фиксация. Здесь мы более подробно рассмотрим причины возникновения этой ошибки и способы ее устранения. ORA-01555 — одна из тех ошибок, которые ставят людей в тупик. Она является основой многих мифов, неточностей и предположений.

На заметку! Ошибка ORA-01555 не связана с повреждением данных либо их утерей. С такой точки зрения это “безопасная” ошибка; единственное последствие заключается в том, что получивший ее запрос не может продолжать обработку.

Ошибка ORA-01555 достаточно прямолинейна и возникает только по двум реальным причинам, но поскольку особый случай одной из них встречается настолько часто, можно считать, что всего причин три:

- сегменты отмены слишком малы для работы, которая выполняется в системе;
- ваши программы выдают пересекающиеся операторы COMMIT (на самом деле это вариация предыдущей ситуации), что подробно было раскрыто в главе 8;
- очистка блоков.

Первые две причины напрямую связаны с моделью согласованного чтения Oracle. Как вы помните из главы 7, результаты запроса *предопределены* — в том смысле, что они определены даже до того, как Oracle извлечет первую строку. СУБД Oracle предоставляет этот согласованный на момент времени “снимок” базы данных за счет применения сегментов undo для отката блоков, которые были изменены после начала вашего запроса. Каждый запускаемый оператор вроде следующего:

```
update t set x = 5 where x = 2;
insert into t select * from t where x = 2;
delete from t where x = 2;
select * from t where x = 2;
```

будет видеть согласованное по чтению представление таблицы T и набор строк, где X=2, независимо от любых других параллельно выполняемых действий в базе данных.

На заметку! Представленные здесь четыре оператора — это просто примеры типов операторов, которые будут видеть согласованное по чтению представление таблицы T. Они не предназначены для запуска в виде одиночной транзакции базы данных, т.к. первое обновление приведет к тому, что следующие за ним три оператора не увидят никаких записей. Операторы приведены просто в целях иллюстрации.

Все операторы, которые “читают” таблицу, используют эту согласованность чтения в своих интересах. В только что показанном примере оператор UPDATE читает таблицу, чтобы найти строки, в которых X=2 (и обновить их). Оператор INSERT чи-

тает таблицу в поиске строк с $X=2$, после чего выполняет их вставку, и т.д. Именно двойное предназначение сегментов undo — откат неудавшихся транзакций и обеспечение согласованности чтения — приводит в результате к ошибке ORA-01555.

Третий пункт в приведенном выше списке описывает более коварную причину возникновения ошибки ORA-01555, т.к. она может произойти в базе данных с единственным сеансом, и этот сеанс не модифицирует таблицы, к которым производились запросы, когда была сгенерирована ошибка ORA-01555! Это кажется невозможным — зачем могут понадобиться данные отмены для таблицы, которая гарантированно не модифицировалась? Вскоре вы узнаете ответ на данный вопрос.

Прежде чем взглянуть на все три случая с сопутствующими иллюстрациями, я хотел бы поделиться решениями проблемы с ошибкой ORA-01555 в общем.

- Правильно установите значение параметра UNDO_RETENTION (оно должно превышать продолжительность самой длительно выполняющейся транзакции). Для определения времени выполнения самых длительных запросов можно применять представление V\$UNDOSTAT. Кроме того, удостоверьтесь в наличии на диске достаточного объема свободного пространства, чтобы сегменты undo могли расти до необходимого размера, исходя из запрошенного значения UNDO_RETENTION.
- Увеличьте размер или добавьте больше сегментов undo, используя ручное управление пространством отмены. Это снизит вероятность того, что данные отмены будут перезаписаны во время длительно выполняющегося запроса. Такой метод направлен на решение всех трех указанных ранее проблем. Обратите внимание, что это определенно не лучший метод; настоятельно рекомендуется применять автоматическое управление пространством отмены.
- Сократите время выполнения запроса (настройте его). По возможности это стоит делать всегда, так что имеет смысл заняться настройкой запроса в первую очередь. В итоге сократится потребность в больших сегментах undo. Данный метод также направлен на решение всех трех описанных выше проблем.
- Соберите статистику об участвующих объектах. Это поможет избежать третьей ситуации из перечисленных ранее. Поскольку очистка блоков является результатом очень крупных массовых обновлений или вставок, сбор статистики необходимо делать в любом случае после массового обновления или объемной загрузки данных.

Мы еще вернемся к этим решениям, поскольку их важно знать. Представляется целесообразным даже выделить их особо, прежде чем мы начнем.

Сегменты отмены фактически очень малы

Сценарий выглядит следующим образом: у вас есть система с небольшими транзакциями. Как результат, вам нужно выделять очень мелкие сегменты отмены. Скажем, пусть справедливы следующие утверждения.

- Каждая транзакция в среднем генерирует 8 Кбайт информации undo.
- В среднем происходит 5 таких транзакций в секунду (40 Кбайт информации undo в секунду, 2400 Кбайт в минуту).

- У вас есть транзакция, генерирующая 1 Мбайт информации undo, которая происходит в среднем 1 раз в минуту. Всего вы генерируете около 3,5 Мбайт информации undo в минуту.
- В системе сконфигурировано 15 Мбайт пространства под информацию undo.

Таким образом, в этой базе данных пространства undo более чем достаточно для обработки транзакций. Сегменты undo задействуются циклически и повторно используют пространство в среднем примерно каждые 3–4 минуты. Если вы определяете размер сегментов undo на основе своих транзакций, которые делают модификации, то поступаете совершенно правильно.

Тем не менее, в этой же среде есть необходимость в формировании отчетов. Некоторые из этих запросов требуют на свое выполнение действительно длительное время — возможно, 5 минут. И здесь возникает проблема. Если эти запросы требуют 5 минут на выполнение и нуждаются в представлении данных в состоянии, в каком они были на момент начала запроса, то имеется довольно высокая вероятность возникновения ошибки ORA-01555. Поскольку сегменты undo будут циклически переключаться во время выполнения этого запроса, вы знаете, что часть информации undo, сгенерированной при запуске запроса, теряется — она оказывается перезаписанной. Если вы наталкиваетесь на блок, который был модифицирован приблизительно во время начала запроса, то информация undo для этого блока будет отсутствовать, и вы получите ошибку ORA-01555.

Рассмотрим небольшой пример. Предположим, что есть таблица с блоками 1, 2, 3, ... 1 000 000 внутри нее. В табл. 9.2 показана последовательность событий, которые могли бы произойти.

Таблица 9.2. Временная шкала длительно выполняющегося запроса

| Время (минуты: секунды) | Действие |
|-------------------------------|---|
| 0:00 | Наш запрос начинается |
| 0:01 | Другой сеанс обновляет блок 1 000 000. Информация undo записывается в какой-то сегмент undo |
| 0:01 | Сеанс, выполнивший обновление, производит фиксацию. Сгенерированные данные undo по-прежнему на месте, но будут перезаписаны, если потребуется пространство |
| 1:00 | Наш запрос продвигается дальше, добравшись до блока 200 000 |
| 1:01 | Происходит высокая активность. К этому моменту сгенерировано чуть больше 14 Мбайт информации undo |
| 3:00 | Наш запрос все еще работает. Он дошел примерно до блока 600 000 |
| 4:00 | Сегменты undo начинают циклически переключаться и повторно использовать пространство, которое было активно, когда запрос начался в момент 0:00. В частности, только что было повторно задействовано пространство сегмента, которым пользовался оператор UPDATE блока 1 000 000 в момент 0:01 |
| 5:00 | Наш запрос, наконец, добрался до блока 1 000 000 и определил, что с момента начала запроса этот блок был модифицирован. Он обращается к сегменту отмены и пытается найти там информацию undo для этого блока, чтобы выполнить его согласованное чтение. В данной точке он обнаруживает, что нужная ему информация больше не существует. Генерируется ошибка ORA-01555 и запрос завершается неудачей |

Вот и все, что происходит. Если размер сегментов undo таков, что есть высокая вероятность их повторного использования во время выполнения ваших запросов, и запросы обращаются к данным, которые, возможно, будут модифицированы, то вы получаете неплохие шансы регулярно сталкиваться с ошибкой ORA-01555. В таком случае вы должны увеличить значение параметра UNDO_RETENTION и позволить Oracle позаботиться о выяснении того, сколько данных undo предохранять (это рекомендуемый подход; он намного проще, чем пытаться определять оптимальный размер сегментов undo самостоятельно), либо увеличить размер сегментов undo (или предусмотреть большее их количество). Вам понадобится сконфигурировать достаточный объем пространства undo, чтобы его хватило на время обработки длительно выполняющихся запросов. Система была настроена на транзакции, которые модифицируют данные — вы забыли установить размеры, которые бы учитывали другие ее компоненты.

В Oracle9i и последующих версиях доступны два метода управления пространством отмены в системе.

- Автоматическое управление пространством отмены. Здесь с помощью параметра UNDO_RETENTION базе данных Oracle сообщается, насколько долго следует удерживать информацию undo. СУБД Oracle будет определять количество сегментов undo, подлежащих созданию, а также их размеры на основе текущей рабочей нагрузки. База данных могла даже перераспределять экстенды между отдельными сегментами undo во время выполнения, чтобы удовлетворить целевое значение параметра UNDO_RETENTION, установленное администратором базы данных. Это рекомендуемый подход к управлению пространством отмены.
- Ручное управление пространством отмены. Здесь работу выполняет администратор базы данных. Он определяет, сколько нужно создать вручную сегментов отмены, основываясь на оценочной и наблюдаемой рабочей нагрузке. Администратор базы данных определяет размеры сегментов на основе объема транзакций (сколько они генерируют информации undo) и продолжительности длительно выполняющихся запросов.

При ручном управлении пространством отмены, когда администратор базы данных определяет количество и размеры сегментов отмены, часто возникает путаница. Люди говорят: “Хорошо, у нас сконфигурировано X Мбайт для информации undo, но этот показатель может расти. У нас есть параметр MAXEXTENTS, установленный в 500, а каждый экстенд имеет размер 1 Мбайт, так что область undo может стать довольно большой”. Проблема в том, что вручную управляемые сегменты отмены никогда не растут из-за запроса; они растут только в случае выдачи операторов INSERT, UPDATE или DELETE. Факт выполнения длительно работающего запроса не заставляет Oracle увеличивать ручной сегмент отмены, чтобы в нем сохранялись данные, которые могут понадобиться. Это может вызвать только длительно выполняющаяся транзакция с оператором UPDATE. В предыдущем примере, даже если ручные сегменты отмены имеют потенциал для роста, увеличиваться они не будут. В такой системе необходимо иметь ручные сегменты отмены, которые уже обладают большими размерами. Вы должны выделить пространство для сегментов отмены на постоянной основе, не предоставляя им возможность расти по собственному усмотрению.

Единственное решение проблемы заключается в том, чтобы либо обеспечить такой размер сегментов отмены, чтобы циклическое переключение происходило каждые 6–10 минут, либо так проектировать запросы, чтобы они никогда не выполнялись дольше 2–3 минут. Первое указание основано на факте наличия запросов, выполнение которых занимает 5 минут. В этом случае администратор базы данных должен обеспечить размер постоянно выделенного пространства отмены в 2–3 раза больше теоретически вычисленного. Второе указание (совершенно правильное) также уместно. Всякий раз, когда появляется возможность ускорить выполнение запросов, вы должны ею воспользоваться. Если сгенерированная с момента запуска запроса информация undo не будет никогда перезаписываться, вы избежите ошибки ORA-01555.

При автоматическом управлении пространством отмены в отношении ошибки ORA-01555 все гораздо проще. Вместо того чтобы заранее выяснять размер пространства отмены и предварительно выделять его, администратор базы данных сообщает базе данных время наиболее длительно выполняющегося запроса и устанавливает это значение в параметре `UNDO_RETENTION`. СУБД Oracle попытается предохранять информацию undo минимум в течение указанного времени. Если доступное пространство будет израсходовано, Oracle расширит сегмент отмены, а не станет повторно применять его, стараясь удовлетворить параметру `UNDO_RETENTION`. В этом главное отличие от ручного управления пространством отмены, при котором выделенное пространство отмены используется повторно, как только это становится возможным. Именно по причине поддержки параметра `UNDO_RETENTION` я настоятельно рекомендую, где только возможно, применять *автоматическое управление пространством отмены*. Этот единственный параметр значительно снижает вероятность появления ошибки ORA-01555 (когда он установлен корректно).

При использовании ручного управления пространством отмены важно помнить также о том, что вероятность возникновения ошибки ORA-01555 диктуется *наименьшим* из сегментов отмены в системе, а не наибольшим или средним. Добавление одного “крупного” сегмента отмены не решит проблему. Достаточно наименьшему сегменту отмены подвергнуться перезаписи во время обработки запроса, и у этого запроса появляется шанс столкнуться с ошибкой ORA-01555. Вот почему я был сторонником создания сегментов отмены с равными размерами, когда применял унаследованные сегменты отмены. В таком случае любой сегмент будет одновременно и наименьшим, и наибольшим. По этой причине я также избегаю использования сегментов отмены с “оптимально” установленными размерами. Если вы сжимаете сегмент отмены, который должен расти, то избавляетесь от большого объема информации undo, которая может понадобиться сразу после этого. Вы отбрасываете самые старые данные undo, сводя риск к минимуму, но риск все же остается. Я предпочитаю вручную сжимать сегменты undo в периоды минимальной нагрузки на базу, если вообще это стоит делать.

Здесь я собираюсь немного углубиться в роль администратора базы данных, так что мы перейдем к рассмотрению следующего случая. Просто важно, чтобы вы понимали, что ошибка ORA-01555 в этом случае обусловлена некорректной настройкой системы для существующей рабочей нагрузки. Единственное решение — правильно выбрать размеры сегментов. Это не ваша ошибка, но ваша проблема, т.к. вы сталкиваетесь с ней. Происходит то же самое, как и при выходе за пределы емкости временного пространства во время выполнения запроса. Вы либо конфигурируете достаточ-

ный объем временного пространства для системы, либо переписываете запросы так, чтобы они применяли план, который не требует временного пространства.

Чтобы продемонстрировать этот эффект, мы можем подготовить небольшой, несколько искусственный тест. Создадим очень маленькое табличное пространство отмены, и в одном сеансе будем генерировать множество мелких транзакций, фактически добиваясь, чтобы выделенное пространство повторно использовалось много раз — независимо от установки `UNDO_RETENTION`, поскольку рост табличного пространства отмены запрещен. Сеанс, применяющий этот сегмент отмены, будет модифицировать таблицу `T`. Он будет использовать полное сканирование `T` и читать ее от начала до конца. В другом сеансе мы запустим запрос, который будет читать таблицу `T` через индекс. В этом случае он будет читать таблицу некоторым случайным образом: сначала прочитает строку 1, затем — строку 1000, потом — строку 500, за ней — строку 20 001 и т.д. Таким образом, в процессе выполнения запроса мы будем посещать блоки в случайном порядке и, возможно, неоднократно. Вероятность получения ошибки `ORA-01555` в данном случае составляет почти 100%. Итак, в одном сеансе мы запускаем следующие команды:

```
EODA@ORA12CR1> create undo tablespace undo_small
2 datafile '/tmp/undo.dbf' size 2m
3 autoextend off
4 /
```

Tablespace created.

Табличное пространство создано.

```
EODA@ORA12CR1> alter system set undo_tablespace = undo_small;
System altered.
```

Система изменена.

Теперь настроим таблицу `T`, предназначенную для запросов и модификации. Обратите внимание, что данные в этой таблице располагаются неупорядочено. Оператору `CREATE TABLE AS SELECT` свойственно помещать строки в блоки в том порядке, в каком он извлекает их из запроса. Мы просто перемешиваем строки, чтобы они не были искусственно отсортированы ни в каком порядке, рандомизировав их распределение:

```
EODA@ORA12CR1> drop table t purge;
```

Table dropped.

Таблица удалена.

```
EODA@ORA12CR1> create table t
2 as
3 select *
4 from all_objects
5 order by dbms_random.random;
```

Table created.

Таблица создана.

```
EODA@ORA12CR1> alter table t add constraint t_pk primary key(object_id);
Table altered.
```

Таблица изменена.

```
EODA@ORA12CR1> exec dbms_stats.gather_table_stats( user, 'T', cascade=> true );
PL/SQL procedure successfully completed.
```

Процедура PL/SQL успешно завершена.

Теперь мы готовы проводить модификации:

```

EODA@ORA12CR1> begin
  2   for x in ( select rowid rid from t )
  3   loop
  4       update t set object_name = lower(object_name) where rowid = x.rid;
  5       commit;
  6   end loop;
  7 end;
  8 /

```

А сейчас, пока этот блок кода PL/SQL выполняется, мы запустим запрос в другом сеансе. Этот запрос будет читать таблицу T и обрабатывать каждую запись. Каждая строка обрабатывается примерно 1/100 секунды перед извлечением следующей строки (эмулируется с помощью DBMS_LOCK.SLEEP(0.01)). В запросе применяется подсказка FIRST_ROWS, чтобы заставить его использовать созданный индекс для чтения записей из таблицы через индекс, отсортированный по OBJECT_ID. Из-за того, что данные вставлялись в таблицу в случайном порядке, блоки таблицы также запрашиваются достаточно случайно. Приведенный ниже блок проработает всего пару секунд и даст сбой:

```

EODA@ORA12CR1> declare
  2   cursor c is
  3   select /*+ first_rows */ object_name
  4       from t
  5       order by object_id;
  6
  7   l_object_name t.object_name%type;
  8   l_rowcnt      number := 0;
  9 begin
 10   open c;
 11   loop
 12       fetch c into l_object_name;
 13       exit when c%notfound;
 14       dbms_lock.sleep( 0.01 );
 15       l_rowcnt := l_rowcnt+1;
 16   end loop;
 17   close c;
 18 exception
 19   when others then
 20       dbms_output.put_line( 'rows fetched = ' || l_rowcnt );
 21       raise;
 22 end;
 23 /

```

rows fetched = 159

declare

*

ERROR at line 1:

ORA-01555: snapshot too old: rollback segment number 16 with name
 "_SYSSMU16_587457654\$" too small

ORA-06512: at line 21

ОШИБКА в строке 1:

*ORA-01555: устаревший снимок: сегмент отката с номером 16 по имени
 _SYSSMU16_587457654\$ очень мал*

ORA-06512: в строке 21

Как видите, до возникновения ошибки ORA-01555: snapshot too old блок успел обработать только 159 строк. Чтобы исправить ситуацию, необходимо удостовериться в том, что сделаны следующие действия.

- Параметр UNDO_RETENTION в базе данных установлен так, что имеет достаточно большое значение, чтобы данный процесс чтения мог быть завершен. Это позволит базе данных увеличить размер табличного пространства отмены настолько, чтобы вместить всю необходимую информацию undo.
- Табличному пространству отмены разрешено расти или же вы вручную выделяете ему больше дискового пространства.

Для этого примера я определил, что длительно работающему процессу потребуется для завершения около 720 секунд (в таблице имеется около 72 000 записей, поэтому при затратах 0,01 секунды на запись получаем 720 секунд). Значение UNDO_RETENTION установлено в 900 (оно выражено в секундах, т.е. информации undo будет удерживаться около 15 минут). Я изменил файл данных табличного пространства отмены так, чтобы разрешить ему расти на 1 Мбайт за раз, вплоть до предельного размера в 2 Гбайт:

```
EODA@ORA12CR1> alter database
  2 datafile '/tmp/undo.dbf'
  3 autoextend on
  4 next 1m
  5 maxsize 2048m;
Database altered.
База данных изменена.
```

Когда я снова запущу эти процессы параллельно, оба они выполнятся до конца. Файл данных табличного пространства отмены на этот раз увеличивается в размере, т.к. ему это позволено, и время удержания информации undo, установленное параметром, соблюдается.

```
EODA@ORA12CR1> select bytes/1024/1024
  2   from dba_data_files
  3  where tablespace_name = 'UNDO_SMALL';
BYTES/1024/1024
-----
                21
```

Итак, вместо получения ошибки процессы завершились успешно, а пространство отмены выросло достаточно для того, чтобы удовлетворить наши потребности. Правда, в этом примере появление ошибки объясняется исключительно тем фактом, что мы читали таблицу T через индекс и производили случайные чтения по всей таблице. Если взамен быстро выполнить полное сканирование, то появится реальный шанс избежать ошибки ORA-01555 *в этом конкретном случае*. Причина в том, что операторам SELECT и UPDATE приходится осуществлять полное сканирование T и, скорее всего, SELECT опередит UPDATE в процессе сканирования (оператор SELECT должен только читать, а UPDATE — как читать, так и обновлять данные, а потому он работает медленнее). Производя случайное чтение, мы увеличиваем вероятность того, что SELECT придется прочитать блок, который UPDATE модифицировал и зафиксировал много строк тому назад. Это как раз и демонстрирует коварную природу

ду ошибки ORA-01555. Ее появление зависит от того, каким образом параллельные сеансы обращаются и манипулируют таблицами базы данных.

Отложенная очистка блоков

Причину возникновения ошибки ORA-01555 трудно устранить полностью, но в любом случае она происходит редко, т.к. подходящие для нее обстоятельства складываются нечасто (по крайней мере, в Oracle8i и последующих версиях). Мы уже обсуждали механизм очистки блоков, но чтобы подытожить, напомним, что это процесс, при котором следующему сеансу для доступа к блоку после его модификации может потребоваться проверить, активна ли до сих пор транзакция, модифицировавшая этот блок в последний раз. Как только процесс определяет, что транзакция больше не активна, он очищает блок, так что следующему сеансу для доступа к нему не придется проходить через этот процесс заново. Чтобы очистить блок, СУБД Oracle определяет сегмент отмены, использованный для предыдущей транзакции (по заголовку блока), и затем выясняет, указывает ли заголовок сегмента отмены на то, что транзакция была зафиксирована, и если это так, то когда она была зафиксирована. Такое подтверждение формируется одним из двух методов. Первый метод состоит в том, что Oracle может определить, что транзакция была зафиксирована очень давно, хотя слот этой транзакции был перезаписан в таблице транзакций сегмента отмены. Другой метод связан с тем, что SCN (системный номер изменения) операции COMMIT по-прежнему находится в таблице транзакций сегмента отмены, а это значит, что транзакция фиксировалась недавно, и ее слот не был перезаписан.

Чтобы получить ошибку ORA-01555 при отложенной очистке блока, должны быть удовлетворены все перечисленные ниже условия.

- Изменение произведено и зафиксировано, но блоки не были очищены автоматически (например, транзакция модифицировала больше блоков, чем умещается в 10% кеша буферов блоков SGA).
- Эти блоки не затрагиваются другим сеансом и не будут затронуты до тех пор, пока их не посетит наш неудачливый запрос (который мы вскоре покажем).
- Начинается длительно выполняющийся запрос. В конечном итоге он прочитает некоторые из блоков, упомянутых в предыдущем условии. Этот запрос запускается с SCN-номером $t1$ — согласованным по чтению SCN, к которому должен быть произведен откат данных, чтобы обеспечить согласованность чтения. На момент начала запроса запись для модифицирующей транзакции все еще находится в таблице транзакций сегмента отмены.
- На протяжении выполнения запроса в системе производится фиксация многих транзакций. Эти транзакции не затрагивают обсуждаемые блоки (если бы они затронули, то мы не имели бы дело с грядущей проблемой, т.к. блоки были бы очищены старой транзакцией, что решило бы проблему очистки).
- Таблицы транзакций в сегментах отмены организованы циклически и повторно используют слоты из-за высокого уровня фиксаций. Наиболее важно то, что запись для исходной модифицирующей транзакции циклически прокручивается и ее место применяется повторно. Кроме того, система повторно использует экстенды сегментов отмены, препятствуя согласованному чтению самого заголовка блока сегмента отмены.

- Вдобавок минимальный SCN-номер, записанный в сегменте отмены, теперь превышает $t1$ (он больше, чем согласованный по чтению SCN-номер запроса) из-за большого количества фиксаций.

Когда наш запрос получает блок, который был модифицирован и зафиксирован до его начала, возникает проблема. Обычно запрос должен перейти к сегменту отмены, указанному блоком, и найти состояние транзакции, которая модифицировала блок (другими словами, найти SCN-номер операции COMMIT этой транзакции). Если обнаруженный номер SCN меньше $t1$, то наш запрос может использовать данный блок. Если же SCN больше $t1$, то запрос должен выполнить откат этого блока. Однако проблема в том, что наш запрос в этом конкретном случае не может определить, больше SCN-номер операции COMMIT этого блока, чем $t1$, или же меньше. Нет уверенности в том, можно ли применять образ этого блока или нет. В результате возникает ошибка ORA-01555.

Чтобы увидеть это, мы создадим много блоков в таблице, которые нужно будет очистить. Затем мы откроем курсор на этой таблице и позволим множеству мелких транзакций выполняться в какой-то другой таблице — не той, которая была только что обновлена, и в которой был открыт курсор. Теперь мы *знаем*, что с данными, нужными курсору, все будет в порядке — мы должны иметь возможность видеть их все, поскольку модификации в таблице были произведены и зафиксированы *до* открытия курсора. Получение ошибки ORA-01555 на этот раз вызвано описанной выше проблемой с отложенной очисткой блоков. Для этого примера используются следующие установки.

- Табличное пространство отмены UNDO_SMALL размером 4 Мбайт.
- Буферный кеш размером 16 Мбайт, которого достаточно для размещения около 2000 блоков. Это позволит получить ряд грязных блоков, сброшенных на диск, чтобы понаблюдать за данным феноменом.

Прежде всего, мы создадим табличное пространство отмены и “большую” таблицу для последующих запросов:

```
EODA@ORA12CR1> create undo tablespace undo_small
  2 datafile '/tmp/undo.dbf' size 4m
  3 autoextend off
  4 /
```

Tablespace created.

Табличное пространство создано.

```
EODA@ORA12CR1> create table big
  2 as
  3 select a.*, rpad(' ',1000,' ') data
  4 from all_objects a;
```

Table created.

Таблица создана.

```
EODA@ORA12CR1> alter table big add constraint big_pk
  2 primary key(object_id);
```

Table altered.

Таблица изменена.

```
EODA@ORA12CR1> exec dbms_stats.gather_table_stats( user, 'BIG' );
PL/SQL procedure successfully completed.
```

Процедура PL/SQL успешно завершена.

На заметку! Вас может удивить, почему я не применил `CASCADE=>TRUE` в вызове сбора статистики по индексу, созданному по умолчанию с ограничением первичного ключа. Дело в том, что, начиная с версии Oracle 10g, оператор `CREATE INDEX` или `ALTER INDEX REBUILD` уже имеет неявно добавленный сбор статистики в случае, когда индексируемая таблица не является пустой. Таким образом, само действие по созданию индекса обладает побочным эффектом в виде сбора статистики. Нет никакой необходимости повторно собирать статистику, которой мы уже располагаем.

Предыдущая таблица будет иметь много блоков, поскольку на блок будет приходиться 6 или 7 строк, учитывая размер поля `data`, а таблица `ALL_OBJECTS` содержит свыше 70 000 строк. Затем мы создадим небольшую таблицу, которую будут модифицировать множество мелких транзакций:

```
EODA@ORA12CR1> create table small ( x int, y char(500) );
Table created.
Таблица создана.

EODA@ORA12CR1> insert into small select rownum, 'x' from all_users;
25 rows created.
25 строк создано.

EODA@ORA12CR1> commit;
Commit complete.
Фиксация завершена.

EODA@ORA12CR1> exec dbms_stats.gather_table_stats( user, 'SMALL' );
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.
```

Теперь внесем в большую таблицу грязные данные. Мы располагаем очень маленьким табличным пространством отмены, поэтому хотим обновить столько блоков большой таблицы, сколько возможно, генерируя при этом минимальный объем информации `undo`. Для этого мы будем использовать необычный оператор `UPDATE`. По существу показанный ниже подзапрос ищет “первый” идентификатор строки (`rowid`) в каждом блоке. Подзапрос возвратит `rowid` каждого блока базы данных, идентифицируя единственную строку в нем. Мы обновим эту строку, установив поле `VARCHAR2(1)`. Это позволит обновить все блоки в таблице (свыше 8000 в этом примере), заполнив буферный кеш грязными блоками, которые нужно будет впоследствии записать (прямо сейчас есть место только для 500 блоков). Мы также обеспечим применение небольшого табличного пространства отмены. Чтобы достичь этого и не превысить емкость табличного пространства отмены, мы задействуем оператор `UPDATE`, который обновит только “первую строку” в каждом блоке. Необходимым средством в этой операции послужит встроенная аналитическая функция `ROW_NUMBER()`; она назначает номер 1 “первой строке” блока базы данных в таблице, которая и будет единственной строкой в блоке, подлежащем обновлению:

```
EODA@ORA12CR1> alter system set undo_tablespace = undo_small;
System altered.
Система изменена.
```



```

EODA@ORA12CR1> update big
2     set temporary = temporary
3   where rowid in
4   (
5     select r
6     from (
7       select rowid r, row_number() over
8         (partition by dbms_rowid.rowid_block_number(rowid) order by rowid) rn
9       from big
10      )
11   where rn = 1
12 )
13 /
3064 rows updated.
3064 строк обновлено.

EODA@ORA12CR1> commit;
Commit complete.
Фиксация завершена.

```

Таким образом, нам известно, что есть много грязных блоков на диске. Мы определенно запишем некоторые из них, т.к. для удержания их всех просто не хватит места. Затем мы откроем курсор, но пока не будем извлекать строку. Вспомните, что при открытии курсора результирующий набор является предопределенным, поэтому даже если СУБД Oracle не обработала строку данных, действие по открытию курсора фиксирует момент времени, на который должны быть представлены результаты. Теперь, поскольку мы будем извлекать данные, которые были только что обновлены и зафиксированы, и знаем, что никто другой не модифицировал эти данные, то должны иметь возможность извлекать строки вообще без необходимости в информации undo. Но здесь дает о себе знать отложенная очистка блоков. Транзакция, которая модифицировала эти блоки, настолько новая, что СУБД Oracle будет обязана проверить, зафиксирована ли она, прежде чем мы начнем, и если мы перезапишем эту информацию (также хранящуюся в табличном пространстве отмены), то запрос завершится неудачей. Итак, вот открытие курсора:

```

EODA@ORA12CR1> variable x refcursor
EODA@ORA12CR1> exec open :x for select * from big where object_id < 100;
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.

EODA@ORA12CR1>
EODA@ORA12CR1> !./run.sh

```

run.sh — это сценарий оболочки; он просто запускает девять сеансов SQL*Plus, используя следующую команду:

```

$ORACLE_HOME/bin/sqlplus eoda/foo @test2 1 &
$ORACLE_HOME/bin/sqlplus eoda/foo @test2 2 &
... (запуск сеансов 3-8) ...
$ORACLE_HOME/bin/sqlplus eoda/foo @test2 9 &

```

Здесь каждому сеансу SQL*Plus передается отличающийся номер (1, 2, 3 и т.д.). В предыдущем сценарии не забудьте заменить eoda/foo корректным именем пользователя и паролем для своей среды.

Сценарий `test2.sql`, который запускает каждая команда, выглядит так:

```
begin
  for i in 1 .. 5000
  loop
    update small set y = i where x= &1;
    commit;
  end loop;
end;
/
exit
```

На заметку! Все сценарии, применяемые в этом примере, доступны для загрузки на веб-сайте издательства. Сценарий `demo11.sql`, находящийся в каталоге `ch09`, автоматизирует данный пример.

В результате мы имеем девять сеансов, которые внутри небольшого цикла иницируют множество транзакций. Сценарий `run.sh` ожидает завершения работы девяти сеансов SQL*Plus, и затем мы возвращаемся в свой сеанс, где открыт курсор. При попытке отобразить его мы получим ошибку `ORA-01555`:

```
EODA@ORA12CR1> print x
ERROR:
ORA-01555: snapshot too old: rollback segment number 17 with name
"_SYSSMU17_452567810$" too small
no rows selected
ОШИБКА:
ORA-01555: устаревший снимок: сегмент отката с номером 17 по имени
"_SYSSMU17_452567810$" слишком мал
нет выбранных строк
```

Как уже упоминалось, приведенный пример представляет редкий случай. Он требует совпадения многочисленных обстоятельств. Необходимы блоки, нуждающиеся в очистке, а такие блоки — нечастое явление в Oracle8i и последующих версиях. Вызов `DBMS_STATS` для сбора статистики избавляется от них, так что большинство распространенных случаев — крупные массовые обновления и пакетные загрузки — можно не учитывать, потому что после таких операций таблицы все равно должны быть проанализированы. Большинство транзакций затрагивают менее 10% блоков в буферном кеше; следовательно, они не генерируют блоки, которые необходимо очистить. Если вы полагаете, что столкнулись с этой проблемой, когда выполнение оператора `SELECT` в отношении таблицы, к которой другие операции DML не применялись, вызывает ошибку `ORA-01555`, попробуйте воспользоваться перечисленными далее рекомендациями.

- Прежде всего, удостоверьтесь, что вы применяете транзакции с “корректными размерами”. Убедитесь, что не производите фиксацию чаще, чем должны.
- Используйте пакет `DBMS_STATS` для сканирования связанных объектов, очищая их после загрузки. Поскольку очистка блоков является результатом выполнения очень больших массовых операторов `UPDATE` или `INSERT`, это нужно делать в любом случае.

- Позвольте табличному пространству отмены возрасть, предоставив ему место для расширения и увеличив время удержания информации undo. В результате снизится вероятность перезаписи слота таблицы транзакций в сегменте отмены на протяжении длительно выполняющегося запроса. Это такое же решение, как и в случае, когда ошибка ORA-01555 вызвана другой причиной (эти две причины очень тесно связаны; вы уже сталкивались с повторным использованием сегмента отмены во время обработки запроса). Фактически я перезапустил предыдущий пример с табличным пространством отмены, настроенным на автоматическое расширение с шагом в 1 Мбайт, с удержанием информации undo в течение 900 секунд. Запрос к таблице BIG завершился успешно.
- Сократите время выполнения запроса — настройте его. Поступать так всегда полезно, когда есть возможность, так что это может быть первым действием, которое следует предпринять.

Резюме

В этой главе рассматривались механизмы повтора и отмены, а также их значение для разработчиков. Я представил здесь ситуации и условия главным образом в целях ознакомления, поскольку эти проблемы должны решаться администраторами баз данных или системными администраторами. Ключевой момент, который вы должны вынести из настоящей главы — это важность повтора и отмены, а также тот факт, что они не должны трактоваться как накладные расходы. Они являются неотъемлемыми компонентами базы данных, необходимыми и обязательными. Как только вы хорошо поймете, как они работают и что делают, то сможете более эффективно их применять. Пожалуй, наиболее важным аспектом можно считать понимание того, что вы ничего не “экономите”, выполняя фиксацию чаще, чем должны (на самом деле вы впустую расходуете ресурсы, т.к. частые фиксации требуют больше процессорного времени, больше дискового пространства и написания большего объема кода). Будьте в курсе того, что база данных должна делать, а затем позвольте ей это делать.

ГЛАВА 10

Таблицы базы данных

В этой главе мы рассмотрим различные типы таблиц базы данных и сценарии их использования (т.е. когда один тип таблиц подходит больше других). Мы сосредоточим внимание на характеристиках физического хранения таблиц: как организованы данные, и каким образом они хранятся.

Когда-то существовал только один тип таблиц — *обычные* таблицы. Они управлялись таким же способом, как и *куча из сущностей*. Со временем в Oracle появились более сложные типы таблиц. Теперь в дополнение к традиционным таблицам имеются кластеризованные таблицы (трех разновидностей), индекс-таблицы, вложенные таблицы, временные таблицы, внешние таблицы и объектные таблицы. Каждый тип таблиц обладает своими характеристиками, которые делают его пригодным для применения в разных прикладных областях.

Типы таблиц

Прежде чем углубляться в детали, мы определим каждый тип таблиц. Всего в Oracle имеется девять основных типов таблиц, которые описаны ниже.

- **Традиционные таблицы (heap organized table).** Это обычные, стандартные таблицы базы данных. Данные управляются в манере, подобной тому, как управление реализовано в куче. При добавлении данных используется первое найденное в сегменте свободное пространство, которое может уместить эти данные. При удалении данных из такой таблицы пространство, которое они занимали, становится доступным для повторного применения последующими операторами INSERT и UPDATE. Это как раз и объясняет наличие слова “heap” (“куча”) в английском названии такого типа таблиц. *Куча* — это область пространства, которое используется отчасти произвольно.
- **Индекс-таблицы (index organized table — IOT).** Эти таблицы хранятся в индексной структуре, что накладывает определенный физический порядок на сами строки. Если в традиционных таблицах данные размещаются там, где они могут уместиться, то в индекс-таблицах данные сохраняются в отсортированном порядке согласно первичному ключу.
- **Кластеризованные индекс-таблицы (index clustered table).** *Кластеры* — это группы, состоящие из одной или более таблиц, которые физически хранятся в тех же самых блоках базы данных, при этом все строки в них разделяют общее значение кластерного ключа и физически находятся близко друг к другу. В этой структуре достигаются две цели. Во-первых, многие таблицы мо-

гут храниться физически вместе. Обычно в блоке базы данных можно ожидать нахождения данных только из одной таблицы, но в случае кластеризованных таблиц в одном блоке могут храниться данные из нескольких таблиц. Во-вторых, все данные, содержащие одно и то же значение кластерного ключа, такое как `DEPTNO=10`, будут храниться физически вместе. Данные кластеризуются вокруг этого значения кластерного ключа. Кластерный ключ строится с применением индекса со структурой В-дерева. Преимущество кластеризованных индекс-таблиц связано с сокращением объема дискового ввода-вывода и улучшением производительности запросов при доступе в таблицы, которые часто соединяются по кластерному ключу.

- **Кластеризованные хеш-таблицы (hash clustered table).** Такие таблицы похожи на кластеризованные индекс-таблицы, но вместо использования индекса со структурой В-дерева для определения местонахождения данных по кластерному ключу хеш-кластер хеширует кластерный ключ, чтобы попасть в блок базы данных, где должны располагаться данные. В хеш-кластере данные являются индексом (фигурально выражаясь). Эти таблицы подходят для хранения данных, которые читаются часто через операцию равенства на ключе.
- **Отсортированные кластеризованные хеш-таблицы (sorted hash clustered table).** Этот тип таблиц появился в версии Oracle 10g и сочетает в себе ряд аспектов кластеризованных хеш-таблиц и индекс-таблиц. Концепция состоит в следующем: вы имеете определенное значение ключа, по которому будут хешироваться строки (например, `CUSTOMER_ID`), и последовательность относящихся к этому ключу записей, которые поступают в отсортированном порядке (записи на основе временных меток) и обрабатываются в таком отсортированном порядке. Например, заказчик размещает свои заказы в вашей системе ввода заказов, и эти заказы извлекаются и обрабатываются методом FIFO (First In, First Out — первым пришел, первым обслужен). В такой системе подходящей структурой данных может оказаться отсортированный хеш-кластер.
- **Вложенные таблицы (nested table).** Эти таблицы являются частью объектно-реляционных расширений Oracle. Они представляют собой просто генерируемые и обслуживаемые системой дочерние таблицы в рамках отношения “родительская—дочерняя”. Они работают во многом так же, как таблицы EMP и DEPT в схеме SCOTT, где EMP является вложенной таблицей. EMP считается дочерней таблицей для DEPT, поскольку она содержит внешний ключ DEPTNO, который указывает на DEPT. Главное отличие в том, что это не автономные традиционные таблицы.
- **Временные таблицы (temporary table).** Эти таблицы хранят промежуточные данные на протяжении транзакции или сеанса. При необходимости они выделяют временные экстенды из временного табличного пространства текущего пользователя. Каждый сеанс будет видеть только свои экстенды; он никогда не будет видеть данные, созданные любым другим сеансом. Временные таблицы позволяют краткосрочно сохранять данные с тем преимуществом, что при этом генерируется намного меньший объем информации redo (и меньше undo в версии Oracle 12c), чем в случае традиционной таблицы (детальное обсуждение поведения временных таблиц в отношении redo и undo приведено в главе 9).

- **Объектные таблицы (object table).** Эти таблицы создаются на основе объектного типа. Они имеют специальные атрибуты, которые не ассоциируются с не-объектными таблицами, такие как столбец `OID` (объектный идентификатор), генерируемый системой для каждой строки. Объектные таблицы на самом деле являются частными случаями традиционных, индекс- и временных таблиц, и они также могут включать вложенные таблицы как часть своей структуры.
- **Внешние таблицы (external table).** Данные этих таблиц не хранятся в самой базе данных; взамен они располагаются за пределами базы данных в обычных файлах операционной системы (со столбцами данных в файле, обычно идентифицируемых разделителями или позициями). В Oracle9i и последующих версиях внешние таблицы предоставляют возможность выдавать запросы к файлу, находящемуся вне базы данных, как если бы он был традиционной таблицей внутри базы данных. Они наиболее полезны в качестве средства помещения информации в базу данных (являются очень мощным инструментом загрузки данных). Более того, в версии Oracle 10g, в которой появилась возможность выгрузки данных во внешние таблицы, они позволяют легко перемещать данные из одной базы данных Oracle в другую, не используя связи баз данных. Внешние таблицы довольно подробно рассматриваются в главе 15.

Представленная ниже информация касается таблиц всех типов.

- Таблица может иметь вплоть до 1000 столбцов, хотя я не рекомендую применять максимальное количество столбцов, если только в этом нет крайней необходимости. Таблицы наиболее эффективно работают тогда, когда в них содержится гораздо меньше 1000 столбцов. Строку с более чем 254 столбцами Oracle будет внутренне хранить в виде отдельных фрагментов, которые указывают друг на друга и должны собираться вместе для воссоздания полного образа строки.
- Таблица может иметь практически неограниченное количество строк, однако это не будут допускать другие ограничения. Например, обычно табличное пространство может содержать максимум 1022 файла (хотя в Oracle 10g и последующих версиях есть табличные пространства `BIGFILE`, которые позволяют выходить за рамки этих пределов). Скажем, вы имеете типичное табличное пространство и используете файлы размером 32 Гбайт, т.е. общий размер составит 32 704 Гбайт (1022 файла по 32 Гбайт). Это соответствует 2 143 289 344 блокам размером по 16 Кбайт. В каждом таком блоке могут уместиться 160 строк размером от 80 до 100 байтов на блок. В итоге получается 342 926 295 040 строк. Однако если секционировать таблицу, то это число можно легко увеличить в несколько раз. Например, в таблице с 1024 хеш-разделами это будет 1024 × 342 926 295 040 строк. Пределы, конечно, существуют, но на практике вы будете сталкиваться с другими ограничениями гораздо раньше, не успевая подойти даже близко к 342 926 295 040 строк на таблицу.
- Таблица может иметь столько индексов, сколько существует перестановок столбцов в ней (а также перестановок функций на этих столбцах и перестановок любого уникального выражения, какое только можно вообразить). С появлением индексов на основе функций реальное количество индексов, которые можно создать, теоретически бесконечно! Однако, опять-таки, практические ограни-

чения, такие как общая производительность (каждый новый индекс добавляет накладные расходы к оператору INSERT для этой таблицы), будут уменьшать действительное количество создаваемых и сопровождаемых индексов.

- Нет никаких ограничений относительно максимального количества таблиц, которое может существовать даже внутри одной базы данных. Но снова практические ограничения будут удерживать это количество в разумных пределах. Миллионов таблиц у вас не будет (такое количество таблиц не удобно ни создавать, ни поддерживать), но тысячи таблиц вполне могут существовать.

В следующем разделе мы рассмотрим некоторые параметры и термины, касающиеся таблиц. После этого мы перейдем к обсуждению традиционных таблиц и затем к исследованию остальных типов таблиц.

Терминология

В этом разделе рассматриваются разнообразные параметры хранения и термины, связанные с таблицами. Не все параметры применяются с каждым типом таблиц. Например, параметр PCTUSED в контексте индекс-таблицы не имеет смысла (причина станет очевидной в главе 11). Параметры, подходящие для каждого типа таблиц, будут раскрыты при обсуждении каждого отдельного типа. Цель заключается в том, чтобы представить вам термины и дать их определение. Более подробные сведения по использованию параметров приведены в последующих разделах.

Сегмент

В СУБД Oracle *сегмент* (segment) — это объект, который занимает определенное место на диске. Существует много типов сегментов, и наиболее популярные из них перечислены ниже.

- **Кластер (cluster).** Этот тип сегментов способен хранить таблицы. Существуют два вида кластеров: со структурой В-дерева и хеш-кластеры. Кластеры обычно применяются для хранения связанных данных из множества таблиц, предварительно соединенных в одном и том же блоке базы данных, и для хранения вместе связанной информации из одиночной таблицы. Термином “кластер” обозначается способность сегментов такого типа физически группировать связанные между собой данные.
- **Таблица (table).** Сегмент таблицы удерживает данные для таблицы базы данных и вероятно является наиболее распространенным типом сегментов, используемым в сочетании с индексным сегментом.
- **Секция (partition) или подсекция (subpartition) таблицы.** Этот тип сегментов применяется при секционировании и очень похож на табличный сегмент. Сегмент типа секции или подсекции таблицы хранит только срез данных из таблицы. Секционированная таблица образована из одного или более сегментов типа секции таблицы, а таблица с составным секционированием — из одного или более сегментов типа подсекции таблицы.
- **Индекс (index).** В сегменте такого типа хранится структура индекса.

- **Секция индекса (index partition).** Подобно секции таблицы, этот тип сегментов содержит некоторый срез индекса. Секционированный индекс состоит из одного или более сегментов типа секции индекса.
- **LOB-секция, LOB-подсекция, LOB-индекс и LOB-сегмент.** Сегменты типа LOB-индекса и LOB-сегмента хранят структуру *большого объекта* (large object — LOB). Когда таблица, содержащая LOB-объект, секционируется, LOB-сегмент тоже разбивается на секции — для этого используется сегмент типа LOB-секции. Интересно отметить, что сегмента типа секции LOB-индекса не существует: по какой-то причине Oracle помечает секционированный LOB-индекс как секцию индекса (интересно, почему для LOB-индекса было предусмотрено специальное название). Объекты LOB подробно обсуждаются в главе 12.
- **Вложенная таблица.** Это тип сегментов, назначенный для вложенных таблиц — специальной разновидности дочерней таблицы в отношении “главный—подчиненный”, о чем более подробно будет рассказываться далее в этой главе.
- **Сегменты отката (rollback) и сегменты отмены Type2 (Type2 undo).** В таких сегментах хранятся данные отмены. Сегменты отката создаются и управляются вручную администратором базы данных. Сегменты Type2 undo создаются и управляются автоматически Oracle.

Итак, например, таблица *может быть* сегментом. Индекс *может быть* сегментом. Я сделал акцент на формулировке “может быть”, потому что мы можем секционировать индекс на отдельные сегменты. Таким образом, объект индекса сам по себе будет просто определением, а не физическим сегментом — индекс будет образован из множества секций и каждая *секция индекса* будет сегментом. Таблица может быть или не быть сегментом. По той же причине мы можем иметь много табличных сегментов из-за секционирования или создать таблицу в сегменте, который называется кластером. Здесь таблица будет располагаться (возможно, с другими таблицами) в том же самом сегменте кластера.

Однако в самом общем случае таблица и индекс будут сегментами. Пока что проще всего думать о них именно так. Когда создается таблица, это обычно означает создание нового табличного сегмента и, как было показано в главе 3, такой сегмент состоит из экстенгов, а экстенги — из блоков. Это нормальная иерархия хранения. Очень важно отметить, что это отношение “один к одному” получается только в общем случае. Например, взгляните на следующий простой оператор CREATE TABLE:

```
EODA@ORA12CR1> create table t ( x int primary key, y clob, z blob );
```

Этот оператор создает шесть сегментов в Oracle 11g Release 1 и предшествующих версиях; в Oracle 11g Release 2 и последующих версиях по умолчанию создание сегментов откладывается до вставки первой строки (ниже приведен синтаксис для немедленного создания сегментов). Выполнение этого оператора CREATE TABLE в схеме без владельца дает такой результат:

```
EODA@ORA12CR1> select segment_name, segment_type from user_segments;
no rows selected
нет выбранных строк

EODA@ORA12CR1> create table t
2 ( x int primary key,
```



```

3   y clob,
4   z blob )
5  SEGMENT CREATION IMMEDIATE
6  /

```

Table created.

Таблица создана.

```
EODA@ORA12CR1> select segment_name, segment_type from user_segments;
```

| SEGMENT_NAME | SEGMENT_TYPE |
|-----------------------------|--------------|
| T | TABLE |
| SYS_LOB0000021096C00003\$\$ | LOBSEGMENT |
| SYS_LOB0000021096C00002\$\$ | LOBSEGMENT |
| SYS_IL0000021096C00003\$\$ | LOBINDEX |
| SYS_IL0000021096C00002\$\$ | LOBINDEX |
| SYS_C005958 | INDEX |

6 rows selected.

6 строк выбрано.

В этом примере сама таблица создала сегмент, что видно по первой строке вывода. Кроме того, ограничение первичного ключа привело в данном случае к созданию индексного сегмента для принудительного обеспечения уникальности.

На заметку! Ограничения уникальности и первичного ключа могут создавать или не создавать новый индекс. Если на столбцах с ограничениями индекс уже существует, и эти столбцы находятся на ведущих позициях в индексе, то ограничение может и будет использовать их.

Вдобавок каждый LOB-столбец создал два сегмента: один для хранения действительных порций данных, на которые указывает указатель большого символьного объекта (character large object — CLOB) или большого двоичного объекта (binary large object — BLOB), и один для организации этих данных. LOB-столбцы могут вмещать очень крупные фрагменты информации, вплоть до многих гигабайтов. Они хранятся порциями в LOB-сегменте, а LOB-индекс применяется для отслеживания их местоположения и порядка, в котором к ним должен осуществляться доступ.

Обратите внимание, что в строке 5 оператора CREATE TABLE используется синтаксис, специфичный для Oracle 11g Release 2 и последующих версий — конструкция SEGMENT CREATION IMMEDIATE. Если вы попытаетесь применить этот синтаксис в более ранних выпусках, то получите следующее сообщение об ошибке:

```

ops$tkyte@ORA11GR1> Create table t
2  ( x int primary key,
3    y clob,
4    z blob )
5  SEGMENT CREATION IMMEDIATE
6  /

```

SEGMENT CREATION IMMEDIATE

*

ERROR at line 5:

ORA-00922: missing or invalid option

ОШИБКА в строке 5:

ORA-00922: пропущенная или недопустимая опция

На заметку! Средство отложенного создания сегментов доступно только в редакции Enterprise системы Oracle. Если вы работаете в среде со смесью баз данных редакций Enterprise и Standard, будьте аккуратны при экспортировании объектов из базы данных Enterprise в базу данных Standard. Если вы попытаетесь экспортировать объекты, не имеющие созданных сегментов, или импортировать их в базу данных Standard, то можете получить сообщение об ошибке ORA-00439: feature not enabled (ORA-00439: средство не доступно). Один из способов обойти эту проблему предполагает изначальное создание таблиц в базе данных Enterprise с использованием конструкции SEGMENT CREATION IMMEDIATE. За дополнительными деталями обращайтесь к примечанию MOS (My Oracle Support — Моя поддержка Oracle) под номером 1087325.1.

Управление пространством сегментов

Начиная с версии Oracle9i, доступны два метода управления пространством сегментов.

- **Ручное управление пространством сегментов (Manual Segment Space Management — MSSM).** Вы самостоятельно устанавливаете значения таких параметров, как FREELISTS, FREELIST GROUPS, PCTUSED и т.д., чтобы управлять выделением, использованием и повторным использованием пространства в сегменте с течением времени. В настоящей главе я буду ссылаться на этот метод управления пространством как на MSSM, но учтите, что это выбранная мною аббревиатура, которую вы не найдете в документации по Oracle.
- **Автоматическое управление пространством сегментов (Automatic Segment Space Management — ASSM).** Вы самостоятельно устанавливаете значение только одного параметра, который имеет отношение к использованию пространства: PCTFREE. Остальные параметры принимаются при создании сегмента, но игнорируются.

Метод MSSM в Oracle является унаследованной реализацией. Он существует на протяжении многих лет во множестве предшествующих версий. Метод ASSM впервые появился в Oracle9i Release 1 и предназначался для устранения потребности в настройке несметного числа параметров, применяемых при управлении выделением пространства, и обеспечения высокой степени параллелизма. Например, установив параметр FREELISTS в 1 (стандартное значение), вы можете обнаружить, что сегменты, подвергающиеся интенсивным вставкам/обновлениям, начинают страдать от конкуренции за выделение свободного пространства. Когда база данных Oracle собирается вставлять строку в таблицу, обновлять запись ключа индекса либо обновлять строку, приводя к ее перемещению в другое место (вскоре мы обсудим это более подробно), ей может понадобиться получить блок из списка свободных блоков, ассоциированного с сегментом. Если есть только один такой список, то просматривать и модифицировать его может только одна транзакция за раз — транзакциям придется ожидать друг друга. Наличие нескольких списков свободных блоков и групп списков свободных блоков позволяет в этом случае увеличить степень параллелизма, т.к. транзакции могут просматривать разные списки и не соперничать между собой.

Когда я вскоре буду описывать параметры хранения, то обязательно расскажу, какие из них предназначены для ручного, а какие — для автоматического управления пространством сегментов, но уже сейчас можно перечислить параметры, которые применимы только к сегментам ASSM:

- BUFFER_POOL
- PCTFREE
- INITTRANS
- MAXTRANS (только Oracle 9i; в Oracle 10g и последующих версиях он игнорируется для всех сегментов)

Оставшиеся параметры хранения и физических атрибутов к сегментам ASSM не применимы.

Управление пространством сегментов — это атрибут, унаследованный от табличного пространства, в котором содержится сегмент (сегменты никогда не охватывают несколько табличных пространств). Чтобы для сегмента использовался метод ASSM, он должен находиться в табличном пространстве, которое поддерживает такой метод управления пространством.

Маркер максимального уровня заполнения

Этот термин применяется к табличным сегментам, хранящимся в базе данных. Скажем, если представить таблицу в виде плоской структуры или последовательности блоков, уложенных друг за другом слева направо в строке, тогда *маркером максимального уровня заполнения* (High-Water Mark — HWM) будет крайний справа блок, который когда-либо содержал данные (рис. 10.1).

На рис. 10.1 показано, что HWM-маркер в только что созданной таблице начинается с первого блока. По мере помещения в эту таблицу данных количество занятых блоков увеличивается и HWM-маркер сдвигается вправо.

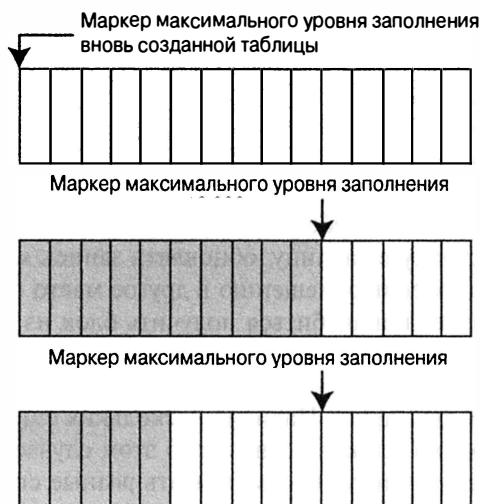


Рис. 10.1. Графическая иллюстрация HWM-маркера

Если мы удалим какие-то (или даже *все*) строки из таблицы, то в итоге может получиться много блоков, которые больше не содержат данных, но они по-прежнему находятся *перед* HWM-маркером и остаются там до тех пор, пока объект не будет перестроен, очищен или сжат. (Сжатие сегментов — это средство, появившееся в версии Oracle 10g, которое поддерживается, только если сегмент находится в табличном пространстве ASSM.)

Маркер HWM важен потому, что во время полного сканирования Oracle будет просматривать все блоки до HWM, даже когда они *не* содержат данных. Это влияет на производительность полного сканирования, особенно если большинство блоков перед HWM-маркером являются пустыми. Чтобы увидеть это, создайте таблицу с 1 000 000 строк (или даже больше) и запустите в отношении нее запрос `SELECT COUNT(*)`. Теперь удалите из таблицы все строки; вы заметите, что запрос `SELECT COUNT(*)` для подсчета 0 строк в таблице выполняется столько же времени, сколько тратилось при подсчете 1 000 000 строк (или дольше, если требуется очистить блоки, как объяснялось в главе 9). Причина в том, что Oracle читает все блоки до HWM-маркера, чтобы выяснить, содержатся ли в них данные. Сравните это с тем, что происходит, когда вместо удаления каждой отдельной строки вы применяете к таблице оператор `TRUNCATE`. Этот оператор сбросит HWM-маркер таблицы в ноль и также очистит все ассоциированные индексы. Когда вы планируете удалять из таблицы все строки, то оператор `TRUNCATE` (если его можно использовать) по этой причине будет наилучшим вариантом для выбора.

Внимание! Имейте в виду, что оператор `TRUNCATE` не поддается откату, равно как и не будет запускать триггеры (если они существуют) на таблице. Вследствие этого перед усечением удостоверьтесь, что точно хотите удалить данные, поскольку отменить такое действие не получится.

В табличном пространстве MSSM сегменты имеют определенный HWM-маркер. Однако в табличном пространстве ASSM имеется просто HWM-маркер и нижний HWM-маркер. В случае MSSM, когда HWM-маркер передвигается дальше (например, по мере вставки новых строк), все блоки форматируются и проверяются, и Oracle может безопасно их читать. С другой стороны, в случае ASSM, когда HWM-маркер передвигается дальше, Oracle не форматирует *все* блоки немедленно: они форматируются и делаются безопасными для чтения при первом их действительном использовании. Это произойдет тогда, когда база данных решит вставить запись в заданный блок. В условиях ASSM данные вставляются в любой блок между нижним HWM-маркером и HWM-маркером, так что многие блоки между этими позициями могут быть не форматированы. Нижний HWM-маркер определен как точка, перед которой все блоки являются форматированными (потому что они содержат данные в текущий момент или содержали их ранее).

Таким образом, при полном сканировании сегмента мы должны знать, что блоки, подлежащие чтению, являются безопасными или неформатированными (т.е. не содержат ничего интересного, поэтому обрабатываться не будут). Чтобы проверке на предмет безопасности/небезопасности подвергся не каждый блок в таблице, Oracle поддерживает два маркера: нижний HWM и HWM. База данных Oracle будет выполнять полное сканирование таблицы вплоть до HWM-маркера, и для всех блоков, находящихся перед нижним HWM-маркером, будет производиться чтение и

обработка. С блоками, расположенными между нижним HWM-маркером и HWM-маркером (рис. 10.2), Oracle придется обращаться более осторожно, и посредством битовой информации ASSM, применяемой для управления этими блоками, выяснять, какие из них должны быть прочитаны, а какие — просто проигнорированы.

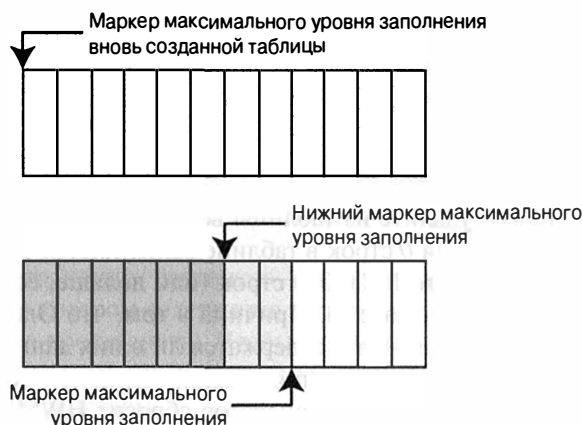


Рис. 10.2. Графическая иллюстрация нижнего HWM-маркера

Списки свободных блоков

Когда используется табличное пространство MSSM, с помощью списка свободных блоков Oracle отслеживает блоки перед HWM-маркером для объектов, которые имеют в них свободное пространство.

На заметку! Списки свободных блоков и группы списков свободных блоков не имеют никакого отношения к табличным пространствам ASSM; этот прием применяется только табличными пространствами MSSM.

С каждым объектом ассоциируется, по меньшей мере, один список свободных блоков и во время использования блоки по мере необходимости добавляются в список свободных блоков либо удаляются из него. Важно отметить, что в списке свободных блоков могут присутствовать только блоки, которые находятся перед HWM-маркером объекта. Блоки, расположенные за HWM-маркером, будут применяться только тогда, когда списки свободных блоков пусты; в этот момент Oracle передвигает HWM-маркер и добавляет эти блоки в список свободных блоков. В подобной манере Oracle откладывает увеличение HWM-маркера для объекта до тех пор, пока не придется это делать.

Объект может иметь более одного списка свободных блоков. Если предполагается значительный объем действий INSERT или UPDATE в отношении объекта со стороны множества параллельно работающих пользователей, то конфигурирование для этого объекта нескольких списков свободных блоков может значительно увеличить его производительность (ценой возможного дополнительного пространства хранения). Наличие достаточного количества списков свободных блоков, покрывающего ваши потребности, является критически важным.

Списки свободных блоков могут оказать значительное положительное (или отрицательное) влияние на производительность в среде с многочисленными параллельными операциями вставки и обновления. Преимущества от корректно установленного количества списков свободных блоков можно продемонстрировать на очень простом примере. Для начала создадим табличное пространство, управляемое методом MSSM. Потребуется указать конструкцию SEGMENT SPACE MANAGEMENT MANUAL. В этом примере создается табличное пространство MSSM (по имени mssm) и затем таблица T внутри него:

```

ODA@ORA12CR1> create tablespace mssm
  2 datafile size 1m autoextend on next 1m
  3 segment space management manual;

```

Tablespace created.

Табличное пространство создано.

А теперь создадим простую таблицу:

```

ODA@ORA12CR1> create table t ( x int, y char(50) ) tablespace mssm;
Table created.

```

Таблица создана.

На заметку! В приведенном выше примере mssm — это имя табличного пространства, а не ключевое слово. Его можно заменить именем любого существующего табличного пространства, которое использует ручное управление пространством сегментов.

Мы начнем интенсивно вставлять в эту таблицу множество строк с применением пяти параллельных сеансов. Если мы измерим системные события ожиданий, связанные с блоками, до и после вставки, то обнаружим длительные периоды ожидания, особенно на блоках данных (попытка вставить данные). Подобное часто случается из-за недостаточного количества списков свободных блоков у таблиц (и индексов, но мы рассмотрим это более подробно в главе 11). Я использовал для этого пакет Statspack — получил statspack.snap, выполнил сценарий, который запустил пять параллельных сеансов SQL*Plus, и подождал, пока они все завершатся, после чего сделал еще один снимок statspack.snap. В сеансах выполнялся простой сценарий, который приведен ниже:

```

begin
  for i in 1 .. 1000000
  loop
    insert into t values ( i, 'x' );
  end loop;
  commit;
end;
/
exit;

```

Это очень простой блок кода, и здесь я являюсь единственным пользователем в базе данных. Значит, я должен получить наилучшую производительность, какую только возможно. Кеш буферов нужным образом сконфигурирован, журналы повтора имеют соответствующие размеры, индексы не замедляют функционирование, и я работаю на машине с двумя гиперпоточными процессорами Xeon — в общем, все должно выполняться быстро. Однако впоследствии обнаруживается такая ситуация:

| Snapshot | Snap Id | Snap Time | Sessions | Curs/Sess | Comment |
|------------------------|-------------|--------------------|-------------|-----------|---------|
| ~~~~~ | ----- | ----- | ----- | ----- | ----- |
| Begin Snap: | 195 | 27-Jan-14 11:58:12 | 36 | 1.0 | |
| End Snap: | 196 | 27-Jan-14 11:58:23 | 36 | 1.0 | |
| Elapsed: | 0.18 (mins) | Av Act Sess: | 4.7 | | |
| DB time: | 0.87 (mins) | DB CPU: | 0.31 (mins) | | |
| ... | | | | | |
| Top 5 Timed Events | | | | Avg | %Total |
| ~~~~~ | | | | wait | Call |
| Event | Waits | Time (s) | | (ms) | Time |
| ----- | ----- | ----- | | ----- | ----- |
| AQPC idle | 1 | 30 | | 30009 | 23.3 |
| buffer busy waits | 51,262 | 28 | | 1 | 21.5 |
| LGWR worker group idle | 599 | 20 | | 33 | 15.5 |
| CPU time | | 18 | | | 14.2 |
| lreg timer | 4 | 12 | | 3000 | 9.3 |
| ----- | ----- | ----- | | ----- | ----- |

На ожидание доступа к буферу (buffer busy waits) было потрачено 28 секунд, или в среднем чуть меньше 6 секунд на сеанс. Это ожидание было вызвано исключительно тем фактом, что при такой параллельной активности для таблицы сконфигурировано недостаточное количество списков свободных блоков. Уменьшить указанное время ожидания легко, просто создав таблицу с несколькими списками свободных блоков:

```
EODA@ORA12CR1> create table t ( x int, y char(50) )
2 storage ( freelists 5 ) tablespace MSSM;
Table created.
Таблица создана.
```

либо изменив объект:

```
EODA@ORA12CR1> alter table t storage ( FREELISTS 5 );
Table altered.
Таблица изменена.
```

Вы заметите, что снизилось не только значение buffer busy waits, но также значения CPU time (поскольку сейчас делается меньше работы; конкуренция за структуру данных, зафиксированную в защелке, может действительно перегрузить процессор) и Elapsed:

| Snapshot | Snap Id | Snap Time | Sessions | Curs/Sess | Comment |
|-------------------------|-------------|--------------------|-------------|-----------|---------|
| ~~~~~ | ----- | ----- | ----- | ----- | ----- |
| Begin Snap: | 197 | 27-Jan-14 12:07:05 | 36 | 1.0 | |
| End Snap: | 198 | 27-Jan-14 12:07:14 | 36 | 1.0 | |
| Elapsed: | 0.15 (mins) | Av Act Sess: | 4.2 | | |
| DB time: | 0.64 (mins) | DB CPU: | 0.25 (mins) | | |
| ... | | | | | |
| Top 5 Timed Events | | | | Avg | %Total |
| ~~~~~ | | | | wait | Call |
| Event | Waits | Time (s) | | (ms) | Time |
| ----- | ----- | ----- | | ----- | ----- |
| LGWR worker group idle | 346 | 279 | | 806 | 85.9 |
| CPU time | | 15 | | | 4.5 |
| log file parallel write | 344 | 9 | | 27 | 2.9 |
| lreg timer | 3 | 9 | | 3000 | 2.8 |
| heartbeat redo informer | 8 | 8 | | 1000 | 2.5 |
| ----- | ----- | ----- | | ----- | ----- |

Все, что вы хотите сделать для таблицы — это попытаться определить максимальное количество параллельных (по-настоящему параллельных) операций вставки или обновления, которые потребуют дополнительного пространства. Количество *по-настоящему параллельных* операций означает то, насколько часто вы ожидаете, что два пользователя будут в точности в один и тот же момент запрашивать свободный блок для данной таблицы. Это не единица измерения пересекающихся транзакций, а показатель количества сеансов, одновременно выполняющих операции вставки безотносительно к границам транзакций. Для увеличения степени параллелизма необходимо иметь примерно столько же списков свободных блоков, сколько производится параллельных операций вставки в таблицу.

Вы должны только сделать количество списков свободных блоков действительно большим и затем не беспокоиться об этом, правильно? Нет — разумеется, все не так просто. Когда списков свободных блоков много, среди них есть главный список и процессные списки. Если сегмент имеет единственный список свободных блоков, то главный список и процессный список представляют собой одно и то же. При наличии двух списков свободных блоков на самом деле будет один главный список свободных блоков и два процессных списка. Заданному сеансу будет назначаться одиночный процессный список свободных блоков на основе хеш-значения его идентификатора сеанса. Каждый процессный список свободных блоков будет содержать совсем немного блоков — оставшиеся свободные блоки находятся в главном списке свободных блоков. Во время работы процессный список свободных блоков по мере необходимости извлекает несколько блоков из главного списка. Если главный список свободных блоков не может удовлетворить требование к пространству, Oracle передвинет HWM-маркер и добавит пустые блоки в главный список свободных блоков. Таким образом, со временем главный список свободных блоков распределит свое пространство хранения между множеством процессных списков (каждый из которых будет иметь только небольшое число блоков). Итак, каждый процесс будет использовать единственный процессный список свободных блоков. Он не переходит от одного процессного списка к другому в поисках свободного пространства. Это означает, что когда таблица имеет 10 процессных списков свободных блоков, и один из процессов израсходовал все доступные блоки в своем списке, то он не будет искать нужное ему пространство в других процессных списках свободных блоков. Даже если каждый из оставшихся девяти процессных списков содержит по 5 блоков (всего 45 блоков), процесс все равно обратится к главному списку свободных блоков. Предполагая, что главный список не может удовлетворить запрос свободного блока, HWM-маркер таблицы будет перемещен либо, если сделать это невозможно (все пространство занято), таблица будет расширена (чтобы получить еще один экстен-тент). Затем процесс продолжит пользоваться пространством только в своем списке свободных блоков (который больше не является пустым). При наличии нескольких списков свободных блоков приходится идти на определенный компромисс. С одной стороны, применение множества списков свободных блоков существенно увеличивает производительность. С другой стороны, оно может привести к тому, что таблица будет занимать несколько больше места на диске, чем в действительности необходимо. Вы должны решить, что будет менее болезненным в вашей среде.

Не стоит недооценивать полезность параметра `FREELISTS`, особенно с учетом того, что в Oracle 8.1.6 и последующих версиях его можно изменять в большую и

меньшую сторону. Для этого параметра можно указать большое значение, когда необходимо произвести параллельную загрузку данных с помощью SQL*Loader в обычном режиме. Это позволит достичь высокого уровня параллелизма загрузки с минимальным временем ожидания. После загрузки данных значение параметра FREELISTS можно уменьшить до числа, подходящего для выполнения повседневных операций. При уменьшении пространства блоки во многих существующих списках свободных блоков будут объединены в один главный список.

Другой способ решения упомянутой ранее проблемы с высокими значениями buffer busy waits предусматривает использование табличного пространства, управляемого ASSM. Давайте возьмем предыдущий пример и создадим таблицу T в табличном пространстве ASSM:

```
EOA@ORA12CR1> create tablespace assm
  2 datafile size 1m autoextend on next 1m
  3 segment space management auto;
Tablespace created.
Табличное пространство создано.

EOA@ORA12CR1> create table t ( x int, y char(50) ) tablespace ASSM;
Table created.
Таблица создана.
```

Вы заметите, что значения buffer busy waits, CPU time и Elapsed в этом случае также снизились подобно тому, как это происходило при конфигурировании подходящего количества списков свободных блоков для сегмента, управляемого MSSM — и без необходимости в подсчете оптимального числа требуемых списков свободных блоков:

| Snapshot | Snap Id | Snap Time | Sessions | Curs/Sess | Comment |
|-------------------------|-------------|--------------------|-------------|-----------|---------|
| ~~~~~ | ----- | ----- | ----- | ----- | ----- |
| Begin Snap: | 199 | 27-Jan-14 12:16:30 | 33 | 1.0 | |
| End Snap: | 200 | 27-Jan-14 12:16:37 | 33 | 1.0 | |
| Elapsed: | 0.12 (mins) | Av Act Sess: | 5.3 | | |
| DB time: | 0.62 (mins) | DB CPU: | 0.25 (mins) | | |
| ... | | | | | |
| Top 5 Timed Events | | | | Avg | %Total |
| ~~~~~ | | | | wait | Call |
| Event | | Waits | Time (s) | (ms) | Time |
| ----- | | ----- | ----- | ----- | ----- |
| LGWR worker group idle | | 341 | 562 | 1647 | 92.9 |
| CPU time | | | 15 | | 2.4 |
| log file parallel write | | 341 | 10 | 29 | 1.6 |
| heartbeat redo informer | | 8 | 8 | 1000 | 1.3 |
| lreg timer | | 2 | 6 | 3000 | 1.0 |
| ----- | | ----- | ----- | ----- | ----- |

Одной из главных целей ASSM является исключение необходимости вручную определять корректные настройки для многих основных параметров хранения. По сравнению с MSSM механизм ASSM в некоторых случаях применяет дополнительное пространство, поскольку старается распределить операции вставки по множеству блоков, но почти во всех ситуациях факт использования дополнительного пространства хранения значительно перевешивается сокращением проблем с параллелизмом. По этой причине среда, в которой коэффициент использования пространства хра-

нения критичен, а параллелизм — нет (на ум сразу приходят хранилища данных), не обязательно выигрывает от хранилища, управляемого методом ASSM.

Параметры PCTFREE и PCTUSED

В общем случае параметр PCTFREE указывает Oracle, сколько пространства должно быть зарезервировано в блоке для будущих обновлений. По умолчанию он равен 10%. Если имеется больший процент свободного пространства, чем указано в PCTFREE, тогда этот блок считается *свободным*. Параметр PCTUSED указывает Oracle процент свободного пространства, который должен быть доступен в блоке, который не является в текущий момент свободным, для того, чтобы он снова стал свободным. Значение по умолчанию составляет 40%.

Как упоминалось ранее, когда параметр PCTFREE применяется с таблицей (но не с индекс-таблицей, что будет показано позже), он сообщает Oracle, какой объем пространства должен быть зарезервирован в блоке для будущих обновлений. Это означает, что если размер блока составляет 8 Кбайт, то в случае, когда добавление новой строки приводит к тому, что в блоке становится меньше 800 байт свободного пространства, Oracle будет использовать вместо этого блока другой из списка свободных блоков. Это и есть 10% пространства для обновления строк в этом блоке.

На заметку! Для разных типов таблиц параметры PCTFREE и PCTUSED реализованы по-разному. В одних типах таблиц задействованы оба параметра, тогда как в других — только PCTFREE, причем лишь во время создания объекта. При создании индекс-таблиц применяется параметр PCTFREE для выделения пространства под будущие обновления строк, однако он не используется, например, для принятия решения о прекращении вставки строк в заданный блок.

Точный эффект от этих двух параметров зависит от того, какое табличное пространство применяется — ASSM или MSSM. В табличном пространстве MSSM эти параметры управляют тем, когда блок будет помещен и изъят из списка свободных блоков. Если для PCTFREE и PCTUSED используются стандартные значения (соответственно, 10 и 40), то блок будет оставаться в списке свободных блоков до тех пор, пока не заполнится на 90% (в нем останется 10% свободного пространства). Заполнившись на 90%, он будет исключен из списка свободных блоков и не возвратится туда до тех пор, пока свободное пространство в нем не превысит 60% объема.

В табличном пространстве ASSM параметр PCTFREE по-прежнему ограничивает возможность вставки новой строки в блок, но не управляет тем, находится ли блок в списке свободных блоков, т.к. в методе ASSM списки свободных блоков вообще не применяются. Параметр PCTUSED в ASSM попросту игнорируется.

Для PCTFREE предусмотрены три типа значений: слишком высокое, слишком низкое и почти правильное. Установив параметр PCTFREE для блоков в слишком высокое значение, вы будете понапрасну расходовать пространство. Если вы установите PCTFREE в 50%, но никогда не будете обновлять данные, то 50% пространства каждого блока будут тратиться впустую. Однако для какой-то другой таблицы значение PCTFREE, равное 50%, может оказаться вполне подходящим. Если строки сначала занимают мало места и постепенно увеличиваются вдвое, то установка PCTFREE в слишком маленькое значение может послужить причиной перемещения строк при их обновлении.

Перемещение строк

Что такое перемещение строк? *Перемещение строк* (row migration) — это ситуация, когда строка вынуждена покинуть блок, в котором создавалась, поскольку она стала слишком большой, чтобы уместиться в нем вместе с остальными строками. Для иллюстрации перемещения строк начнем с блока, который выглядит, как показано на рис. 10.3.



Рис. 10.3. Блок данных перед обновлением

Почти одна седьмая блока является свободным пространством. Тем не менее, предположим, что посредством оператора UPDATE необходимо более чем вдвое увеличить объем пространства, используемого строкой 4 (в текущий момент она занимает одну седьмую часть блока). В таком случае, даже если Oracle объединит все доступное в этом блоке свободное пространство (рис. 10.4), места для увеличенной вдвое строки 4 все равно не хватит, потому что размер свободного пространства меньше текущего размера строки 4.



Рис. 10.4. Блок данных после объединения свободного пространства

Если бы строка уместилась в объединенное пространство, то объединение произошло бы. Однако в данной ситуации Oracle не будет выполнять такое объединение, и блок останется в том виде, в каком есть. Поскольку строке 4 пришлось бы выйти за пределы этого блока, если оставить ее в нем, Oracle переместит ее. Тем не менее, Oracle не может просто переместить эту строку — должен быть оставлен адрес для перенаправления. Дело в том, что могут существовать индексы, которые физически указывают на данный адрес для строки 4. Простое обновление не модифицирует также и эти индексы.

На заметку! Имеется особый случай: в секционированных таблицах идентификатор строки (rowid), т.е. адрес строки, будет изменяться. Мы рассмотрим этот случай в главе 13. Вдобавок другие административные операции, такие как FLASHBACK TABLE и ALTER TABLE SHRINK, также могут изменять назначенные идентификаторы строк.

Следовательно, когда база данных Oracle перемещает строку, она оставляет указатель на место, где эта строка действительно находится. После обновления блоки могут выглядеть так, как представлено на рис. 10.5.

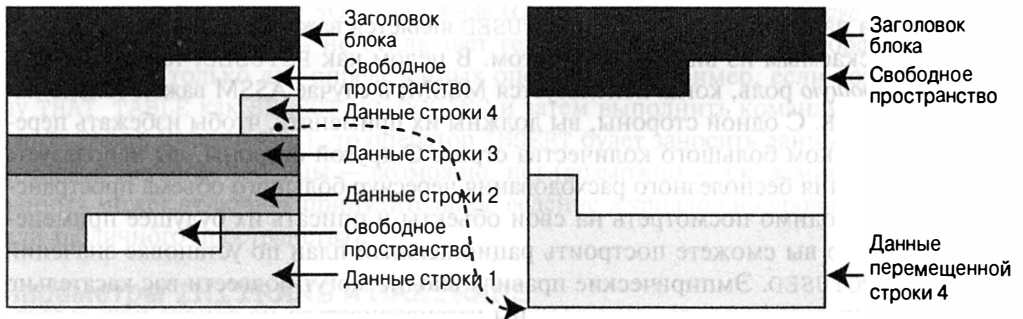


Рис. 10.5. Графическая иллюстрация перемещенной строки

Итак, *перемещенная строка* — это строка, которую пришлось переместить из блока, куда она была вставлена, в какой-то другой блок. Почему это может стать проблемой? Ваше приложение никогда о них не узнает; на SQL-операторах это тоже никак не отразится. Наличие таких строк имеет значение только с точки зрения производительности. При чтении перемещенной строки индекс будет указывать на исходный блок, а этот блок, в свою очередь — на новый блок. Вместо того чтобы выполнить две или около того операции ввода-вывода для чтения индекса плюс одну операцию для чтения таблицы, придется делать дополнительную операцию ввода-вывода, чтобы получить действительную строку данных. Когда речь идет об одной такой строке, сложностей не возникает — вы даже не заметите этого. Однако в ситуации, когда в подобном состоянии находится значительная часть строк и доступ к ним осуществляют многие пользователи, вы начнете замечать этот побочный эффект. Доступ к этим данным замедлится (время доступа увеличится из-за дополнительных операций ввода-вывода и связанного с ними защелкивания данных), снизится эффективность буферного кеша (в буфер должны помещаться два блока, а не один, как было бы при отсутствии перемещения строк), а также возрастет размер и сложность таблицы. По этим причинам вы вряд ли захотите иметь дело с перемещенными строками (но не лишайтесь сна из-за того, что пара сотен/тысяч строк в таблице, содержащей многие тысячи строк, окажутся перемещенными).

Интересно посмотреть, что Oracle будет делать, если строку, которая была перемещена из блока слева в блок справа на рис. 10.5, в какой-то момент в будущем *снова* придется переместить. Такая ситуация может произойти из-за добавления в этот блок других строк и обновления самой строки, в результате чего она станет еще больше по размеру.

База данных Oracle в действительности переместит эту строку *обратно* в исходный блок и, если в нем есть достаточно места, оставит ее там (строка может стать *неперемещенной*). Если же пространства недостаточно, Oracle переместит эту строку в другой блок и изменит ее адрес в *исходном* блоке. Таким образом, перемещение строк будет всегда включать в себя один уровень косвенности.

Вернемся к параметру PCTFREE и к тому, для чего он применяется: это настройка, которая в случае корректной установки поможет минимизировать количество перемещений строк.

Установка параметров PCTFREE и PCTUSED

Установка параметров PCTFREE и PCTUSED является важным — и в значительной степени упускаемым из виду — предметом. В целом как PCTUSED, так и PCTFREE играют *решающую* роль, когда используется MSSM; в случае ASSM важен только параметр PCTFREE. С одной стороны, вы должны их применять, чтобы избежать перемещения слишком большого количества строк. С другой стороны, вы используете их для устранения бесполезного расходования чересчур большого объема пространства. Вам необходимо посмотреть на свои объекты и описать их будущее применение, после чего вы сможете построить рациональный план по установке значений PCTFREE и PCTUSED. Эмпирические правила вполне могут подвести вас касательно этих настроек; они на самом деле должны устанавливаться на основе использования. Вы можете принять во внимание следующие конфигурации (памятуя о том, что термины “высокое значение” и “низкое значение” являются *относительными*, а в случае ASSM применим только параметр PCTFREE).

- **Высокое значение PCTFREE, низкое значение PCTUSED.** Такая конфигурация предназначена для случаев, когда вы вставляете большой объем данных, которые будут обновляться, и операции обновления будут часто увеличивать размер строк. Эта конфигурация резервирует много пространства в блоке после операций вставки (высокое значение PCTFREE) и обеспечивает то, что блок должен быть почти пустым перед возвращением в список свободных блоков (низкое значение PCTUSED).
- **Низкое значение PCTFREE, высокое значение PCTUSED.** Такая конфигурация предназначена для случаев, когда вы планируете производить только вставки или удаления из таблицы, или если будут делаться обновления, то в основном такие, которые сокращают размер строки.

И снова жесткие правила о том, что считается высоким значением, а что низким, для этих параметров отсутствуют. При установке PCTFREE и PCTUSED вам придется принимать во внимание поведение вашего приложения. Значение PCTFREE может варьироваться в пределах от 0 до 99. Высоким значением PCTFREE может считаться, скажем, 70, что соответствует резервированию 70% блока под обновления. Низким значением PCTFREE может быть, например, 5, которое говорит о том, что вы оставляете совсем немного места в блоке для будущих обновлений (это увеличивает размер строки). Параметр PCTUSED может принимать значения в пределах от 0 до 99. Высоким значением PCTFREE может считаться выбранное из диапазона 70–80. Низким значением PCTFREE может быть величина около 10.

Параметры LOGGING и NOLOGGING

Обычно объекты создаются с параметром LOGGING, т.е. все выполняемые над ними операции, которые могут генерировать информацию redo, будут ее генерировать. Параметр NOLOGGING позволяет определенным операциям выполняться в отношении объекта без генерации redo; эта тема была довольно подробно раскрыта в главе 9. Параметр NOLOGGING влияет только на несколько специфических операций, таких как начальное создание объекта, загрузка данных в прямом режиме с использованием SQL*Loader или перестройка (чтобы узнать, какие операции применимы к объекту базы данных, с которым вы работаете, обращайтесь в соответствующий раздел руководства по языку SQL для Oracle (*Oracle Database SQL Language Reference*)).

Параметр NOLOGGING не отключает генерацию журналов повтора для объекта в целом, а только для определенных операций. Например, если создать таблицу THAT_TABLE как SELECT NOLOGGING и затем выполнить команду INSERT INTO THAT_TABLE VALUES (1), то оператор INSERT будет заносить данные в журнал, но оператор создания таблицы — возможно, нет (возможно — т.к. администратор базы данных может включить принудительное ведение журналов на уровне базы данных или табличного пространства).

Параметры INITRANS и MAXTRANS

Каждый блок в сегменте имеет заголовок. Часть этого заголовка представляет собой таблицу транзакций. В эту таблицу будут заноситься записи, которые описывают, какие строки/элементы в этом блоке транзакциями заблокированы. Начальный размер этой таблицы определяется параметром INITRANS для объекта (для таблиц и индексов этот параметр имеет стандартное значение 2). Таблица транзакций будет по мере необходимости динамически увеличиваться до тех пор, пока количество записей в ней не достигнет числа, указанного в параметре MAXTRANS (т.е. пока в блоке есть достаточный объем свободного пространства). Каждая выделенная запись транзакции может потреблять 23 или 24 байта в заголовке блока. Обратите внимание, что, начиная с Oracle 10g, параметр MAXTRANS игнорируется — для всех сегментов он равен 255.

Традиционные таблицы

Традиционные таблицы используются в приложениях вероятно 99% (или больше) времени. Традиционная таблица — это тип таблиц, который вы получаете по умолчанию, когда выдаете оператор CREATE TABLE. При желании создать таблицу любого другого типа вы должны специально указать нужный тип в операторе CREATE.

Куча (heap) — это классическая структура данных в вычислительной технике. Она по существу представляет собой большую область пространства на диске или в памяти (конечно же, диск отражает случай таблицы базы данных), которая управляется в предположительно произвольной манере. Данные будут помещаться туда, где они лучше всего умещаются, а не в каком-то специфическом порядке. Многие ожидают, что данные будут поступать из таблицы в том же порядке, в каком они туда заносились, но в случае кучи это определенно не гарантируется. На самом деле гарантируется противоположное: строки будут извлекаться в совершенно непредсказуемом порядке. Продемонстрировать это довольно легко.

В следующем примере я создам таблицу, которая позволяет в моей базе данных содержать в каждом блоке (размер блока составляет 8 Кбайт) одну полную строку. Вам необязательно делать так, чтобы в каждом блоке могла уместиться только одна строка: я поступаю подобным образом исключительно ради того, чтобы продемонстрировать предсказуемую последовательность событий. Показанное ниже поведение (отсутствие порядка у строк) можно наблюдать на таблицах всех размеров, в базах данных с любым размером блоков:

```

EODA@ORA12CR1> create table t
  2 ( a int,
  3   b varchar2(4000) default rpad(' ',4000, '*'),
  4   c varchar2(3000) default rpad(' ',3000, '*')
  5 )
  6 /
Table created.
Таблица создана.
EODA@ORA12CR1> insert into t (a) values ( 1 );
1 row created.
1 строка создана.
EODA@ORA12CR1> insert into t (a) values ( 2 );
1 row created.
1 строка создана.
EODA@ORA12CR1> insert into t (a) values ( 3 );
1 row created.
1 строка создана.
EODA@ORA12CR1> delete from t where a = 2 ;
1 row deleted.
1 строка удалена.
EODA@ORA12CR1> insert into t (a) values ( 4 );
1 row created.
1 строка создана.
EODA@ORA12CR1> select a from t;
-----
A
-----
1
4
3

```

Если вы хотите воспроизвести этот пример, настройте столбцы В и С в соответствии со своим размером блока. Например, для размера блока 2 Кбайт столбец С не нужен, а столбец В должен иметь тип `varchar2(1500) default rpad(' ',1500, '*')`. Поскольку в таблице наподобие этой данные управляются в куче, то после того, как пространство становится доступным, оно применяется повторно.

На заметку! При использовании метода ASSM или MSSM вы обнаружите, что строки попадают в разные места. Лежащие в основе этих методов процедуры управления пространством существенно отличаются, поэтому одни и те же операции, выполняемые в отношении таблицы в ASSM и MSSM, могут давать результаты с разным физическим порядком. Логически данные будут теми же самыми, но будут храниться по-разному.

При полном сканировании таблицы данные будут извлекаться в том порядке, в каком они встречаются, а не в порядке вставки. С таблицами базы данных связана важная ключевая концепция: в общем случае они представляют собой неупорядоченные коллекции данных. Также обратите внимание, что для наблюдения этого эффекта применять оператор `DELETE` не обязательно; тех же результатов можно добиться с использованием *только* операторов `INSERT`. Если я вставлю сначала небольшую строку, потом очень большую строку, которая не сможет уместиться в том же блоке вместе с небольшой строкой, а затем вставлю снова небольшую строку, то с высокой вероятностью смогу увидеть, что строки по умолчанию будут поступать в порядке “небольшая строка, небольшая строка, большая строка”. Строки не будут извлекаться в порядке их вставки — Oracle размещает данные там, где они уместаются, а не в порядке по дате или по транзакции.

Если запросу необходимо извлекать данные в порядке их вставки, к таблице понадобится добавить столбец и применять его для упорядочивания данных во время извлечения. Это может быть, например, числовой столбец, управляемый с помощью увеличивающейся последовательности (объекта `SEQUENCE`). Далее можно *приблизительно оценить* порядок вставки строк, используя оператор `SELECT` с конструкцией `ORDER BY` на этом столбце. Результат будет приближением, потому что строка с номером в последовательности 55 вполне может быть зафиксирована раньше строки с номером 54, вследствие чего она формально окажется в базе данных первой.

Вы должны воспринимать традиционную таблицу как большую неупорядоченную коллекцию строк. Эти строки будут возвращаться, по-видимому, в случайном порядке, а в зависимости от других применяемых опций (параллельный запрос, различные режимы оптимизатора и т.д.) могут возвращаться в разном порядке даже для одного и того же запроса. Никогда не рассчитывайте получить строки в каком-то определенном порядке, если только запрос не включает конструкцию `ORDER BY`!

Помимо сказанного, что еще важно знать о традиционных таблицах? Описание синтаксиса оператора `CREATE TABLE` в руководстве *Oracle Database SQL Language Reference* от Oracle занимает 87 страниц, так что есть еще немало опций. Опций настолько много, что охватить их все довольно трудно. Описание одних лишь коммутационных диаграмм (или *диаграмм цепочек*) находится на 20 страницах. Чтобы увидеть большинство опций, доступных в операторе `CREATE TABLE` для какой-то таблицы, я создаю таблицу настолько просто, насколько возможно, например:

```
EODA@ORA12CR1> create table t
2  ( x int primary key,
3    y date,
4    z clob
5  )
6  /
Table created.
Таблица создана.
```

Затем с использованием стандартного пакета `DBMS_METADATA` я запрашиваю описание этой таблицы и получаю расширенную версию синтаксиса:

```
EODA@ORA12CR1> select dbms_metadata.get_ddl( 'TABLE', 'T' ) from dual;
DBMS_METADATA.GET_DDL('TABLE','T')
-----
```



```

CREATE TABLE "EODA"."T"
(
  "X" NUMBER(*,0),
  "Y" DATE,
  "Z" CLOB,
  PRIMARY KEY ("X")
  USING INDEX PCTFREE 10 INITRANS 2 MAXTRANS 255
  TABLESPACE "USERS" ENABLE
) SEGMENT CREATION DEFERRED
PCTFREE 10 PCTUSED 40 INITRANS 1 MAXTRANS 255
NOCOMPRESS LOGGING
TABLESPACE "USERS"
LOB ("Z") STORE AS SECUREFILE (
  TABLESPACE "USERS" ENABLE STORAGE IN ROW CHUNK 8192
  NOCACHE LOGGING NOCOMPRESS KEEP_DUPLICATES )

```

Такой прием позволяет просмотреть множество опций для оператора CREATE TABLE. Понадобится только указать типы данных и тому подобное, а Oracle сгенерирует расширенную версию синтаксиса этого оператора. Затем можно настроить расширенную версию, скажем, изменив ENABLE STORAGE IN ROW на DISABLE STORAGE IN ROW, что отключит возможность хранения LOB-данных в строке со структурированными данными, в результате чего они будут храниться в отдельном сегменте. Я применяю этот прием постоянно, чтобы избежать расшифровки огромных коммутационных диаграмм. Кроме того, он позволяет узнать опции, доступные для оператора CREATE TABLE в разных обстоятельствах.

Теперь, когда вы знаете, как увидеть доступные для заданного оператора CREATE TABLE опции, возникает вопрос: какие из них являются наиболее важными для традиционных таблиц? По моему мнению, таких опций три в случае ASSM и пять — в случае MSSM.

- **FREELISTS.** Только для MSSM. Каждая таблица управляет блоками, которые были выделены для нее в куче, через список свободных блоков. Таблица может иметь несколько списков свободных блоков. Если ожидается интенсивная вставка в таблицу множеством параллельно работающих пользователей, то конфигурирование для таблицы нескольких списков свободных блоков может существенно улучшить производительность (ценой возможного расхода дополнительного пространства). Более подробно об этом рассказывалось в разделе “Списки свободных блоков” ранее в главе.
- **PCTFREE.** Для ASSM и MSSM. Степень заполнения блока, допускаемая во время процесса INSERT. Как упоминалось ранее, это применяется для управления тем, может ли строка быть добавлена к блоку, на основе того, насколько он полон в текущий момент. Эта опция также управляет перемещениями строк, вызванными последовательными операциями обновления, и должна устанавливаться на базе использования таблицы.
- **PCTUSED.** Только для MSSM. Показывает, насколько пустым должен стать блок, прежде чем он снова станет кандидатом на вставку данных. Блок, в котором занято меньше пространства, чем указано в PCTUSED, становится кандидатом для вставки новых строк. Как и PCTFREE, значение для PCTUSED должно подбираться на основе использования таблицы.

- **INITRANS.** Для ASSM и MSSM. Количество слотов транзакций, изначально выделенных блоку. Если для INITRANS установлено слишком низкое значение (стандартное значение равно 2), могут возникнуть проблемы с параллелизмом в блоке, к которому получают доступ много пользователей. Если блок базы данных почти заполнен, и список транзакций не может быть динамически расширен, то сеансы будут выстраиваться в очередь к этому блоку, т.к. каждой параллельной транзакции необходим слот транзакции. Если предполагается наличие множества параллельных операций обновления одних и тех же блоков, рассмотрите возможность увеличения INITRANS.
- **COMPRESS/NOCOMPRESS.** Для ASSM и MSSM. Включает или отключает сжатие табличных данных на время выполнения операций либо прямого маршрута, либо обычного маршрута (“нормальных” операций, если хотите), таких как INSERT. До выхода версии Oracle9i Release 2 эта опция была недоступной. Начиная с Oracle9i Release 2 и до Oracle 10g Release 2, стала доступной опция COMPRESS или NOCOMPRESS для включения или отключения сжатия таблиц на период выполнения только операций прямого маршрута. В этих выпусках преимуществами сжатия могли воспользоваться только операции прямого маршрута, такие как CREATE TABLE AS SELECT, INSERT /*+ APPEND */ , ALTER TABLE T MOVE, и прямые загрузки SQL*Loader. В Oracle 11g Release 1 и последующих версиях доступны опции NOLOGGING, COMPRESS FOR OLTP и COMPRESS BASIC. Опция NOLOGGING отключает любое сжатие, COMPRESS FOR OLTP включает сжатие для всех операций (прямого или обычного маршрута), а COMPRESS BASIC включает сжатие только для операций прямого маршрута. Начиная с версии Oracle 12c Release 1, опции сжатия синтаксически указываются как ROW STORE COMPRESS BASIC (включает сжатие на время выполнения операций прямого маршрута) и ROW STORE COMPRESS ADVANCED (включает сжатие для всех операций).

На заметку! На LOB-данные, хранящиеся в отдельном LOB-сегменте, действие параметров PCTFREE и PCTUSED, установленных для таблицы, не распространяется. Такие блоки LOB управляются по-другому: они всегда заполняются до отказа и возвращаются в список свободных блоков только когда становятся совершенно пустыми.

На описанные выше параметры следует обращать особое внимание. С появлением локально управляемых табличных пространств, которые настоятельно рекомендуется применять, все остальные параметры хранения (такие как PCTINCREASE, NEXT и т.д.), на мой взгляд, больше не являются существенными.

Индекс-таблицы

Индекс-таблицы — это довольно простые таблицы, которые хранятся в индексной структуре. В то время как в традиционной таблице данные никак не организованы (т.е. попадают туда, где есть свободное пространство), в индекс-таблице данные сохранены и отсортированы по первичному ключу. В том, что касается приложений, индекс-таблицы ведут себя точно так же как и “нормальные” таблицы; для доступа к ним используются обычные SQL-операторы. Они особенно подходят для информационно-поисковых, пространственных и OLAP-приложений.

Какая польза от индекс-таблицы? На самом деле лучше поинтересоваться пользой от традиционной таблицы. Поскольку предполагается, что все таблицы реляционной базы данных в любом случае должны иметь первичный ключ, не является ли традиционная таблица просто неэффективным расходом пространства? При работе с традиционной таблицей необходимо выделить место и под саму таблицу, и под индекс по первичному ключу таблицы. В ситуации с индекс-таблицей накладные расходы, связанные с пространством под индекс для первичного ключа, исчезают, т.к. индекс — это данные, а данные — это индекс. Индекс представляет собой сложную структуру данных, для управления и поддержания которой необходимо проводить немалую работу, и требования по его обслуживанию растут по мере увеличения ширины строки, подлежащей хранению. С другой стороны, управлять кучей гораздо проще. Таким образом, традиционные таблицы обладают более эффективными характеристиками, чем индекс-таблицы. Однако индекс-таблицам присущи определенные преимущества по сравнению с традиционными таблицами. Например, однажды мне пришлось строить инвертированный списочный индекс по текстовым данным (это было до появления *interMedia* и связанных технологий). Я имел дело с таблицей, заполненной документами; нужно было выполнить разбор этих документов и найти в них слова. Таблица выглядела следующим образом:

```
create table keywords
( word varchar2(50),
  position int,
  doc_id int,
  primary key(word,position,doc_id)
);
```

Эта таблица состоит исключительно из столбцов первичного ключа. Накладные расходы составили *более 100%*; размеры таблицы и индекса по первичному ключу оказались сравнимыми (в действительности, индекс по первичному ключу был больше, т.к. в нем физически хранились идентификаторы строк, на которые он указывал, а в таблице идентификаторы строк не сохраняются — они выводятся). Таблица применялась только в конструкции *WHERE* со столбцом *WORD* или столбцами *WORD* и *POSITION*. Это значит, что я никогда не использовал таблицу, а только индекс на ней. Сама таблица была не более чем накладными расходами. Требовалось найти все документы, содержащие заданное слово (или похожее на него слово и т.д.). Традиционная таблица *KEYWORDS* была бесполезной; она только замедляла работу приложения во время обслуживания таблицы *KEYWORDS* и вдвое увеличивала требования к хранилищу. Здесь индекс-таблица пришлась бы очень кстати.

Еще одна реализация, при которой уместно применять индекс-таблицу — справочная таблица для кодов. Например, может существовать справочная таблица для определения штата по почтовому индексу. В такой ситуации следует отказаться от традиционной таблицы и просто работать с самой индекс-таблицей. Любая таблица, доступ в которую производится исключительно через ее первичный ключ, является кандидатом на то, чтобы быть индекс-таблицей.

Когда желательно обеспечить близкое расположение данных либо их физическое хранение в специфичном порядке, то индекс-таблица является именно той структурой, которая нужна. Пользователи баз данных *Sybase* и *SQL Server* в таком случае использовали бы кластеризованный индекс, но индекс-таблицы в этом смыс-

ле лучше. Кластеризованный индекс в таких базах данных может повлечь за собой накладные расходы вплоть до 110% (подобно тому, как было в примере с таблицей KEYWORDS). В случае индекс-таблиц накладные расходы нулевые, т.к. данные сохраняются только один раз. Классическим примером, когда может понадобиться физически близкое расположение данных, является ситуация с отношением “родительская–дочерняя”. Предположим, что таблица сотрудников EMP имеет дочернюю таблицу адресов. В таблице адресов может быть записан домашний адрес сотрудника, добавленный, когда ему в начале было отправлено письмо с приглашением на работу. Позже сотрудник добавляет свой рабочий адрес. Затем по причине смены места жительства сотрудник вводит новый домашний адрес. Также сотрудник может добавить адрес учебного заведения, в котором он повышал квалификацию, и т.д. То есть сотрудник может иметь три, четыре и даже более записей с адресами, которые добавляются произвольным образом с течением времени. В нормальной традиционной таблице эти записи будут размещены вразброс. Шансы, что хотя бы две из них будут расположены в одном и том же блоке базы данных, практически равны нулю. Однако, запрашивая информацию о сотруднике, вы всегда также будете получать все записи с адресами этого сотрудника. Строки, добавляемые в разное время, всегда извлекаются вместе. Чтобы сделать операцию извлечения более эффективной, в качестве типа дочерней таблицы с адресами может быть выбрана индекс-таблица, благодаря чему все записи для каждого сотрудника при вставке будут размещаться близко друг к другу, что снижает объем работ по их извлечению.

Увидеть преимущества использования индекс-таблицы для физически близкого расположения информации в дочерней таблице поможет пример. Давайте создадим таблицу EMP и заполним ее данными:

```
EODA@ORA12CR1> create table emp
2 as
3 select object_id      empno,
4        object_name    ename,
5        created        hiredate,
6        owner          job
7   from all_objects
8 /
```

Table created.

Таблица создана.

```
EODA@ORA12CR1> alter table emp add constraint emp_pk primary key(empno);
```

Table altered.

Таблица изменена.

```
EODA@ORA12CR1> begin
2   dbms_stats.gather_table_stats( user, 'EMP', cascade=>true );
3 end;
4 /
```

PL/SQL procedure successfully completed.

Процедура PL/SQL успешно завершена.

А теперь реализуем дочернюю таблицу в двух вариантах: как обычную традиционную таблицу и как индекс-таблицу:

```

EODA@ORA12CR1> create table heap_addresses
 2 ( empno    references emp(empno) on delete cascade,
 3   addr_type varchar2(10),
 4   street    varchar2(20),
 5   city      varchar2(20),
 6   state     varchar2(2),
 7   zip       number,
 8   primary key (empno,addr_type)
 9 )
10 /
Table created.
Таблица создана.

```

```

EODA@ORA12CR1> create table iot_addresses
 2 ( empno    references emp(empno) on delete cascade,
 3   addr_type varchar2(10),
 4   street    varchar2(20),
 5   city      varchar2(20),
 6   state     varchar2(2),
 7   zip       number,
 8   primary key (empno,addr_type)
 9 )
10 ORGANIZATION INDEX
11 /
Table created.
Таблица создана.

```

Я заполнил эти таблицы информацией, вставив в них для каждого сотрудника рабочий адрес, домашний адрес, предыдущий адрес и адрес учебного заведения. В традиционной таблице данные будут размещаться ближе к концу таблицы; по мере поступления данные будут добавляться в конец таблицы из-за того, что данные только поступают, но ничего не удаляется. Со временем, если какие-то адреса будут удалены, новые вставляемые адреса станут размещаться в этой таблице еще более случайно. Достаточно сказать, что шансы нахождения рабочего и домашнего адресов сотрудника в одном и том же блоке в традиционной таблице практически равны нулю. Однако поскольку в индекс-таблице применяется первичный ключ по EMPNO и ADDR_TYPE, можно иметь уверенность в том, что все адреса для любого заданного сотрудника EMPNO будут находиться в одном или, возможно, двух расположенных рядом индексных блоках. Для заполнения этих таблиц адресами использовались следующие операторы INSERT:

```

EODA@ORA12CR1> insert into heap_addresses
 2 select empno, 'WORK', '123 main street', 'Washington', 'DC', 20123
 3 from emp;
72075 rows created.
72075 строк создано.

EODA@ORA12CR1> insert into iot_addresses
 2 select empno, 'WORK', '123 main street', 'Washington', 'DC', 20123
 3 from emp;
72075 rows created.
72075 строк создано.

```

Я поступил подобным образом еще трижды, изменяя WORK по очереди на HOME, PREV и SCHOOL. После этого я собрал статистику:

```
EODA@ORA12CR1> exec dbms_stats.gather_table_stats( user, 'HEAP_ADDRESSES' );
PL/SQL procedure successfully completed.
```

Процедура PL/SQL успешно завершена.

```
EODA@ORA12CR1> exec dbms_stats.gather_table_stats( user, 'IOT_ADDRESSES' );
PL/SQL procedure successfully completed.
```

Процедура PL/SQL успешно завершена.

Теперь можно посмотреть на измеримые отличия с применением средства AUTOTRACE:

```
EODA@ORA12CR1> set autotrace traceonly
EODA@ORA12CR1> select *
  2   from emp, heap_addresses
  3   where emp.empno = heap_addresses.empno
  4   and emp.empno = 42;
```

Execution Plan

Plan hash value: 775524973

| Id | Operation | Name | Rows | Bytes | Cost (%CPU) |
|-----|-------------------------------------|----------------|------|-------|-------------|
| 0 | SELECT STATEMENT | | 4 | 292 | 8 (0) |
| 1 | NESTED LOOPS | | 4 | 292 | 8 (0) |
| 2 | TABLE ACCESS BY INDEX ROWID | EMP | 1 | 27 | 2 (0) |
| * 3 | INDEX UNIQUE SCAN | EMP_PK | 1 | | 1 (0) |
| 4 | TABLE ACCESS BY INDEX ROWID BATCHED | HEAP_ADDRESSES | 4 | 184 | ... |
| * 5 | INDEX RANGE SCAN | SYS_C0032863 | 4 | | 2 (0) |

Predicate Information (identified by operation id):

```
3 - access("EMP"."EMPNO"=42)
5 - access("HEAP_ADDRESSES"."EMPNO"=42)
```

Statistics

```
1 recursive calls
0 db block gets
11 consistent gets
0 physical reads
0 redo size
1361 bytes sent via SQL*Net to client
543 bytes received via SQL*Net from client
2 SQL*Net roundtrips to/from client
0 sorts (memory)
0 sorts (disk)
4 rows processed
```

Это довольно распространенный план: переход в таблицу EMP по первичному ключу, получение строки, переход в таблицу адресов с использованием EMPNO и выбор дочерних записей с применением индекса. Для извлечения этих данных было выполнено 11 операций ввода-вывода. Теперь запустим тот же самый запрос, но к индекс-таблице адресов:

```
EODA@ORA12CR1> select *
2   from emp, iot_addresses
3   where emp.empno = iot_addresses.empno
4   and emp.empno = 42;
```

Execution Plan

Plan hash value: 252066017

| Id | Operation | Name | Rows | Bytes | Cost (%CPU) | Time |
|-----|-----------------------------|--------------------|------|-------|-------------|----------|
| 0 | SELECT STATEMENT | | 4 | 292 | 4 (0) | 00:00:01 |
| 1 | NESTED LOOPS | | 4 | 292 | 4 (0) | 00:00:01 |
| 2 | TABLE ACCESS BY INDEX ROWID | EMP | 1 | 27 | 2 (0) | 00:00:01 |
| * 3 | INDEX UNIQUE SCAN | EMP_PK | 1 | | 1 (0) | 00:00:01 |
| * 4 | INDEX RANGE SCAN | SYS_IOT_TOP_182459 | 4 | 184 | 2 (0) | 00:00:01 |

Predicate Information (identified by operation id):

- 3 - access("EMP"."EMPNO"=42)
- 4 - access("IOT_ADDRESSES"."EMPNO"=42)

Statistics

- 1 recursive calls
- 0 db block gets
- 7 consistent gets**
- 0 physical reads
- 0 redo size
- 1361 bytes sent via SQL*Net to client
- 543 bytes received via SQL*Net from client
- 2 SQL*Net roundtrips to/from client
- 0 sorts (memory)
- 0 sorts (disk)
- 4 rows processed

В этом случае операций ввода-вывода стало на четыре меньше (должно быть понятно, что это за четыре операции); были пропущены четыре шага TABLE ACCESS (BY INDEX ROWID BATCHED). Чем больше имеется дочерних записей, тем больше будет пропускаться операций ввода-вывода.

Итак, что это за четыре операции ввода-вывода? В данном случае они составили треть всего объема ввода-вывода, выполненного для запроса, а если выдать этот запрос повторно, то количество таких операций увеличится.

Каждая операция ввода-вывода и операция согласованного чтения требует доступа к кешу буферов, и хотя чтение данных из кеша буферов происходит быстрее, чем с диска, правда также и то, что операции чтения из кеша буферов *не являются бесплатными и обходятся далеко не дешево в смысле ресурсов*. Каждая из них требует множества зашелок на буферном кеше, а зашелки — это механизм сериализации, который будет ограничивать возможности масштабирования. Мы можем измерить снижение числа операций ввода-вывода и зашелок, выполнив показанный ниже блок кода PL/SQL:

```

EODA@ORA12CR1> begin
2   for x in ( select empno from emp )
3   loop
4       for y in ( select emp.ename, a.street, a.city, a.state, a.zip
5                   from emp, heap_addresses a
6                   where emp.empno = a.empno
7                   and emp.empno = x.empno )
8       loop
9           null;
10          end loop;
11      end loop;
12 end;
13 /

```

PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.

Здесь мы эмулируем период занятости и запускаем запрос приблизительно 72 000 раз, по одному разу для каждого значения EMPNO. Если выполнить этот блок кода для таблиц HEAP_ADDRESSES и IOT_ADDRESSES, то утилита TKPROF отобразит следующие сведения:

```

SELECT EMP.ENAME, A.STREET, A.CITY, A.STATE, A.ZIP
FROM EMP, HEAP_ADDRESSES A WHERE EMP.EMPNO = A.EMPNO AND EMP.EMPNO = :B1

```

| call | count | cpu | elapsed | disk | query | current | rows |
|---------|--------|------|---------|------|--------|---------|--------|
| Parse | 1 | 0.00 | 0.00 | 0 | 0 | 0 | 0 |
| Execute | 72110 | 1.02 | 1.01 | 0 | 0 | 0 | 0 |
| Fetch | 72110 | 2.16 | 2.11 | 0 | 722532 | 0 | 288440 |
| total | 144221 | 3.18 | 3.12 | 0 | 722532 | 0 | 288440 |

...

| Rows (1st) | Rows (avg) | Rows (max) | Row Source Operation |
|---------------------------------|------------|------------|--|
| 4 | 4 | 4 | NESTED LOOPS (cr=10 pr=0 pw=0 time=40 |
| us cost=8 size=228 card=4) | | | |
| 1 | 1 | 1 | TABLE ACCESS BY INDEX ROWID EMP (cr=3 |
| pr=0 pw=0 time=11 us cost=2... | | | |
| 1 | 1 | 1 | INDEX UNIQUE SCAN EMP_PK (cr=2 pr=0 |
| pw=0 time=7 us cost=1 size=0... | | | |
| 4 | 4 | 4 | TABLE ACCESS BY INDEX ROWID BATCHED |
| HEAP_ADDRESSES (cr=7... | | | |
| 4 | 4 | 4 | INDEX RANGE SCAN SYS_C0032863 (cr=3 pr=0 |
| pw=0 time=10 us cost=2... | | | |

```
SELECT EMP.ENAME, A.STREET, A.CITY, A.STATE, A.ZIP
FROM EMP, IOT_ADDRESSES A WHERE EMP.EMPNO = A.EMPNO AND EMP.EMPNO = :B1
```

| call | count | cpu | elapsed | disk | query | current | rows |
|---------|--------|------|---------|------|--------|---------|--------|
| Parse | 1 | 0.00 | 0.00 | 0 | 0 | 0 | 0 |
| Execute | 72110 | 1.04 | 1.01 | 0 | 0 | 0 | 0 |
| Fetch | 72110 | 1.64 | 1.63 | 0 | 437360 | 0 | 288440 |
| total | 144221 | 2.69 | 2.64 | 0 | 437360 | 0 | 288440 |

```
...
Rows (1st) Rows (avg) Rows (max) Row Source Operation
-----
4 4 4 NESTED LOOPS (cr=7 pr=0 pw=0 time=28 us
cost=4 size=228 card=4)
1 1 1 TABLE ACCESS BY INDEX ROWID EMP (cr=3
pr=0 pw=0 time=11 us cost=2...)
1 1 1 INDEX UNIQUE SCAN EMP_PK (cr=2 pr=0
pw=0 time=7 us cost=1 size=0...)
4 4 4 INDEX RANGE SCAN SYS_IOT_TOP_182459
(cr=4 pr=0 pw=0 time=15 us...)
Rows Row Source Operation
```

```
4 NESTED LOOPS (cr=7 pr=3 pw=0 time=9 us cost=4 size=280 card=4)
1 TABLE ACCESS BY INDEX ROWID EMP (cr=3 pr=0 pw=0 time=0 us cost=2 size=30...)
1 INDEX UNIQUE SCAN EMP_PK (cr=2 pr=0 pw=0 time=0 us cost=1 size=0...)
4 INDEX RANGE SCAN SYS_IOT_TOP_93124 (cr=4 pr=3 pw=0 time=3 us cost=2...)
```

Оба запроса извлекли в точности то же количество строк, но в таблице HEAP_ADDRESSES происходило гораздо больше логического ввода-вывода. Поскольку степень параллелизма возросла, вполне можно ожидать, что таблица HEAP_ADDRESSES более часто использовала процессор, а запрос в это время, скорее всего, ожидал защелки в буферном кеше. С помощью runstat (утилита моей собственной разработки; за подробностями обращайтесь в раздел “Настройка среды” в самом начале книги) можно оценить разницу в количестве защелок. В своей системе я увидел такой вывод:

| Name | Run1 | Run2 | Diff |
|--|---------------|---------------|----------------|
| STAT...buffer is pinned count | 216,342 | 0 | -216,342 |
| STAT...consistent gets | 723,461 | 438,275 | -285,186 |
| STAT...consistent gets from ca | 723,461 | 438,275 | -285,186 |
| STAT...consistent gets pin (fa | 362,888 | 77,700 | -285,188 |
| STAT...consistent gets pin | 362,888 | 77,700 | -285,188 |
| STAT...no work - consistent re | 362,870 | 77,682 | -285,188 |
| STAT...session logical reads | 723,538 | 438,332 | -285,206 |
| STAT...table fetch by rowid | 360,570 | 72,114 | -288,456 |
| STAT...buffer is not pinned co | 649,026 | 288,456 | -360,570 |
| STAT...session pga memory | 393,216 | 0 | -393,216 |
| STAT...session pga memory max | 393,216 | 0 | -393,216 |
| LATCH.cache buffers chains | 1,091,314 | 518,788 | -572,526 |
| STAT...logical read bytes from | 5,927,223,296 | 3,590,815,744 | -2,336,407,552 |
| Run1 latches total versus runs -- difference and pct | | | |
| Run1 | Run2 | Diff | Pct |
| 1,235,153 | 620,581 | -614,572 | 199.03% |

Здесь Run1 относится к таблице HEAP_ADDRESSES, а Run2 — к таблице IOT_ADDRESSES. Как видите, разница в количестве защелок значительна и заметно снижение этого количества в результате повторения, в основном из-за защелки цепочек кеша буферов (защищающей буферный кеш). В рассматриваемом случае индекс-таблица обеспечивает перечисленные ниже преимущества.

- Возросшая эффективность буферного кеша, потому что любому заданному запросу приходится иметь в кеше меньше блоков.
- Сниженная потребность в доступе к буферному кешу, что улучшает масштабируемость.
- Меньший объем работы для извлечения данных, т.к. данные извлекаются гораздо быстрее.
- Меньший объем физического ввода-вывода на каждый запрос, поскольку для любого заданного запроса требуется меньше отдельных блоков и одна физическая операция ввода-вывода адресов, скорее всего, извлечет их все (а не только один, как в случае традиционной таблицы).

То же самое было бы справедливо в ситуации, когда часто применяются запросы BETWEEN по первичному или уникальному ключу. Наличие данных, физически хранящихся в отсортированном виде, увеличит производительность этих запросов. Например, в своей базе данных я поддерживаю таблицу с котировками акций. Каждый день для сотен акций я собираю вместе такие данные, как тикер, дата, цена на момент закрытия биржи, максимальную цену за день, минимальную цену за день, количество акций и другую связанную информацию. Эта таблица выглядит следующим образом:

```
EODA@ORA12CR1> create table stocks
2  ( ticker   varchar2(10),
3    day      date,
4    value    number,
5    change   number,
6    high     number,
7    low      number,
8    vol      number,
9    primary key(ticker,day)
10 )
11 organization index
12 /
Table created.
Таблица создана.
```

Я часто просматриваю сведения об одной акции за определенный диапазон дней (например, для вычисления скользящего среднего). Если бы использовалась традиционная таблица, то вероятность существования двух строк для тикера ORCL в одном и том же блоке базы данных была бы практически нулевой. Причина в том, что каждый вечер я вставляю записи за день для всех акций. Это полностью заполняет, по крайней мере, один блок в базе данных (на самом деле много блоков). Таким образом, ежедневно добавляется новая запись ORCL, которая сохраняется в другом блоке, а в не том же, где хранится любая существующая запись ORCL. Выполним следующий запрос:

```
Select * from stocks
where ticker = 'ORCL'
and day between sysdate-100 and sysdate;
```

База данных Oracle прочитает индекс и затем произведет доступ к таблице по идентификатору строки, чтобы получить остальные данные строки. Каждая из ста извлекаемых строк будет находиться в другом блоке базы данных из-за способа загрузки данных в таблицу — каждая, возможно, потребует физической операции ввода-вывода. Теперь представим, что те же самые данные хранятся в индекс-таблице. Тому же самому запросу придется лишь прочитать соответствующие блоки индекса, которые уже содержат все данные. При этом не только устраняется потребность в доступе к таблице, но еще и все строки для тикера ORCL в заданном диапазоне дат будут физически храниться близко друг к другу. Понадобится меньший объем логического и физического ввода-вывода.

Теперь вы понимаете, когда может возникать необходимости в индекс-таблицах и как с ними работать. Следующее, что нужно узнать — это опции, доступные для таких таблиц. В чем здесь подвох? Опции очень похожи на опции для традиционных таблиц. Чтобы ознакомиться с деталями, давайте снова воспользуемся пакетом DBMS_METADATA. Начнем с создания трех базовых вариаций индекс-таблицы:

```
EODA@ORA12CR1> create table t1
2 ( x int primary key,
3   y varchar2(25),
4   z date
5 )
6 organization index;
Table created.
Таблица создана.
```

```
EODA@ORA12CR1> create table t2
2 ( x int primary key,
3   y varchar2(25),
4   z date
5 )
6 organization index
7 OVERFLOW;
Table created.
Таблица создана.
```

```
EODA@ORA12CR1> create table t3
2 ( x int primary key,
3   y varchar2(25),
4   z date
5 )
6 organization index
7 overflow INCLUDING y;
Table created.
Таблица создана.
```

Позже мы выясним, что делают опции OVERFLOW и INCLUDING, а пока взглянем на детализированную версию SQL-оператора для создания первой таблицы:

```

EODA@ORA12CR1> select dbms_metadata.get_ddl( 'TABLE', 'T1' ) from dual;
DBMS_METADATA.GET_DDL('TABLE','T1')
-----
CREATE TABLE "EODA"."T1"
(
  "X" NUMBER(*,0),
  "Y" VARCHAR2(25),
  "Z" DATE,
  PRIMARY KEY ("X") ENABLE
) ORGANIZATION INDEX NOCOMPRESS PCTFREE 10 INITRANS 2 MAXTRANS 255 LOGGING
STORAGE(INITIAL 65536 NEXT 1048576 MINEXTENTS 1 MAXEXTENTS 2147483645
PCTINCREASE 0 FREELISTS 1 FREELIST GROUPS 1
BUFFER_POOL DEFAULT FLASH_CACHE DEFAULT CELL_FLASH_CACHE DEFAULT)
TABLESPACE "USERS"
PCTTHRESHOLD 50

```

Здесь встречается новая опция PCTTHRESHOLD, которая вскоре будет описана. Вы могли заметить, что по сравнению с синтаксисом предыдущего оператора CREATE TABLE здесь кое-чего не хватает — отсутствует конструкция PCTUSED, хотя PCTFREE на месте. Причина в том, что индекс — это сложная структура данных, которая не организуется произвольным образом подобно куче, так что данные должны размещаться там, где им место. В отличие от кучи, где блоки доступны для вставок лишь иногда, в индексе блоки доступны для добавления новых записей всегда. Если данные относятся к определенному блоку из-за своих значений, то они поступят в этот блок независимо от того, насколько он полон или пуст. Вдобавок параметр PCTFREE применяется только тогда, когда объект создается и заполняется данными в индексной структуре. Здесь он используется не так, как в традиционной таблице: PCTFREE резервирует пространство для вновь созданного индекса, а не для последующих операций над ним. Те же соображения, которые были сделаны касательно списков свободных блоков в традиционных таблицах, в целом применимы к индекс-таблицам.

Для начала рассмотрим опцию NOCOMPRESS. Эта опция отличается по реализации от сжатия таблиц, которое обсуждалось ранее. Она работает для любой операции над индекс-таблицей (в противоположность сжатию таблиц, которое может как действовать, так и нет, для операций прямого маршрута). Наличие NOCOMPRESS сообщает Oracle о необходимости сохранять каждое значение в индексной записи (т.е. не сжимать его). Если бы первичный ключ объекта был организован на столбцах A, B и C, то физически сохранялось бы каждое вхождение A, B и C. Противоположностью NOCOMPRESS является COMPRESS N, где N — целое число, представляющее количество столбцов для сжатия. Эта опция удаляет повторяющиеся значения и факторизует (выносит) их на уровне блоков, так что, например, повторяющиеся значения A и возможно B больше физически не сохраняются. Взгляните, к примеру, на следующую таблицу:

```

EODA@ORA12CR1> create table iot
2 ( owner, object_type, object_name,
3   primary key(owner,object_type,object_name)
4 )
5 organization index
6 NOCOMPRESS
7 as
8 select distinct owner, object_type, object_name from all_objects
/
Table created.

```

Если задуматься, то значение OWNER повторяется много сотен раз. Каждая схема (OWNER) склонна владеть множеством объектов. Даже пара значений OWNER и OBJECT_TYPE повторяется много раз, поэтому заданная схема будет иметь десятки таблиц, пакетов и т.д. Никогда не повторяются только все три столбца вместе. Мы можем заставить Oracle подавлять эти повторяющиеся значения. Вместо индексного блока со значениями, показанными в табл. 10.1, мы можем использовать конструкцию COMPRESS 2 (факторизуем два первых столбца) и получить блок со значениями, представленными в табл. 10.2.

Таблица 10.1. Листовой блок индекса, NOCOMPRESS

| | | | |
|------------------|------------------|------------------|------------------|
| Sys, table, t1 | Sys, table, t2 | Sys, table, t3 | Sys, table, t4 |
| Sys, table, t5 | Sys, table, t6 | Sys, table, t7 | Sys, table, t8 |
| ... | ... | ... | ... |
| Sys, table, t100 | Sys, table, t101 | Sys, table, t102 | Sys, table, t103 |

Таблица 10.2. Листовой блок индекса, COMPRESS 2

| | | | |
|------------|------|------|------|
| Sys, table | t1 | t2 | t3 |
| t4 | t5 | ... | ... |
| ... | t103 | t104 | ... |
| t300 | t301 | t302 | t303 |

Таким образом, значения SYS и TABLE появляются только один раз, а сохраняется только третий столбец. В данном случае мы можем иметь в индексном блоке гораздо больше записей, чем было бы иначе. При этом степень параллелизма не увеличивается — все операции по-прежнему выполняются на уровне строк — и функциональность не меняется. Эта опция *может* задействовать немного больше вычислительных возможностей процессора, т.к. Oracle придется выполнять дополнительную работу по объединению ключей. С другой стороны, опция может значительно сократить количество операций ввода-вывода и позволить находиться в кеше буферов большему объему данных, поскольку каждый блок будет содержать больше данных. Мы достигаем вполне разумного компромисса.

Мы продемонстрируем указанную экономию, проведя краткий тест предыдущего оператора CREATE TABLE с опциями NOCOMPRESS, COMPRESS 1 и COMPRESS 2. Для начала создадим индекс-таблицу без сжатия:

```
EODA@ORA12CR1> create table iot
2  ( owner, object_type, object_name,
3    constraint iot_pk primary key(owner,object_type,object_name)
4  )
5  organization index
6  NOCOMPRESS
7  as
8  select distinct owner, object_type, object_name
9    from all_objects
10 /
Table created.
Таблица создана.
```

Теперь можно измерить используемое пространство. Для этого мы применим команду `ANALYZE INDEX VALIDATE STRUCTURE`. Эта команда заполняет динамическое представление производительности по имени `INDEX_STATS`, которое будет содержать самое большое одну строку с информацией из последнего выполнения команды `ANALYZE`:

```
EODA@ORA12CR1> analyze index iot_pk validate structure;
Index analyzed.
Индекс проанализирован.
```

```
EODA@ORA12CR1> select lf_blks, br_blks, used_space,
2      opt_cmpr_count, opt_cmpr_pctsave
3      from index_stats;
```

| LF_BLKs | BR_BLKs | USED_SPACE | OPT_CMPR_COUNT | OPT_CMPR_PCTSAVE |
|---------|---------|------------|----------------|------------------|
| 240 | 1 | 1726727 | 2 | 37 |

Вывод показывает, что индекс в текущий момент использует 240 листовых блоков (где находятся наши данные) и 1 блок ветвления (такие блоки в Oracle применяются для навигации по структуре индекса) для нахождения листовых блоков. Используемое пространство составляет около 1,7 Мбайт (1 726 727 байтов). Есть еще два столбца со странными названиями, которые тоже пытаются о чем-то сообщить. Столбец `OPT_CMPR_COUNT` (optimum compression count — оптимальная степень сжатия) говорит следующее: “Если вы сделаете этот индекс `COMPRESS 2`, то сможете получить лучшую степень сжатия”. Столбец `OPT_CMPR_PCTSAVE` (optimum compression percentage saved — процент экономии при оптимальной степени сжатия) говорит о том, что если указать `COMPRESS 2`, то можно было бы сэкономить примерно одну треть объема хранилища, а индекс занимал бы всего две трети того дискового пространства, которое он занимает сейчас.

На заметку! Структура индексов более подробно рассматривается в главе 11.

Чтобы проверить эту теорию, перестроим индекс-таблицу с опцией `COMPRESS 1`:

```
EODA@ORA12CR1> alter table iot move compress 1;
Table altered.
Таблица изменена.
```

```
EODA@ORA12CR1> analyze index iot_pk validate structure;
Index analyzed.
Индекс проанализирован.
```

```
EODA@ORA12CR1> select lf_blks, br_blks, used_space,
2      opt_cmpr_count, opt_cmpr_pctsave
3      from index_stats;
```

| LF_BLKs | BR_BLKs | USED_SPACE | OPT_CMPR_COUNT | OPT_CMPR_PCTSAVE |
|---------|---------|------------|----------------|------------------|
| 213 | 1 | 1529506 | 2 | 28 |

Как видите, индекс действительно стал меньше: приблизительно 1,5 Мбайт с меньшим числом листовых блоков. Но столбец `OPT_CMPR_PCTSAVE` все равно говорит о том, что можно высвободить еще 22% пространства, поскольку мы так и не указали оптимальную степень сжатия.

Давайте перестроим индекс-таблицу с опцией COMPRESS 2:

```
EODA@ORA12CR1> alter table iot move compress 2;
```

Table altered.

Таблица изменена.

```
EODA@ORA12CR1> analyze index iot_pk validate structure;
```

Index analyzed.

Индекс проанализирован.

```
EODA@ORA12CR1> select lf_blks, br_blks, used_space,
```

```
2      opt_cmpr_count, opt_cmpr_pctsave
```

```
3      from index_stats;
```

| LF_BLKs | BR_BLKs | USED_SPACE | OPT_CMPR_COUNT | OPT_CMPR_PCTSAVE |
|---------|---------|------------|----------------|------------------|
| 151 | 1 | 1086483 | 2 | 0 |

Теперь мы значительно сократили размер за счет как количества листовых блоков, так и общего объема занимаемого пространства до примерно 1 Мбайт. Если вспомнить исходные цифры:

```
EODA@ORA12CR1> select (1-.37) * 1726727 from dual;
```

```
(1-.37) * 1726727
```

```
-----
1087838.01
```

то несложно заметить, что столбец OPT_CMPR_PCTSAVE ничуть нас не обманывал. Предыдущий пример раскрывает один очень интересный факт об индекс-таблицах: они являются таблицами, но только по названию; их сегмент на самом деле представляет собой индексный сегмент.

В этом месте я собираюсь отложить обсуждение опции PCTTHRESHOLD, т.к. она связана со следующими двумя опциями для индекс-таблиц — OVERFLOW и INCLUDING. Давайте взглянем на расширенную версию SQL-синтаксиса двух наборов таблиц T2 и T3 (я воспользовался процедурой DBMS_METADATA для подавления конструкций хранения, потому что в этом примере они не важны):

```
EODA@ORA12CR1> begin
```

```
2      dbms_metadata.set_transform_param
```

```
3      ( DBMS_METADATA.SESSION_TRANSFORM, 'STORAGE', false );
```

```
4 end;
```

```
5 /
```

PL/SQL procedure successfully completed.

Процедура PL/SQL успешно завершена.

```
EODA@ORA12CR1> select dbms_metadata.get_ddl( 'TABLE', 'T2' ) from dual;
```

```
DBMS_METADATA.GET_DDL('TABLE','T2')
```

```
-----
CREATE TABLE "EODA"."T2"
```

```
(  "X" NUMBER(*,0),
```

```
   "Y" VARCHAR2(25),
```

```
   "Z" DATE,
```

```
   PRIMARY KEY ("X") ENABLE
```

```
) ORGANIZATION INDEX NOCOMPRESS PCTFREE 10 INITRANS 2 MAXTRANS 255 LOGGING
TABLESPACE "USERS"
```

```
PCTTHRESHOLD 50 OVERFLOW
PCTFREE 10 PCTUSED 40 INITRANS 1 MAXTRANS 255 LOGGING
TABLESPACE "USERS"

EODA@ORA12CR1> select dbms_metadata.get_ddl( 'TABLE', 'T3' ) from dual;
DBMS_METADATA.GET_DDL('TABLE','T3')
```

```
CREATE TABLE "EODA"."T3"
(
  "X" NUMBER(*,0),
  "Y" VARCHAR2(25),
  "Z" DATE,
  PRIMARY KEY ("X") ENABLE
) ORGANIZATION INDEX NOCOMPRESS PCTFREE 10 INITRANS 2 MAXTRANS 255 LOGGING
TABLESPACE "USERS"
PCTTHRESHOLD 50 INCLUDING "Y" OVERFLOW
PCTFREE 10 PCTUSED 40 INITRANS 1 MAXTRANS 255 LOGGING
TABLESPACE "USERS"
```

Итак, нам осталось рассмотреть опции PCTTHRESHOLD, OVERFLOW и INCLUDING. Они связаны между собой и цель их состоит в том, чтобы листовые блоки индекса (блоки, которые содержат фактические данные индекса) могли эффективно хранить данные. Индекс обычно включает подмножество столбцов. Обычно в блоке индекса будет во много раз больше записей для строк, чем в блоке традиционной таблицы. Индекс рассчитывает на возможность получить много строк на блок. В противном случае на обслуживание индекса у Oracle уходила бы масса времени, т.к. каждая вставка или обновление могла приводить к разделению блока индекса для учета новых данных.

Конструкция OVERFLOW позволяет настроить еще один сегмент переполнения (превращая индекс-таблицу в многосегментный объект во многом подобно тому, как это делает наличие столбца CLOB), в который могут перемещаться строки, когда они становятся слишком большими.

На заметку! Данные столбцов, входящих в состав первичного ключа, не могут находиться в сегментах переполнения: они должны размещаться непосредственно в листовых блоках.

Обратите внимание, что OVERFLOW позволяет применять с индекс-таблицей конструкцию PCTUSED, когда используется MSSM. Конструкции PCTFREE и PCTUSED имеют для сегмента переполнения то же значение, что и для традиционной таблицы. Условия, когда должен применяться сегмент переполнения, могут быть заданы двумя способами.

- PCTTHRESHOLD. Если количество данных в строке превышает указанный здесь в процентах объем пространства блока, заключительные столбцы этой строки будут сохранены в сегменте переполнения. То есть, если для PCTTHRESHOLD задано 10%, а размер блоков составляет 8 Кбайт, то любая строка, занимающая более 800 байтов в длину, будет частично сохранена в каком-то другом месте, за пределами блока индекса.
- INCLUDING. Все столбцы, включая тот, который указан в конструкции INCLUDING, сохраняются в блоке индекса, а оставшиеся столбцы помещаются в сегмент OVERFLOW.

Для примера рассмотрим следующую таблицу, размер блоков в которой составляет 2 Кбайт:

```

EODA@ORA12CR1> create table iot
2 ( x int,
3   y date,
4   z varchar2(2000),
5   constraint iot_pk primary key (x)
6 )
7 organization index
8 pctthreshold 10
9 overflow
10 /

```

Table created.

Таблица создана.

Графически эту таблицу можно изобразить так, как показано на рис. 10.6.



Рис. 10.6. Индекс-таблица с сегментом переполнения, конструкция PCTTHRESHOLD

Прямоугольники серого цвета — это записи индекса, представляющие собой часть более крупной структуры индекса (индексы подробно рассматриваются в главе 11). Вкратце структура индекса представляет собой дерево, а листовые блоки (хранящие данные) в действительности являются двусвязным списком, упрощающим перемещение по узлам по порядку после того, как мы найдем в индексе место, с которого хотим начать. Прямоугольником белого цвета обозначен сегмент OVERFLOW. Именно здесь будут сохраняться данные, объем которых превышает порог, указанный в параметре PCTTHRESHOLD. Чтобы выяснить, какие столбцы должны быть сохранены в сегменте переполнения, Oracle будет обрабатывать столбцы в обратном порядке от последнего столбца в строке до завершающего столбца в первичном ключе, не включая его. В рассматриваемом примере числовой столбец X и столбец даты Y всегда будут уместиться в блоке индекса. Последний столбец, Z, имеет варьирующуюся длину. Когда он меньше приблизительно 190 байтов (10% от блока 2 Кбайт составляет около 200 байтов, из которых нужно вычесть 7 байтов для даты и от 3 до 5 байтов для числа), он будет сохраняться в блоке индекса. Когда он превышает 190 байтов, Oracle будет сохранять данные столбца Z в сегменте переполнения и настраивать указатель на него (фактически идентификатор строки).

Другой способ предполагает использование конструкции `INCLUDING`. С ее помощью мы явно указываем, какие столбцы необходимо хранить в блоке индекса, а какие должны быть сохранены в сегменте переполнения. Для примера создадим следующую таблицу:

```
EODA@ORA12CR1> create table iot
2 ( x int,
3   y date,
4   z varchar2(2000),
5   constraint iot_pk primary key (x)
6 )
7 organization index
8 including y
9 overflow
10 /
Table created.
Таблица создана.
```

Графическое представление этой таблицы приведено на рис. 10.7.

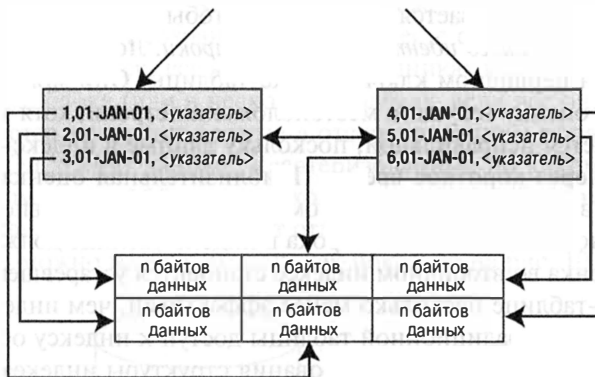


Рис. 10.7. Индекс-таблица с сегментом переполнения, конструкция `INCLUDING`

В этой ситуации столбец `Z` всегда будет сохраняться за пределами блока индекса, в сегменте переполнения, независимо от размера содержащихся в нем данных (все столбцы, не входящие в первичный ключ, которые следуют за столбцом, указанным в конструкции `INCLUDING`, сохраняются в сегменте переполнения).

Что же тогда лучше: `PCTTHRESHOLD`, `INCLUDING` либо комбинация этих конструкций? Это зависит от ваших потребностей. При наличии приложения, которое всегда (или почти всегда) получает доступ к первым четырем столбцам таблицы и редко — к последним пяти столбцам, лучше применять `INCLUDING`. Вы можете включить столбцы вплоть до четвертого и позволить остальным пяти храниться за пределами блока индекса. Если во время выполнения эти пять столбцов понадобятся, они будут извлекаться во многом подобно тому, как извлекается расщепленная строка. База данных Oracle будет читать заголовок строки, указатель на остаток строки и затем читать эти данные. С другой стороны, если вы не можете утверждать, что приложение почти всегда получает доступ к одним столбцам и очень редко — к другим, то лучше использовать конструкцию `PCTTHRESHOLD`.

После определения количества строк, которые в среднем желательно сохранять в каждом блоке индекса, установка `PCTTHRESHOLD` осуществляется легко. Предположим, что вы хотите иметь 20 строк на блок индекса. Это означает, что каждая строка должна составлять одну двадцатую (5%) блока. Параметр `PCTTHRESHOLD` должен быть установлен в 5 и каждая порция строки, которая остается в листовом блоке индекса, должна занимать не более 5% блока.

Последняя деталь, касающаяся индекс-таблиц, которую необходимо принять во внимание — это индексация. Допускается иметь индекс на самих индекс-таблицах — своего рода индекс на индексе. Такие индексы называются *вторичными индексами*. Обычно индекс содержит физический адрес строки, на которую он указывает, т.е. идентификатор строки (`rowid`). Вторичный индекс индекс-таблицы поступать так не может; для обращения к строке он должен применять какой-то другой способ. Причина в том, что строка в индекс-таблице может часто перемещаться, и она это делает не так, как в традиционной таблице. Ожидается, что строка в индекс-таблице должна находиться в определенной позиции внутри структуры индекса, зависящей от ее значения первичного ключа; строка будет перемещаться только из-за изменения размера и формы самого индекса. (Поддержание структур индексов более подробно раскрывается в главе 11.) Чтобы обеспечить это, в Oracle было введено понятие *логического идентификатора строки*. Логические идентификаторы строк основаны на первичном ключе индекс-таблицы. Они могут также содержать приблизительную оценку текущего местоположения строки, хотя такая оценка почти всегда оказывается неправильной, поскольку данные в индекс-таблице склонны к перемещению через короткое время. Приблизительная оценка — это физический адрес строки в индекс-таблице на момент, когда она была впервые помещена в структуру вторичного индекса. Если строка в индекс-таблице должна переместиться в другой блок, оценка во вторичном индексе становится устаревшей. Следовательно, индекс на индекс-таблице несколько менее эффективен, чем индекс на традиционной таблице. В случае традиционной таблицы доступ к индексу обычно требует одной операции ввода-вывода для сканирования структуры индекса и одного чтения для считывания данных таблицы. В случае индекс-таблицы обычно выполняются два сканирования: одно для структуры вторичного индекса и одно для самой индекс-таблицы. В остальном индексы на индекс-таблицах обеспечивают быстрый и эффективный доступ к данным в индекс-таблицах с использованием столбцов, не относящихся к первичному ключу.

Заключительные соображения по поводу индекс-таблиц

Выбор правильной комбинации данных, размещаемых в листовых блоках индекса, и данных, находящихся в сегменте переполнения, является самой важной частью процесса настройки индекс-таблицы. Оцените разнообразные сценарии с различными условиями переполнения и посмотрите, какое влияние они оказывают на выдаваемые запросы `INSERT`, `UPDATE`, `DELETE` и `SELECT`. Если структура строится однажды, но читается часто, размещайте в блоке индекса как можно больше данных. Если структура часто модифицируется, вы должны постараться достичь баланса между размещением всех данных в блоке индекса (хорошо для извлечения) и частой реорганизацией данных в индексе (плохо для модификаций). Соображения по поводу списков свободных блоков, касающиеся традиционных таблиц, применимы также и к индекс-таблицам. Параметры `PCTFREE` и `PCTUSED` в индекс-таблице играют две

роли. Параметр PCTFREE для индекс-таблицы почти так же важен, как для традиционной таблицы, а параметр PCTUSED обычно в игру не вступает. Однако когда речь заходит о сегменте переполнения, параметры PCTFREE и PCTUSED интерпретируются для индекс-таблицы точно так же, как для традиционной таблицы; устанавливайте их в соответствии с той же самой логикой, что и для традиционной таблицы.

Кластеризованные индекс-таблицы

По моим наблюдениям, люди часто не совсем правильно понимают, что собой представляет кластер в Oracle. Многие склонны путать его с “кластеризованным индексом” из SQL Server или Sybase. Это совсем не одно и то же. Кластер представляет способ хранения группы таблиц, которые совместно используют ряд общих столбцов, в одних и тех же блоках базы данных, а также хранения связанных данных вместе в том же самом блоке. Кластеризованный индекс в SQL Server вынуждает строки храниться в отсортированном порядке согласно ключу индекса подобно только что описанным индекс-таблицам. В случае кластера один блок данных может содержать данные из многих таблиц. Концептуально данные сохраняются как “предварительно соединенные”. Кластер также может применяться с одиночными таблицами, где он позволяет сохранять данные вместе сгруппированными по определенному столбцу. Например, все данные о сотрудниках в отделе 10 могут храниться в одном и том же блоке (или в нескольких блоках, если все они не помещаются в один блок). Кластер не хранит данные в отсортированном порядке — это задача индекс-таблицы. Он хранит данные, кластеризованные по какому-то ключу, но в куче. Таким образом, отдел 100 может находиться рядом с отделом 1 и в то же время очень далеко (физически на диске) от отделов 101 и 99.

Графически это можно изобразить так, как показано на рис. 10.8.

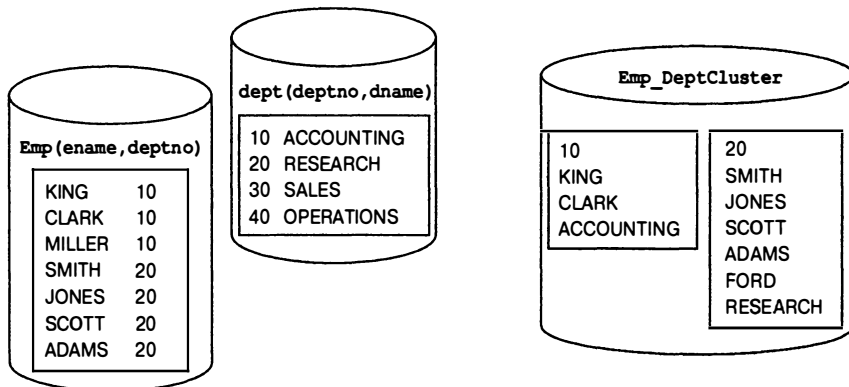


Рис. 10.8. Кластеризованные индексные данные

В левой части этого рисунка находятся обычные таблицы. Таблица EMP хранится в своем сегменте, как и таблица DEPT. Они могут располагаться в разных файлах и разных табличных пространствах, и они определенно будут пребывать в отдельных экстентах. В правой части рисунка видно, что произойдет, если объединить эти две таблицы в кластер. Прямоугольники представляют блоки базы данных. Значение 10 факторизуется и сохраняется один раз. Затем в этом блоке сохраняются все данные

из всех таблиц кластера для отдела 10. Если все данные для отдела 10 не умещаются в этот блок, то с ним сцепляются дополнительные блоки, которые будут содержать переполнение; это похоже на блоки переполнения в индекс-таблицах.

Итак, давайте посмотрим, как можно построить кластеризованный объект. Создать кластер таблиц в объекте довольно легко. Определение хранилища для объекта (PCTFREE, PCTUSED, INITIAL и т.д.) связывается с кластером, а не таблицами. Это имеет смысл, т.к. в кластере будет много таблиц, и они будут находиться в том же самом блоке. Наличие разных значений PCTFREE не имело бы смысла. Следовательно, оператор CREATE CLUSTER во многом похож на CREATE TABLE с небольшим количеством столбцов (только столбцы кластерного ключа):

```
EODA@ORA12CR1> create cluster emp_dept_cluster
2 ( deptno number(2) )
3 size 1024
4 /
Cluster created.
Кластер создан.
```

Здесь мы создали *индексный кластер* (другим типом является *хеш-кластер*, который более подробно рассматривается сразу после раздела “Заключительные соображения по поводу кластеризованных индекс-таблиц”). Столбцом кластерного ключа для этого кластера будет DEPTNO. Чтобы соответствовать этому определению, столбцы в таблицах не обязаны называться DEPTNO, но они *должны* иметь тип NUMBER(2). В определении кластера присутствует опция SIZE 1024. Она используется для сообщения Oracle о том, что мы ожидаем ассоциирования с каждым значением кластерного ключа около 1024 байтов данных. Эту информацию Oracle применяет для подсчета *максимального* количества кластерных ключей, которые могут уместиться в блоке. Учитывая размер блока 8 Кбайт, получится до семи кластерных ключей на блок базы данных (но может быть и меньше, если размер данных окажется больше ожидаемого). Например, данные для отделов 10, 20, 30, 40, 50, 60 и 70, скорее всего, попадут в один блок, а при вставке данных для отдела 80 будет использоваться новый блок. Это вовсе не означает, что данные сохраняются в отсортированном виде; дело просто в том, что если мы вставляем данные отделов в таком порядке, то они естественным образом помещаются вместе. В случае вставки данных отделов в порядке 10, 80, 20, 30, 40, 50, 60 и затем 70 данные последнего отдела (70) вероятно окажутся в новом добавленном блоке. Как вы увидите далее, на количество ключей, которые могут быть сохранены в блоке, будут влиять размер данных и порядок их вставки.

Таким образом, параметр SIZE управляет максимальным количеством кластерных ключей на блок. Он оказывает наибольшее влияние на утилизацию пространства кластера. Если установить для него очень высокое значение, мы получим слишком мало ключей на блок, и будет расходоваться больше пространства, чем необходимо. Если установить для него очень низкое значение, то данные будут чрезмерно сцепляться, что уведет в сторону от главного предназначения кластера, которое заключается в хранении всех данных вместе в одном блоке. Этот параметр является наиболее важным для кластера.

Далее нам необходимо проиндексировать кластер, прежде чем мы сможем помещать в него данные. Мы могли бы создать таблицы в кластере прямо сейчас, но мы собираемся создавать и заполнять таблицы одновременно, поэтому нуждаемся в кластерном индексе *до того, как* иметь дело с любыми данными. Работа кластерного

индекса состоит в том, чтобы взять значение кластерного ключа и вернуть адрес блока, который содержит этот ключ. В сущности, он представляет собой первичный ключ, где каждое значение кластерного ключа указывает на одиночный блок в самом кластере. То есть, когда мы запрашиваем данные для отдела 10, база данных Oracle прочитает кластерный ключ, определит адрес блока для него и затем прочитает данные. Индекс по кластерному ключу создается следующим образом:

```

EODA@ORA12CR1> create index emp_dept_cluster_idx
  2 on cluster emp_dept_cluster
  3 /
Index created.
Индекс создан.

```

Он может иметь все нормальные параметры хранения индекса и может располагаться в другом табличном пространстве. Это просто обычный индекс, так что он может быть создан на множестве столбцов; просто так случилось, что он организован в кластере, и он может также включать запись для полностью нулевого значения (об этом будет рассказываться в главе 11). Обратите внимание, что в приведенном операторе CREATE INDEX мы *не* указываем список столбцов — он выводится из самого определения кластера. Теперь все готово к созданию таблиц в кластере:

```

EODA@ORA12CR1> create table dept
  2 ( deptno    number(2) primary key,
  3   dname     varchar2(14),
  4   loc       varchar2(13)
  5 )
  6 cluster emp_dept_cluster(deptno)
  7 /
Table created.
Таблица создана.

```

```

EODA@ORA12CR1> create table emp
  2 ( empno     number primary key,
  3   ename     varchar2(10),
  4   job       varchar2(9),
  5   mgr       number,
  6   hiredate  date,
  7   sal       number,
  8   comm      number,
  9   deptno    number(2) references dept(deptno)
 10 )
 11 cluster emp_dept_cluster(deptno)
 12 /
Table created.
Таблица создана.

```

Единственное отличие от нормальной таблицы связано с применением ключевого слова CLUSTER и сообщением Oracle о том, какой столбец базовой таблицы будет отображаться на кластерный ключ в самом кластере. Вспомните, что кластер здесь является сегментом, следовательно, данная таблица никогда не будет иметь атрибуты сегмента, такие как TABLESPACE, PCTFREE и т.д. — это атрибуты кластерного сегмента, а не только что созданной таблицы. Теперь можно загрузить в эти таблицы начальный набор данных:

```

EODA@ORA12CR1> insert into dept
2  ( deptno, dname, loc )
3  select deptno+r, dname, loc
4    from scott.dept,
5       (select level r from dual connect by level < 10);

36 rows created.
36 строк создано.

EODA@ORA12CR1> insert into emp
2  (empno, ename, job, mgr, hiredate, sal, comm, deptno)
3  select rownum, ename, job, mgr, hiredate, sal, comm, deptno+r
4    from scott.emp,
5       (select level r from dual connect by level < 10);

126 rows created.
126 строк создано.

```

На заметку! В этом примере для генерации данных использовался трюк SQL. Для демонстрации того, что Oracle ограничит количество ключей отделов на блок на основе параметра `SIZE`, понадобилось более семи отделов. Таким образом, необходимо более четырех строк отделов в `SCOTT.DEPT`. Я сгенерировал девять строк с применением приема `connect by level` в отношении таблицы `DUAL` и выполнил декартово соединение этих девяти строк с четырьмя строками из `DEPT`, что дало в результате 36 уникальных строк. Подобный трюк был проделан с `EMP` для генерации данных по этим отделам.

Теперь, когда данные загружены, давайте ознакомимся с их организацией на диске. С помощью пакета `DBMS_ROWID` мы заглянем в идентификатор строки и увидим, какие блоки данных там находятся. Для начала исследуем таблицу `DEPT` и узнаем, сколько имеется строк `DEPT` в блоке:

```

EODA@ORA12CR1> select min(count(*)), max(count(*)), avg(count(*))
2    from dept
3    group by dbms_rowid.rowid_block_number(rowid)
4    /

```

| MIN(COUNT(*)) | MAX(COUNT(*)) | AVG(COUNT(*)) |
|---------------|---------------|---------------|
| ----- | ----- | ----- |
| 1 | 7 | 6 |

Итак, хотя таблица `DEPT` загружалась данными первой — и строки `DEPT` очень маленькие (в блок 8 Кбайт их могут уместиться сотни) — мы видим, что максимальное количество строк `DEPT` в блоке этой таблицы равно всего семи. Это соответствует тому, чего мы ожидали ранее, когда устанавливали `SIZE` в 1024. Мы рассчитывали, что с блоком 8 Кбайт и 1024 байтами данных на кластерный ключ для комбинированных записей `EMP` и `DEPT` мы получим примерно семь уникальных значений кластерного ключа на блок, и это в точности то, что мы здесь видим. А теперь посмотрим на таблицы `EMP` и `DEPT` вместе. Мы взглянем на идентификаторы строк каждой из этих таблиц и сравним количества блоков после соединения по `DEPTNO`. Если количества блоков окажутся теми же самыми, мы будем знать, что строка `EMP` и строка `DEPT` хранятся в одном физическом блоке базы данных вместе, а если количества будут отличаться, то нет. В этом случае мы наблюдаем, что все данные хранятся должным образом. Отсутствуют ситуации, когда запись из таблицы `EMP` хранится в блоке отдельно от соответствующей ей записи `DEPT`:

```

EODA@ORA12CR1> select *
2   from (
3   select dept_blk, emp_blk,
4           case when dept_blk <> emp_blk then '*' end flag,
5           deptno
6   from (
7   select dbms_rowid.rowid_block_number(dept.rowid) dept_blk,
8          dbms_rowid.rowid_block_number(emp.rowid) emp_blk,
9          dept.deptno
10  from emp, dept
11  where emp.deptno = dept.deptno
12  )
13  )
14  where flag = '*'
15  order by deptno
16  /
no rows selected
нет выбранных строк

```

Это именно то, к чему мы стремились — каждая строка в таблице EMP хранится в одном блоке с соответствующей ей строкой DEPT. Но что произойдет, если наша оценка оказалась ошибочной и 1024 не хватит? Что если данные некоторых отделов займут близкое к 1024 значение, а другие превысят его? Тогда, очевидно, данные не смогут уместиться в один блок, и некоторые записи EMP придется поместить в отдельный от соответствующей записи DEPT блок. Это легко увидеть, изменив предыдущий пример (я начну с таблиц в состоянии перед загрузкой, сразу после их создания). На этот раз мы загрузим каждую запись EMP восемь раз, чтобы увеличить количество записей о сотрудниках в каждом отделе:

```

EODA@ORA12CR1> insert into dept
2   ( deptno, dname, loc )
3   select deptno+r, dname, loc
4   from scott.dept,
5        (select level r from dual connect by level < 10);
36 rows created.
36 строк создано.

EODA@ORA12CR1> insert into emp
2   ( empno, ename, job, mgr, hiredate, sal, comm, deptno)
3   select rownum, ename, job, mgr, hiredate, sal, comm, deptno+r
4   from scott.emp,
5        (select level r from dual connect by level < 10),
6        (select level r2 from dual connect by level < 8);
882 rows created.
882 строк создано.

EODA@ORA12CR1> select min(count(*)), max(count(*)), avg(count(*))
2   from dept
3   group by dbms_rowid.rowid_block_number(rowid)
4   /

```

| MIN(COUNT(*)) | MAX(COUNT(*)) | AVG(COUNT(*)) |
|---------------|---------------|---------------|
| 1 | 7 | 6 |

Пока что все выглядит в точности как в предыдущем примере, но давайте сравним блоки, в которых находятся записи EMP, с блоками, где размещены записи DEPT:

```

EODA@ORA12CR1> select *
2   from (
3   select dept_blk, emp_blk,
4           case when dept_blk <> emp_blk then '*' end flag,
5           deptno
6   from (
7   select dbms_rowid.rowid_block_number(dept.rowid) dept_blk,
8          dbms_rowid.rowid_block_number(emp.rowid) emp_blk,
9          dept.deptno
10  from emp, dept
11  where emp.deptno = dept.deptno
12  )
13  )
14  where flag = '*'
15  order by deptno
16  /

```

| DEPT_BLK | EMP_BLK | F | DEPTNO |
|----------|---------|---|--------|
| 24845 | 22362 | * | 12 |
| 24845 | 22362 | * | 12 |
| 24845 | 22362 | * | 12 |
| ... | | | |
| 24844 | 22362 | * | 39 |
| 24844 | 22362 | * | 39 |
| 24844 | 22362 | * | 39 |

46 rows selected.
46 строк выбрано.

Можно заметить, что 46 из 882 строк таблицы EMP расположены в блоке, отличающемся от блока, в котором находятся строки таблицы DEPT, соответствующие им по значению DEPTNO. С учетом того, что размер кластера занижен (параметр SIZE оказался слишком мал для реальных данных), мы могли бы пересоздать кластер с SIZE, равным 1200, и тогда получить следующие результаты:

```

EODA@ORA12CR1> select min(count(*)), max(count(*)), avg(count(*))
2   from dept
3   group by dbms_rowid.rowid_block_number(rowid)
4   /

```

| MIN(COUNT(*)) | MAX(COUNT(*)) | AVG(COUNT(*)) |
|---------------|---------------|---------------|
| 6 | 6 | 6 |

```

EODA@ORA12CR1> select *
2   from (
3   select dept_blk, emp_blk,
4           case when dept_blk <> emp_blk then '*' end flag,
5           deptno
6   from (
7   select dbms_rowid.rowid_block_number(dept.rowid) dept_blk,

```

```

8      dbms_rowid.rowid_block_number(emp.rowid) emp_blk,
9      dept.deptno
10     from emp, dept
11     where emp.deptno = dept.deptno
12     )
13     )
14     where flag = '*'
15     order by deptno
16 /

no rows selected
нет выбранных строк

```

Теперь мы храним только шесть значений DEPTNO на блок, оставляя достаточно места для сохранения всех данных EMP в том же блоке, что и соответствующие им записи DEPT.

Есть небольшая головоломка, которой вы можете удивить и поразить своих друзей. Многие люди ошибочно полагают, что идентификатор строки (rowid) уникальным образом идентифицирует строку в базе данных, и по заданному идентификатору строки можно определить, из какой таблицы поступила та или иная строка. В действительности это *не удастся*. Вы можете и будете получать дублированные идентификаторы строк в кластере. Например, после выполнения показанного выше кода вы должны обнаружить такое:

```

EODA@ORA12CR1> select rowid from emp
2 intersect
3 select rowid from dept;

ROWID
-----
AAAE+/AAEAAABErAAA
AAAE+/AAEAAABErAAB
...
AAAE+/AAGAAAFdvAAE
AAAE+/AAGAAAFdvAAF

36 rows selected.
36 строк выбрано.

```

Каждый идентификатор строки, назначенный строкам в таблице DEPT, был также назначен строкам в таблице EMP. Причина в том, что для уникальной идентификации строки необходим идентификатор таблицы и идентификатор строки. Псевдостолбец ROWID является уникальным только в пределах таблицы.

Кроме того, согласно моим наблюдениям, многие считают, что объект кластера является объектом, известным лишь посвященным, и его на самом деле никто не использует — все имеют дело с обычными таблицами. В действительности мы применяем кластеры каждый раз, когда работаем в Oracle. Скажем, большая часть словаря данных хранится в разнообразных кластерах; для примера выполните следующий код от имени SYS:

```

SYS@ORA12CR1> break on cluster_name
SYS@ORA12CR1> select cluster_name, table_name
2     from user_tables
3     where cluster_name is not null
4     order by 1;

```

| CLUSTER_NAME | TABLE_NAME |
|----------------------|---------------|
| C_COBJ# | CDEF\$ |
| | CCOL\$ |
| C_FILE#_BLOCK# | SEG\$ |
| | UET\$ |
| C_MLOG# | SLOG\$ |
| | MLOG\$ |
| C_OBJ# | LIBRARY\$ |
| | ASSEMBLY\$ |
| | ATTRCOL\$ |
| | TYPE_MISC\$ |
| | VIEWTRCOL\$ |
| | OPQTYPE\$ |
| | ICOL\$ |
| | IND\$ |
| | CLU\$ |
| | TAB\$ |
| | COL\$ |
| | LOB\$ |
| | COLTYPE\$ |
| | SUBCOLTYPE\$ |
| | NTAB\$ |
| | REFCON\$ |
| | ICOLDEP\$ |
| C_OBJ#_INTCOL# | HISTGRM\$ |
| C_RG# | RGROUP\$ |
| | RGCHILD\$ |
| C_TOID_VERSION# | RESULT\$ |
| | PARAMETER\$ |
| | METHOD\$ |
| | ATTRIBUTE\$ |
| | COLLECTION\$ |
| | TYPE\$ |
| C_TS# | TS\$ |
| | FET\$ |
| C_USER# | TSQ\$ |
| | USER\$ |
| SMON_SCN_TO_TIME_AUX | SMON_SCN_TIME |
| 37 rows selected. | |
| 37 строк выбрано. | |

Как видите, большинство связанных с объектами данных хранится в единственном кластере (C_OBJ#) — в 17 таблицах, совместно использующих один и тот же блок. В них содержится преимущественно информация, связанная со столбцами, так что вся информация о наборе столбцов любой таблицы или индекса хранится физически в одном и том же блоке. В этом есть смысл, поскольку во время разбора запроса Oracle необходим доступ к данным для всех столбцов в таблице, на которую ссылается запрос. Если бы эти данные были разбросаны по разным местам, требовалось бы определенное время на то, чтобы собрать их вместе. Но данные обычно находятся в единственном блоке и в любой момент доступны.

Когда необходимо применять кластер? Пожалуй, проще перечислить ситуации, в которых он *не* должен использоваться.

- *Если вы ожидаете, что таблицы в кластере будут интенсивно модифицироваться.* Вы должны знать, что наличие индексного кластера будет негативно отражаться на производительности операторов DML, в частности, INSERT. Индексный кластер требует большего объема работы по управлению данными в кластере. Данные должны помещаться более аккуратно, на что уходит больше времени.
- *Если необходимо выполнять полные сканирования таблиц в кластере.* Вместо полного сканирования данных в одной таблице придется полностью сканировать данные (возможно) множества таблиц. Просмотру будут подлежать гораздо больше данных, поэтому полное сканирование будет выполняться дольше.
- *Если необходимо секционировать таблицы.* Таблицы в кластере не могут быть секционированы, как не может быть секционирован сам кластер.
- *Если вы уверены, что часто будете применять оператор TRUNCATE и загружать данные в таблицу.* Таблицы в кластере не могут быть усечены. Это очевидно — поскольку кластер хранит более одной таблицы в блоке, строки из кластерной таблицы должны удаляться посредством оператора DELETE.

Итак, кластер является подходящим вариантом, если есть данные, которые по большей части читаются (это вовсе *не* означает, что они никогда не перезаписываются; кластерные таблицы вполне допускают модификацию), читаются через индексы — либо индекс на кластерном ключе, либо другие индексы, которые вы предусмотрите на таблицах в кластере — и часто соединяются вместе. Ищите таблицы, которые логически связаны между собой и всегда используются вместе, как поступили разработчики словаря данных Oracle, когда поместили в кластер всю информацию, связанную со столбцами.

Заключительные соображения по поводу кластеризованных индекс-таблиц

Кластеризованные индекс-таблицы предоставляют возможность физически предварительно соединять данные. Кластеры применяются для хранения связанных между собой данных из многих таблиц в одном и том же блоке базы данных. Кластеры могут ускорять выполнение интенсивных в плане чтения операций, которые всегда соединяют данные вместе или получают доступ к связанным наборам данных (например, к записям сотрудников из отдела 10).

Кластеризованные индекс-таблицы сокращают количество блоков, которые Oracle приходится кешировать. Вместо того чтобы выделять десять блоков для десяти сотрудников из одного отдела, Oracle поместит их в один блок и тем самым увеличит эффективность кеша буферов. Что касается недостатков: если значение для параметра SIZE подсчитано некорректно, то кластеры могут оказаться неэффективными в отношении утилизации пространства и замедлять выполнение операций DML.

Кластеризованные хеш-таблицы

Кластеризованные хеш-таблицы концептуально очень похожи на описанные выше кластеризованные индекс-таблицы за одним главным исключением: индекс по кластерному ключу заменяется хеш-функцией. Данные в такой таблице являются индексом; какого-либо физического индекса не существует. База данных Oracle будет брать значение ключа для строки, хешировать его с помощью либо внутренней, либо предоставленной вами функции, и затем использовать его для выяснения, где данные должны находиться на диске. Однако одним из недостатков применения алгоритма хеширования для определения местонахождения данных является то, что выполнить просмотр таблицы в хеш-кластере по диапазону, не добавив в нее обычный индекс, не удастся. В индексном кластере следующий запрос смог бы воспользоваться индексом по кластерному ключу для поиска этих строк:

```
select * from emp where deptno between 10 and 20
```

В хеш-кластере этот запрос в результате приводит к полному сканированию таблицы, если только нет индекса на столбце DEPTNO. Лишь поиск с условием точного равенства (включая IN со списками и подзапросами) может выполняться по хеш-ключу без участия индекса, который поддерживает просмотры по диапазону.

В идеальном мире с эффективно распределенными значениями хеш-ключа и хеш-функцией, которая распределяет их равномерно по всем выделенным под хеш-кластер блокам, мы могли бы запрашивать и получать данные с помощью единственной операции ввода-вывода. Однако в реальных условиях обычно получается больше значений ключа хеш-кластера, указывающих на один и тот же блок в базе данных, чем физически может уместиться в этом блоке. В результате Oracle будет сцеплять блоки вместе в связный список, чтобы сохранить все строки, хеш-ключ для которых указывает на этот блок. В таком случае, когда нужно извлечь строки, соответствующие одному и тому же хеш-ключу, возможно понадобится посетить более одного блока.

Подобно хеш-таблице в языке программирования, хеш-таблицы в базе данных имеют фиксированный размер. При создании таблицы вы должны раз и навсегда определить, какое количество хеш-ключей будет иметь таблица. Это не ограничивает количество строк, которые можно в нее поместить.

На рис. 10.9 показано графическое представление хеш-кластера с созданной в нем таблицей EMP. Когда клиент отправляет запрос с ключом хеш-кластера в предикате, Oracle применяет хеш-функцию, чтобы определить, в каком блоке находятся запрашиваемые данные, и затем считывает один этот блок с целью нахождения нужных данных. Если возникает слишком много коллизий или значение параметра SIZE в CREATE CLUSTER было занижено, Oracle выделит блоки переполнения, которые будут сцеплены с исходным блоком.

Для создания хеш-кластера используется тот же самый оператор CREATE CLUSTER, который применялся для создания индексного кластера, только с другими опциями. К нему добавляется опция HASHKEYS, в которой указывается размер хеш-таблицы. База данных Oracle возьмет ваше значение HASHKEYS и округлит его в большую сторону до ближайшего простого числа; количество ключей хеш-кластера всегда должно быть простым числом.

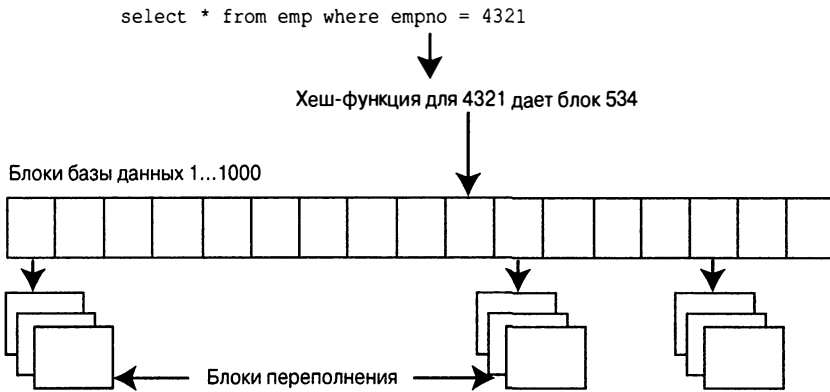


Рис. 10.9. Графическая иллюстрация хеш-кластера

Затем Oracle умножит это измененное значение HASHKEYS на значение параметра SIZE и выделит для кластера столько байтов пространства, сколько получится в результате такого умножения. В этом заключается крупное отличие хеш-кластера от рассмотренного ранее индексного кластера, для которого пространство выделяется динамически по мере надобности. Для хеш-кластера пространство выделяется заранее в объеме, достаточном для хранения $\text{HASHKEYS}/\text{trunc}(\text{blocksize}/\text{SIZE})$ байтов данных. Например, если параметр SIZE установлен в 1500 байтов, а размер блоков составляет 4 Кбайт, то Oracle будет ожидать сохранения двух ключей на блок. То есть, если планируется наличие 1000 хеш-ключей, то будет выделено 500 блоков.

Интересно отметить, что в отличие от обычной хеш-таблицы в языке программирования хеш-таблицы в Oracle допускают возникновение хеш-коллизий — в действительности во многих случаях такая ситуация даже желательна. В приведенном ранее примере с таблицами DEPT и EMP можно создать хеш-кластер на основе столбца DEPTNO. Очевидно, что многие строки дадут в результате хеширования одно и то же значение, и как раз это от них и ожидается (они имеют то же самое значение DEPTNO). В некотором отношении именно это и является задачей кластера: объединение вместе похожих данных. По этой причине Oracle предлагает указать значение HASHKEYS (сколько номеров отделов ожидается в будущем) и SIZE (размер данных, которые будут ассоциированы с каждым номером отдела). Затем Oracle создаст хеш-таблицу, рассчитанную на HASHKEYS отделов размером по SIZE байтов каждый. Чего желательно избегать — так это незапланированных хеш-коллизий. Должно быть ясно, что если установить размер хеш-таблицы в 1000 (на самом деле он будет равен 1009, т.к. размер хеш-таблицы должен быть простым числом, и Oracle произведет округление) и поместить в эту таблицу 1010 отделов, то возникнет минимум одна коллизия (два разных отдела будут давать одинаковые хеш-значения). Незапланированных хеш-коллизий следует избегать, потому что они приносят накладные расходы и увеличивают вероятность сцепления блоков.

Чтобы посмотреть, как задействуется пространство хеш-кластерами, мы применим небольшую хранимую процедуру SHOW_SPACE (за деталями обращайтесь в раздел “Настройка среды” в самом начале книги). Эта процедура использует пакет DBMS_SPACE для получения сведений о пространстве, занимаемом сегментами в базе данных.

Выполнив показанный ниже оператор CREATE CLUSTER, мы сможем увидеть, какое пространство он выделит:

```

EODA@ORA12CR1> create cluster hash_cluster
2 ( hash_key number )
3 hashkeys 1000
4 size 8192
5 tablespace mssm
6 /

```

Cluster created.

Кластер создан.

```

EODA@ORA12CR1> exec show_space( 'HASH_CLUSTER', user, 'CLUSTER' )
Free Blocks..... 0
Total Blocks..... 1,024
Total Bytes..... 8,388,608
Total MBytes..... 8
Unused Blocks..... 14
Unused Bytes..... 114,688
Last Used Ext FileId..... 7
Last Used Ext BlockId..... 1,024
Last Used Block..... 114

```

PL/SQL procedure successfully completed.

Процедура PL/SQL успешно завершена.

Несложно заметить, что под таблицу было выделено всего 1024 блока. Четырнадцать из этих блоков не используются (свободны). Один блок относится к накладным расходам таблицы и предназначен для управления экстенентами. Следовательно, перед HWM-маркером этого объекта находится 1009 блоков, и именно они применяются кластером. Поскольку получилось так, что 1009 является следующим простым числом, которое больше 1000, а размер блоков составляет 8 Кбайт, мы видим, что на самом деле Oracle было выделено и отформатировано 1009 блоков. Результирующее число немного больше, чем это, из-за способа округления экстенентов и/или за счет использования локально управляемых табличных пространств с экстенентами одинаковых размеров.

Приведенный пример указывает на связанные с хеш-кластерами сложности, о которых обязательно следует знать. Обычно при создании пустой таблицы количество блоков, находящихся перед HWM-маркером этой таблицы, равно нулю. При полном сканировании Oracle доходит до HWM-маркера и останавливается. В случае хеш-кластера создаваемые таблицы с самого начала будут иметь большой размер, и их создание будет занимать больше времени, т.к. Oracle придется инициализировать каждый блок — действие, которое обычно выполняется при добавлении данных в таблицу. Потенциально таблицы могут содержать какие-то данные в первом и последнем блоках и ничего между ними. Полное сканирование практически пустого хеш-кластера будет занимать столько же времени, сколько и полное сканирование хеш-кластера, заполненного данными. Это необязательно является недостатком, т.к. хеш-кластер создается для очень быстрого доступа к данным за счет поиска по хеш-ключу, а не для частого выполнения полного сканирования.

Теперь можно приступать к добавлению таблиц к хеш-кластеру в той же манере, как это делалось для индексного кластера:

```

EODA@ORA12CR1> create table hashed_table
  2 ( x number, data1 varchar2(4000), data2 varchar2(4000) )
  3 cluster hash_cluster(x);
Table created.
Таблица создана.

```

Чтобы посмотреть, в чем особенность хеш-кластера, проведем небольшой тест. Я создал хеш-кластер, загрузил в него кое-какие данные, скопировал эти данные в обычную таблицу с обычным индексом и затем выполнил произвольные операции чтения в каждой таблице (одинаковые “произвольные” операции чтения в каждой). С использованием runstats, SQL_TRACE и TKPROF я смог определить характеристики каждой из них. Ниже приведен код, который применялся при этом:

```

EODA@ORA12CR1> create cluster hash_cluster
  2 ( hash_key number )
  3 hashkeys 75000
  4 size 150
  5 /
Cluster created.
Кластер создан.

```

```

EODA@ORA12CR1> create table t_hashed
  2 cluster hash_cluster(object_id)
  3 as
  4 select *
  5 from all_objects
  6 /
Table created.
Таблица создана.

```

```

EODA@ORA12CR1> alter table t_hashed add constraint
  2 t_hashed_pk primary key(object_id)
  3 /
Table altered.
Таблица изменена.

```

```

EODA@ORA12CR1> begin
  2 dbms_stats.gather_table_stats( user, 'T_HASHED' );
  3 end;
  4 /
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.

```

Хеш-кластер создан с параметром SIZE, равным 150 байтов. Такое значение объясняется тем, что размер строк в таблице в среднем будет составлять около 100 байтов, но в зависимости от данных иногда будет чуть большим или чуть меньшим, но не превышать 150 байтов. Затем в кластере была создана таблица и заполнена данными из таблицы ALL_OBJECTS.

Далее я создал обычный клон этой таблицы:

```

EODA@ORA12CR1> create table t_heap
  2 as
  3 select *
  4 from t_hashed
  5 /
Table created.

```



```

EODA@ORA12CR1> alter table t_heap add constraint
  2  t_heap_pk primary key(object_id)
  3  /

```

Table altered.

Таблица изменена.

```

EODA@ORA12CR1> begin
  2  dbms_stats.gather_table_stats( user, 'T_HEAP' );
  3  end;
  4  /

```

PL/SQL procedure successfully completed.

Процедура PL/SQL успешно завершена.

Затем понадобились случайные данные, чтобы с их помощью выбрать строки из каждой таблицы. Для этого я просто выбрал все значения OBJECT_ID в массив и отсортировал их в случайном порядке, чтобы разбросать по всей таблице. Для определения и объявления массива использовался пакет PL/SQL, а с помощью небольшого PL/SQL-кода массив был заполнен простыми числами:

```

EODA@ORA12CR1> create or replace package state_pkg
  2  as
  3      type array is table of t_hashed.object_id%type;
  4      g_data array;
  5  end;
  6  /

```

Package created.

Пакет создан.

```

EODA@ORA12CR1> begin
  2  select object_id bulk collect into state_pkg.g_data
  3  from t_hashed
  4  order by dbms_random.random;
  5  end;
  6  /

```

PL/SQL procedure successfully completed.

Процедура PL/SQL успешно завершена.

Чтобы посмотреть, как работает каждая из таблиц, я подготовил следующий блок кода (замена всех вхождений HEAP словом HASHED даст блок кода для второй таблицы):

```

EODA@ORA12CR1> declare
  2  l_rec t_heap%rowtype;
  3  begin
  4  for i in 1 .. state_pkg.g_data.count
  5  loop
  6      select * into l_rec from t_heap
  7      where object_id = state_pkg.g_data(i);
  8  end loop;
  9  end;
 10  /

```

PL/SQL procedure successfully completed.

Процедура PL/SQL успешно завершена.

Далее я запустил этот блок кода три раза (а также его вторую версию, в которой вместо HEAP указано HASHED). Первый запуск был предназначен для разогрева системы и проведения всех необходимых полных разборов. Во время второго запуска

блоков кода я воспользовался утилитой runstats, чтобы посмотреть, чем они реально отличаются: сначала для таблицы T_HASHED, а затем — для T_HEAP. В третий раз блоки кода выполнялись при включенной опции SQL_TRACE, что позволило просмотреть отчет TKPROF. Отчет утилиты runstats выглядел так, как показано ниже:

```
EODA@ORA12CR1> exec runstats_pkg.rs_stop(10000);
```

```
Run1 ran in 198 cpu hsecs
Run2 ran in 206 cpu hsecs
run 1 ran in 96.12% of the time
```

| Name | Run1 | Run2 | Diff |
|--|-------------|---------------|---------------|
| STAT...redo size | 21,896 | 23,716 | 1,820 |
| STAT...table scan rows gotten | 0 | 4,611 | 4,611 |
| LATCH.simulator hash latch | 4,326 | 9,114 | 4,788 |
| LATCH.cache buffers chains | 145,070 | 217,054 | 71,984 |
| STAT...Cached Commit SCN refer | 72,056 | 0 | -72,056 |
| STAT...consistent gets pin | 72,119 | 39 | -72,080 |
| STAT...consistent gets pin (fa | 72,119 | 39 | -72,080 |
| STAT...no work - consistent re | 72,105 | 24 | -72,081 |
| STAT...cluster key scans | 72,105 | 1 | -72,104 |
| STAT...cluster key scan block | 72,105 | 1 | -72,104 |
| STAT...rows fetched via callba | 18 | 72,123 | 72,105 |
| STAT...table fetch by rowid | 18 | 72,123 | 72,105 |
| STAT...index fetch by key | 19 | 72,126 | 72,107 |
| STAT...buffer is not pinned co | 72,141 | 216,354 | 144,213 |
| STAT...session logical reads | 72,320 | 216,554 | 144,234 |
| STAT...consistent gets | 72,175 | 216,419 | 144,244 |
| STAT...consistent gets from ca | 72,175 | 216,419 | 144,244 |
| STAT...consistent gets examina | 56 | 216,380 | 216,324 |
| STAT...consistent gets examina | 56 | 216,380 | 216,324 |
| STAT...session pga memory | 262,144 | -65,536 | -327,680 |
| STAT...logical read bytes from | 592,445,440 | 1,774,010,368 | 1,181,564,928 |
| Run1 latches total versus runs -- difference and pct | | | |
| Run1 | Run2 | Diff | Pct |
| 223,979 | 299,841 | 75,862 | 74.70% |

PL/SQL procedure successfully completed.

Процедура PL/SQL успешно завершена.

Здесь видно, что эти два блока кода были выполнены за почти одинаковое время по часам процессора. Однако важным отличием, на которое следует обратить внимание, является значительное снижение количества защелок цепочек кеша буферов. В первой реализации (T_HASHED) защелок было значительно меньше, а это означает, что хешированная реализация должна лучше масштабироваться в среде с интенсивным чтением, т.к. она использует меньше ресурсов, требующих определенного уровня сериализации. Это произошло исключительно потому, что хешированной реализации необходимо гораздо меньше операций ввода-вывода, чем традиционной таблице (T_HEAP) — в приведенном отчете это подтверждается количеством согласованных чтений (consistent gets). Отчет TKPROF демонстрирует это даже яснее:

```
SELECT * FROM T_HASHED WHERE OBJECT_ID = :B1
```

| call | count | cpu | elapsed | disk | query | current | rows |
|---------|--------|------|---------|------|-------|---------|-------|
| Parse | 1 | 0.00 | 0.00 | 0 | 0 | 0 | 0 |
| Execute | 72105 | 0.75 | 0.75 | 0 | 2 | 0 | 0 |
| Fetch | 72105 | 0.74 | 0.71 | 0 | 72105 | 0 | 72105 |
| total | 144211 | 1.50 | 1.47 | 0 | 72107 | 0 | 72105 |

...
Rows (1st) Rows (avg) Rows (max) Row Source Operation

1 1 1 TABLE ACCESS HASH T_HASHED (cr=1 pr=0
⚡pw=0 time=19 us)

SELECT * FROM T_HEAP WHERE OBJECT_ID = :B1

| call | count | cpu | elapsed | disk | query | current | rows |
|---------|--------|------|---------|------|--------|---------|-------|
| Parse | 1 | 0.00 | 0.00 | 0 | 0 | 0 | 0 |
| Execute | 72105 | 0.81 | 0.81 | 0 | 0 | 0 | 0 |
| Fetch | 72105 | 0.75 | 0.74 | 0 | 216315 | 0 | 72105 |
| total | 144211 | 1.56 | 1.55 | 0 | 216315 | 0 | 72105 |

...
Rows (1st) Rows (avg) Rows (max) Row Source Operation

1 1 1 TABLE ACCESS BY INDEX ROWID T_HEAP
⚡(cr=3 pr=0 pw=0 ...
1 1 1 INDEX UNIQUE SCAN T_HEAP_PK (cr=2 pr=0
⚡pw=0 time=14...

Реализация T_HASHED просто преобразовывала переданное в запросе значение OBJECT_ID в файл или блок для чтения и считывала его — индекс не применялся. Однако таблица T_HEAP должна была выполнять для каждой строки по две операции ввода-вывода в индексе. Значение cr=2 в столбце Row Source Operation отчета TKPROF показывает, сколько точно согласованных чтений было выполнено на индексе. Каждый раз, когда запускался запрос с условием OBJECT_ID = :B1, Oracle приходилось читать корневой блок индекса и затем находить листовой блок, содержащий адрес этой строки. Затем необходимо было брать информацию из листового блока, которая включала идентификатор строки, и получать доступ к этой строке в таблице, выполняя третью операцию ввода-вывода. В таблице T_HEAP делалось в три раза больше операций ввода-вывода, чем в таблице T_HASHED.

Интересно отметить следующие моменты.

- Хеш-кластер выполнил значительно меньше операций ввода-вывода (столбец query отчета). Именно этого мы ожидали. Запрос просто брал произвольные значения OBJECT_ID, применял к ним хеш-функцию и переходил в блок. Для извлечения данных хеш-кластер должен был делать, по крайней мере, одну операцию ввода-вывода. Обычной таблице с индексом для этого приходилось сканировать индекс и получать доступ к таблице по идентификатору строки. Чтобы получить данные, индексированная таблица в этом случае должна выполнить, по меньшей мере, три операции ввода-вывода.

- Запрос хеш-кластера использовал одинаковое количество ресурсов процессора для всех целей, даже несмотря на то, что он обращался к кешу буферов в три раза меньше. Это также можно было ожидать. Операция хеширования интенсивно использует процессор. Операция индексного поиска насыщена операциями ввода-вывода. Был достигнут компромисс. Однако при увеличении количества пользователей можно ожидать, что запрос хеш-кластера будет лучше масштабироваться, т.к. получать доступ к кешу буферов ему придется гораздо реже.

Последний момент очень важен. В мире компьютеров все вращается вокруг ресурсов и их утилизации. Если вы ограничены в плане ввода-вывода и запускаете запросы, которые предполагают выполнение множества операций чтения по ключу, то хеш-кластер может улучшить производительность. Если вы ограничены возможностями процессора, то хеш-кластер может привести к снижению производительности, потому что ему требуется больше ресурсов ЦП для выполнения операции хеширования. Но если лишние ресурсы процессора тратятся на защелки цепочек кеша буферов, тогда наоборот — хеш-кластер может существенно сократить количество используемых ресурсов. Это одна из основных причин неработоспособности эмпирических правил в реальных системах: то, что работает у вас, может не работать у других, находящихся в похожих, однако отличающихся условиях.

Существует особый вид хеш-кластера, который называется *однотабличным хеш-кластером*. Это оптимизированная версия общего хеш-кластера, который был только что рассмотрен. Он поддерживает в кластере только одну таблицу одновременно (перед созданием новой таблицы существующая таблица в однотабличном хеш-кластере должна быть удалена с помощью DROP). Вдобавок, если есть однозначное соответствие между хеш-ключами и строками данных, то доступ к строкам также будет быстрее. Такие хеш-кластеры предназначены для ситуаций, когда необходимо получать доступ к таблице по первичному ключу, не объединяя ее в кластер с другими таблицами. Когда требуется быстрый доступ к записям сотрудников по EMPNO, однотабличный хеш-кластер может оказаться как раз тем, что нужно. Я провел описанный выше тест для однотабличного хеш-кластера и обнаружил, что производительность стала даже лучше, чем в случае применения обычного хеш-кластера. В этом примере вы можете пойти еще дальше и воспользоваться тем, что Oracle позволяет создавать собственную специализированную хеш-функцию (вместо применения функции по умолчанию). В этом случае можно использовать только столбцы, которые доступны в таблице, и только встроенные функции Oracle (т.е. никакого кода PL/SQL). Получая преимущество от того факта, что значение OBJECT_ID в предыдущем примере представляло собой число от 1 до 75 000, я построил собственную хеш-функцию так, что она возвращает в качестве значения сам столбец OBJECT_ID. В этом случае хеш-коллизия никогда не возникнет. Чтобы продемонстрировать все сказанное, я создам однотабличный хеш-кластер с собственной хеш-функцией:

```
EODA@ORA12CR1> create cluster hash_cluster
2 ( hash_key number(10) )
3 hashkeys 75000
4 size 150
5 single table
6 hash is HASH_KEY
7 /
Cluster created.
Кластер создан.
```

Чтобы сделать однотабличный хеш-кластер, я просто добавил к оператору ключевые слова `SINGLE TABLE`. В конструкции `HASH IS` в данном случае указан кластерный ключ `HASH KEY`. Это SQL-функция, так что при желании можно было бы использовать выражение вроде `trunc(mod(hash_key/324+278555)/abs(hash_key+1))` (приведенное выражение не следует считать хорошей хеш-функцией — оно всего лишь демонстрирует возможность создания функций произвольной сложности). Вместо обычного числа применяется `NUMBER(10)`. Из-за того, что хеш-значение должно быть целочисленным, оно не может иметь дробные компоненты. Теперь в этом кластере можно создать таблицу, чтобы построить хеш-таблицу:

```

EODA@ORA12CR1> create table t_hashed
2 cluster hash_cluster(object_id)
3 as
4 select OWNER, OBJECT_NAME, SUBOBJECT_NAME,
5        cast( OBJECT_ID as number(10) ) object_id,
6        DATA_OBJECT_ID, OBJECT_TYPE, CREATED,
7        LAST_DDL_TIME, TIMESTAMP, STATUS, TEMPORARY,
8        GENERATED, SECONDARY
9 from all_objects
10 /
Table created.

```

Обратите внимание, что я воспользовался встроенной функцией `CAST`, чтобы привести тип данных столбца `OBJECT_ID` к тому, каким он должен быть. Затем было проведено тестирование, как и раньше (три запуска каждого блока кода), и на этот раз вывод `runstats` выглядел гораздо более впечатляющим:

```

Run1 ran in 183 cpu hsecs
Run2 ran in 195 cpu hsecs
run 1 ran in 93.85% of the time

```

| Name | Run1 | Run2 | Diff |
|--------------------------------|-------------|---------------|---------------|
| STAT...Cached Commit SCN refer | 42,970 | 0 | -42,970 |
| LATCH.cache buffers chains | 165,638 | 216,945 | 51,307 |
| STAT...cluster key scans | 72,105 | 1 | -72,104 |
| STAT...table fetch by rowid | 13 | 72,118 | 72,105 |
| STAT...rows fetched via callba | 13 | 72,118 | 72,105 |
| STAT...index fetch by key | 14 | 72,121 | 72,107 |
| STAT...consistent gets pin (fa | 82,562 | 39 | -82,523 |
| STAT...consistent gets pin | 82,562 | 39 | -82,523 |
| STAT...cluster key scan block | 82,548 | 1 | -82,547 |
| STAT...buffer is not pinned co | 82,574 | 216,344 | 133,770 |
| STAT...session logical reads | 82,732 | 216,516 | 133,784 |
| STAT...consistent gets | 82,603 | 216,404 | 133,801 |
| STAT...consistent gets from ca | 82,603 | 216,404 | 133,801 |
| STAT...session pga memory | 0 | 196,608 | 196,608 |
| STAT...consistent gets examina | 41 | 216,365 | 216,324 |
| STAT...consistent gets examina | 41 | 216,365 | 216,324 |
| STAT...logical read bytes from | 677,740,544 | 1,773,699,072 | 1,095,958,528 |

```
Run1 latches total versus runs -- difference and pct
```

| Run1 | Run2 | Diff | Pct |
|---------|---------|--------|--------|
| 244,074 | 299,493 | 55,419 | 81.50% |

```
PL/SQL procedure successfully completed.
```

Процедура PL/SQL успешно завершена.

Этот однотабличный хеш-кластер требует даже меньше зашелок в кеше буферов для обработки (он быстрее заканчивает поиск данных и у него имеется больше информации). В результате показатели использования ЦП в отчете TKPROF на этот раз тоже оказались значительно ниже:

```
SELECT * FROM T_HASHED WHERE OBJECT_ID = :B1
```

| call | count | cpu | elapsed | disk | query | current | rows |
|---------|--------|------|---------|------|-------|---------|-------|
| Parse | 1 | 0.00 | 0.00 | 0 | 0 | 0 | 0 |
| Execute | 72105 | 0.70 | 0.70 | 0 | 2 | 0 | 0 |
| Fetch | 72105 | 0.63 | 0.64 | 0 | 82548 | 0 | 72105 |
| total | 144211 | 1.33 | 1.35 | 0 | 82550 | 0 | 72105 |

```
...
Rows (1st) Rows (avg) Rows (max) Row Source Operation
```

```
-----
1 1 1 TABLE ACCESS HASH T_HASHED (cr=1 pr=0)
⌚pw=0 time=25 us)
```

```
*****
```

```
SELECT * FROM T_HEAP WHERE OBJECT_ID = :B1
```

| call | count | cpu | elapsed | disk | query | current | rows |
|---------|--------|------|---------|------|--------|---------|-------|
| Parse | 1 | 0.00 | 0.00 | 0 | 0 | 0 | 0 |
| Execute | 72105 | 0.87 | 0.84 | 0 | 0 | 0 | 0 |
| Fetch | 72105 | 0.70 | 0.71 | 0 | 216315 | 0 | 72105 |
| total | 144211 | 1.58 | 1.55 | 0 | 216315 | 0 | 72105 |

```
...
Rows (1st) Rows (avg) Rows (max) Row Source Operation
```

```
-----
1 1 1 TABLE ACCESS BY INDEX ROWID T_HEAP (cr=3
```

```
⌚pr=0 pw=0 time=22...
```

```
1 1 1 INDEX UNIQUE SCAN T_HEAP_PK (cr=2 pr=0
```

```
⌚pw=0 ...
```

Заключительные соображения по поводу кластеризованных хеш-таблиц

Выше были раскрыты все основные аспекты хеш-кластера. Концептуально хеш-кластеры похожи на индексные кластеры за исключением того, что кластерный индекс в них не применяется. В этом случае *индексом являются сами данные*. Кластерный ключ хешируется в адрес блока, в котором должны находиться данные. Относительно хеш-кластеров важно понимать перечисленные ниже моменты.

- Пространство под хеш-кластер выделяется с самого начала. База данных Oracle будет брать значение `HASHKEYS/trunk(blocksize/SIZE)` и сразу же выделять и форматировать пространство соответствующего объема. Как только в этот кластер помещается первая таблица, любое полное сканирование будет попадать в каждый выделенный блок. В этом отношении хеш-кластер отличается от всех остальных таблиц.
- Количество ключей `HASHKEY` в хеш-кластере является фиксированным. Изменить его без перестройки кластера нельзя. Это значение никак не ограничивает объем данных, которые можно сохранять в кластере, а просто огра-

ничивает количество уникальных хеш-ключей, которое может генерироваться для кластера. Слишком низкое значение будет приводить к возникновению непредвиденных хеш-коллизий, что может повлиять на производительность.

- Производить поиск диапазонов строк по кластерному ключу нельзя. Предикаты, такие как `WHERE cluster_key BETWEEN 50 AND 60`, не могут использовать алгоритм хеширования. В диапазоне между 50 и 60 может находиться бесконечное количество значений, и серверу пришлось бы генерировать их все, чтобы применить к каждому из них хеш-функцию и посмотреть, есть ли там какие-нибудь данные. Понятно, что это невозможно. Кластер будет сканироваться полностью, если вы используете диапазон по кластерному ключу и не индексируете его с помощью обычного индекса.

Хеш-кластеры подходят в следующих ситуациях.

- Когда довольно точно известно, сколько строк будет иметь таблица за все время существования, или хотя бы какой-то обоснованный верхний предел. Правильно выбранные значения для параметров `HASHKEYS` и `SIZE` позволяют избежать повторного создания таблицы.
- Когда DML-операции, особенно вставки, являются облегченными по сравнению с операциями извлечения. Это значит, что имеется баланс между количеством извлекаемых данных и количеством добавляемых данных. У одного пользователя может быть 100 000 операций вставки за единицу времени, а у другого — 100; все зависит от того, какую схему они применяют для извлечения данных. Операции обновления не приносят значительные накладные расходы, если только вы не обновляете `HASHKEY`, что является не очень хорошей идеей, т.к. может привести к перемещению строк.
- Когда вы постоянно получаете доступ к данным по значению `HASHKEY`. Например, это может быть таблица частей, и доступ к этим частям осуществляется по номеру части. Справочные таблицы особенно хорошо подходят для использования в хеш-кластерах.

Отсортированные кластеризованные хеш-таблицы

Отсортированные хеш-кластеры доступны в Oracle 10g и последующих версиях. Они сочетают в себе качества только что описанных кластеризованных хеш-таблиц и индекс-таблиц. Их лучше всего применять в ситуации, когда данные постоянно извлекаются с помощью запроса типа:

```
Select *
  From t
 Where КЛЮЧ=:x
 Order by ОТСОРТИРОВАННЫЙ_СТОЛБЕЦ
```

То есть вы извлекаете данные по ключу и нуждаетесь в их упорядочении по какому-то другому столбцу. Когда используется отсортированный хеш-кластер, Oracle может возвращать данные, вообще не выполняя какой-либо сортировки. Это достигается хранением вставляемых данных в физически отсортированном виде по ключу. Предположим, что есть следующая таблица заказов клиентов:

```

EODA@ORA12CR1> select cust_id, order_dt, order_number
2 from cust_orders
3 order by cust_id, order_dt;

```

| | CUST_ID | ORDER_DT | ORDER_NUMBER |
|---|-----------|--------------------|--------------|
| 1 | 31-MAR-05 | 09.13.57.000000 PM | 21453 |
| | 11-APR-05 | 08.30.45.000000 AM | 21454 |
| | 28-APR-05 | 06.21.09.000000 AM | 21455 |
| 2 | 08-APR-05 | 03.42.45.000000 AM | 21456 |
| | 19-APR-05 | 08.59.33.000000 AM | 21457 |
| | 27-APR-05 | 06.35.34.000000 AM | 21458 |
| | 30-APR-05 | 01.47.34.000000 AM | 21459 |

7 rows selected.

7 строк выбрано.

Таблица хранится в отсортированном хеш-кластере, в котором хеш-ключом является CUST_ID, а полем, по которому должна выполняться сортировка — ORDER_DT. Графически ее можно изобразить так, как показано на рис. 10.10, где 1, 2, 3, 4, ... — это записи, которые хранятся в каждом блоке в отсортированном виде.

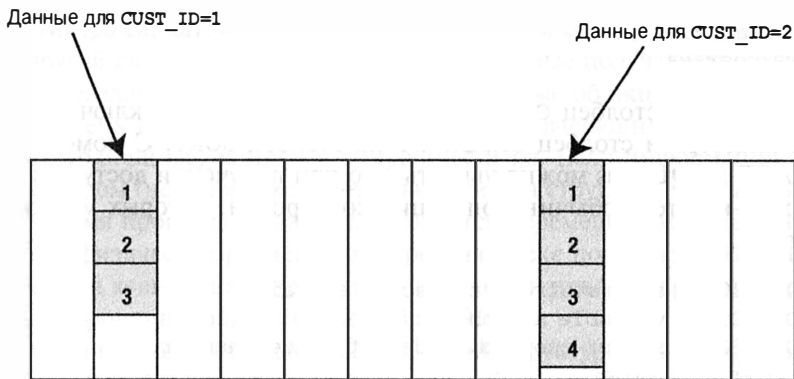


Рис. 10.10. Графическая иллюстрация отсортированного хеш-кластера

Отсортированный хеш-кластер создается почти так же, как и любой другой кластер. Чтобы создать отсортированный хеш-кластер для хранения показанных выше данных, можно было бы применить такой код:

```

EODA@ORA12CR1> CREATE CLUSTER shc
2 (
3   cust_id    NUMBER,
4   order_dt   timestamp SORT
5 )
6 HASHKEYS 10000
7 HASH IS cust_id
8 SIZE 8192
9 /

```

Cluster created.

Кластер создан.

Здесь появилось новое ключевое слово SORT. При создании кластера использовалась конструкция HASH IS CUST_ID, а также конструкция ORDER_DT с типом

timestamp и ключевым словом SORT. Это означает, что данные будут обнаруживаться по столбцу CUST_ID (where CUST_ID=:X) и физически извлекаться отсортированными по ORDER_DT. Формально это означает сохранение данных, которые будут извлекаться через столбец NUMBER и сортироваться по TIMESTAMP. Имена столбцов здесь значения не имеют, как это было в индексных (со структурой В-дерева) и обычных хеш-кластерах, но по соглашению столбцы именуются в соответствии с тем, что они представляют.

Оператор CREATE TABLE для таблицы CUST_ORDERS мог бы выглядеть так:

```

EODA@ORA12CR1> CREATE TABLE cust_orders
2  ( cust_id      number,
3    order_dt     timestamp SORT,
4    order_number number,
5    username     varchar2(30),
6    ship_addr    number,
7    bill_addr    number,
8    invoice_num  number
9  )
10 CLUSTER shc ( cust_id, order_dt )
11 /

```

Table created.
Таблица создана.

Мы отобразили столбец CUST_ID этой таблицы на хеш-ключ отсортированного хеш-кластера и столбец ORDER_DT на столбец SORT. С помощью средства AUTOTRACE в SQL*Plus можно увидеть, что при получении доступа к отсортированному хеш-кластеру обычные операции сортировки, которых мы ожидаем, отсутствуют:

```

EODA@ORA12CR1> set autotrace traceonly explain
EODA@ORA12CR1> variable x number
EODA@ORA12CR1> select cust_id, order_dt, order_number
2    from cust_orders
3   where cust_id = :x
4   order by order_dt;

```

| Id | Operation | Name | Rows | Bytes | Cost (%CPU) | Time |
|-----|-------------------|-------------|------|-------|-------------|----------|
| 0 | SELECT STATEMENT | | 1 | 39 | 1 (0) | 00:00:01 |
| * 1 | TABLE ACCESS HASH | CUST_ORDERS | 1 | 39 | 1 (0) | 00:00:01 |

```

EODA@ORA12CR1> select job, hiredate, empno
2    from scott.emp
3   where job = 'CLERK'
4   order by hiredate;

```

| Id | Operation | Name | Rows | Bytes | Cost (%CPU) | Time |
|-----|-------------------------------------|---------|------|-------|-------------|----------|
| 0 | SELECT STATEMENT | | 4 | 112 | 2 (0) | 00:00:01 |
| 1 | SORT ORDER BY | | 4 | 112 | 2 (0) | 00:00:01 |
| 2 | TABLE ACCESS BY INDEX ROWID BATCHED | EMP | 4 | 112 | 2 (0) | 00:00:01 |
| * 3 | INDEX RANGE SCAN | JOB_IDX | 4 | | 1 (0) | 00:00:01 |

```
EODA@ORA12CR1> set autotrace off
```

Я добавил запрос к обычной таблице SCOTT.EMP (после индексирования столбца JOB в целях демонстрации) для сравнения плана запроса SCOTT.EMP с тем, что может сделать отсортированный хеш-кластер, когда необходим доступ к данным в режиме FIFO (подобно очереди). Как видите, отсортированный хеш-кластер предпринимает один шаг: он берет CUST_ID=:X, хеширует входные данные, находит первую строку и просто начинает читать строки, т.к. они уже отсортированы в требуемом порядке. Обычная таблица работает совсем по-другому: она находит все строки JOB='CLERK' (которые в традиционной таблице могут располагаться где угодно), сортирует их и затем возвращает первую из них.

Таким образом, отсортированный хеш-кластер в отношении извлечения данных обладает всеми теми же качествами, что и обычный хеш-кластер, в том смысле, что может читать данные без обхода индекса, а также многими качествами, которые присущи индекс-таблице, в частности, данные в нем сортируются в пределах ключа по выбранному полю. Такая структура данных работает хорошо, когда данные поступают в порядке поля сортировки и по ключу. То есть, с течением времени данные поступают в порядке сортировки по возрастанию для любого заданного значения ключа. Примером могут служить данные об акциях. Каждую ночь вы получаете новый файл с тикерами, датой (дата может быть ключом сортировки, а тикер — хеш-ключом) и другой связанной информацией. Эти данные получаются и загружаются в порядке, определяемом ключом сортировки. Данные об акциях для тикера ORCL за вчерашний день не могут поступить завтра — за вчерашним значением всегда загружается сегодняшнее, за ним — завтрашнее и т.д. Если информация поступает случайным образом (не в порядке сортировки), то эта структура данных быстро разрушается во время процесса вставки, т.к. придется перемещать многие данные, чтобы обеспечить физическое размещение строк на диске по порядку. В таком случае отсортированный хеш-кластер применять не рекомендуется (с другой стороны, для таких данных хорошо подойдет индекс-таблица).

Рассматривая возможность использования этой структуры, вы должны задействовать те же самые соображения из раздела, посвященного обычным хеш-кластерам, в дополнение к ограничению относительно того, что с течением времени данные должны поступать отсортированными для каждого значения ключа.

Вложенные таблицы

Вложенные таблицы являются частью объектно-реляционных расширений Oracle. Вложенная таблица, относящаяся к одному из двух доступных в Oracle типов коллекций, очень похожа на дочернюю таблицу в традиционной паре родительская таблица/дочерняя таблица, применяемой в реляционной модели. Вложенная таблица — это неупорядоченный набор элементов данных одного и того же типа, которым может быть либо встроенный, либо объектный тип данных. Однако вдобавок она спроектирована так, чтобы создавать иллюзию того, что каждая строка в родительской таблице имеет *собственную* дочернюю таблицу. Например, при наличии 100 строк в родительской таблице *виртуально* имеется 100 вложенных таблиц. Физически существует только одна родительская и одна дочерняя таблица. Между вложенными и дочерними таблицами также имеются крупные синтаксические и семантические отличия, и они будут обсуждаться в этом разделе.

Вложенные таблицы используются двумя способами. Первый — их применение в PL/SQL-коде в качестве метода для расширения языка PL/SQL. Второй — их использование как физического механизма хранения для постоянного хранения коллекций. Я применяю вложенные таблицы в PL/SQL-коде все время, но еще ни разу не использовал их как механизм постоянного хранения.

В этом разделе кратко описывается синтаксис создания, запрашивания и модификации вложенных таблиц. Затем рассматриваются некоторые детали реализации и важные сведения о том, каким образом Oracle в действительности хранит вложенные таблицы.

Синтаксис вложенных таблиц

Процедура создания таблицы с вложенной таблицей довольно проста, но синтаксис для манипулирования ими немного сложнее. Давайте в целях демонстрации поработаем с таблицами EMP и DEPT. Вы уже знакомы с этой небольшой моделью данных, которая реализована реляционным образом с первичным и внешним ключами:

```
EODA@ORA12CR1> create table dept
  2  (deptno number(2) primary key,
  3   dname   varchar2(14),
  4   loc     varchar2(13)
  5  );
```

Table created.

Таблица создана.

```
EODA@ORA12CR1> create table emp
  2  (empno   number(4) primary key,
  3   ename   varchar2(10),
  4   job     varchar2(9),
  5   mgr     number(4) references emp,
  6   hiredate date,
  7   sal     number(7, 2),
  8   comm    number(7, 2),
  9   deptno  number(2) references dept
 10 );
```

Table created.

Таблица создана.

Теперь мы построим эквивалентную реализацию с представлением EMP в виде вложенной таблицы:

```
EODA@ORA12CR1> create or replace type emp_type
  2  as object
  3  (empno   number(4),
  4   ename   varchar2(10),
  5   job     varchar2(9),
  6   mgr     number(4),
  7   hiredate date,
  8   sal     number(7, 2),
  9   comm    number(7, 2)
 10 );
 11 /
```

Type created.

Тип создан.

```

EODA@ORA12CR1> create or replace type emp_tab_type
  2 as table of emp_type
  3 /
Type created.
Тип создан.

```

Чтобы создать таблицу с вложенной таблицей, нам необходим тип вложенной таблицы. В предыдущем выше коде создается тип сложного объекта EMP_TYPE, а затем из него — тип вложенной таблицы EMP_TAB_TYPE. В PL/SQL этот тип будет трактоваться во многом подобно массиву. В SQL он приведет к созданию физической вложенной таблицы. Ниже показан простой оператор CREATE TABLE, в котором применяется этот тип:

```

EODA@ORA12CR1> create table dept_and_emp
  2 (deptno number(2) primary key,
  3  dname varchar2(14),
  4  loc   varchar2(13),
  5  emps emp_tab_type
  6 )
  7 nested table emps store as emps_nt;
Table created.
Таблица создана.

EODA@ORA12CR1> alter table emps_nt add constraint
  2 emps_empno_unique unique(empno)
  3 /
Table altered.
Таблица изменена.

```

Важной частью этого оператора CREATE TABLE является включение столбца EMPS типа EMP_TAB_TYPE и соответствующей конструкции NESTED TABLE EMPS STORE AS EMPS_NT. В результате создается реальная физическая вложенная таблица EMPS_NT, отдельная и дополняющая таблицу DEPT_AND_EMP. Прямо к этой вложенной таблице добавляется ограничение на столбце EMPNO, чтобы обеспечить его уникальность, как это было в исходной реляционной модели. Реализовать нашу полную модель данных не получится; тем не менее, попробуем установить рефлексивное ограничение ссылочной целостности (self-referential constraint):

```

EODA@ORA12CR1> alter table emps_nt add constraint mgr_fk
  2 foreign key(mgr) references emps_nt(empno);
alter table emps_nt add constraint mgr_fk
*
ERROR at line 1:
ORA-30730: referential constraint not allowed on nested table column
ОШИБКА в строке 1:
ORA-30730: ограничение целостности на столбце вложенной таблицы не разрешено

```

Это просто не работает. Вложенные таблицы не поддерживают рефлексивные ограничения ссылочной целостности, т.к. они не могут ссылаться на любые другие таблицы — и даже на самих себя. Таким образом, в настоящей демонстрации мы это требование пропустим (естественно, в реальных условиях так поступать нельзя).

Теперь заполним таблицу DEPT_AND_EMP существующими данными из таблиц EMP и DEPT:

```
EODA@ORA12CR1> insert into dept_and_emp
2  select dept.*,
3     CAST( multiset( select empno, ename, job, mgr, hiredate, sal, comm
4                     from SCOTT.EMP
5                     where emp.deptno = dept.deptno ) AS emp_tab_type )
6  from SCOTT.DEPT
7  /
4 rows created.
```

Здесь необходимо обратить внимание на следующие два момента.

- **Было создано только четыре строки.** В таблице DEPT_AND_EMP действительно находятся всего лишь четыре строки. Четырнадцать строк EMP не существуют независимо.
- **Синтаксис становится довольно причудливым.** Ключевые слова CAST и MULTISSET относятся к синтаксису, которым большинство никогда не пользуется. При работе с объектно-реляционными компонентами в базе данных придется встречать много экзотического синтаксиса. Ключевое слово MULTISSET служит для указания Oracle о том, что подзапрос должен возвращать более одной строки (подзапросы в операторе SELECT ранее были ограничены возвратом одной строки). Ключевое слово CAST применяется для сообщения Oracle о том, что возвращаемый набор строк должен трактоваться как тип коллекции. В этом случае MULTISSET с помощью CAST приводится к типу EMP_TAB_TYPE. Универсальная процедура CAST не ограничивается использованием в коллекциях. Например, если требуется извлечь столбец EMPNO из таблицы EMP как относящийся к типу VARCHAR2 (20), а не NUMBER (4), то можно применить запрос `select cast(empno as VARCHAR2(20)) e from emp.`

Теперь все готово для запрашивания данных. Давайте посмотрим, как может выглядеть одна строка:

```
EODA@ORA12CR1> select deptno, dname, loc, d.emps AS employees
2  from dept_and_emp d
3  where deptno = 10
4  /
```

| DEPTNO | DNAME | LOC | EMPLOYEES (EMPNO, ENAME, JOB, |
|--------|------------|----------|---|
| 10 | ACCOUNTING | NEW YORK | EMP_TAB_TYPE (EMP_TYPE (7782, 'CLARK', 'MANAGER', 7839, '09-JUN-81', 2450, NULL), EMP_TYPE (7839, 'KING', 'PRESIDENT', NULL, '17-NOV-81', 5000, NULL), EMP_TYPE (7934, 'MILLER', 'CLERK', 7782, '23-JAN-82', 1300, NULL)) |

Все данные здесь находятся в единственном столбце. Многие приложения, если только они специально не разрабатывались с учетом поддержки объектно-реляционной функциональности, не смогут иметь дела с этим конкретным столбцом.

Например, ODBC не располагает средствами работы с вложенными таблицами (но их имеют JDBC, OCI, Pro*C, PL/SQL и многие другие API-интерфейсы и языки). В таких случаях Oracle предоставляет способ отмены вложенности коллекции и ее трактовки в основном как реляционной таблицы:

```

EODA@ORA12CR1> select d.deptno, d.dname, emp.*
  2   from dept_and_emp D, table(d.emps) emp
  3   /

```

| DEPTNO | DNAME | EMPNO | ENAME | JOB | MGR | HIREDATE | SAL | COMM |
|--------|------------|-------|--------|-----------|------|-----------|------|------|
| 10 | ACCOUNTING | 7782 | CLARK | MANAGER | 7839 | 09-JUN-81 | 2450 | |
| 10 | ACCOUNTING | 7839 | KING | PRESIDENT | | 17-NOV-81 | 5000 | |
| 10 | ACCOUNTING | 7934 | MILLER | CLERK | 7782 | 23-JAN-82 | 1300 | |
| 20 | RESEARCH | 7369 | SMITH | CLERK | 7902 | 17-DEC-80 | 800 | |
| 20 | RESEARCH | 7566 | JONES | MANAGER | 7839 | 02-APR-81 | 2975 | |
| 20 | RESEARCH | 7788 | SCOTT | ANALYST | 7566 | 09-DEC-82 | 3000 | |
| 20 | RESEARCH | 7876 | ADAMS | CLERK | 7788 | 12-JAN-83 | 1100 | |
| 20 | RESEARCH | 7902 | FORD | ANALYST | 7566 | 03-DEC-81 | 3000 | |
| 30 | SALES | 7499 | ALLEN | SALESMAN | 7698 | 20-FEB-81 | 1600 | 300 |
| 30 | SALES | 7521 | WARD | SALESMAN | 7698 | 22-FEB-81 | 1250 | 500 |
| 30 | SALES | 7654 | MARTIN | SALESMAN | 7698 | 28-SEP-81 | 1250 | 1400 |
| 30 | SALES | 7698 | BLAKE | MANAGER | 7839 | 01-MAY-81 | 2850 | |
| 30 | SALES | 7844 | TURNER | SALESMAN | 7698 | 08-SEP-81 | 1500 | 0 |
| 30 | SALES | 7900 | JAMES | CLERK | 7698 | 03-DEC-81 | 950 | |

```

14 rows selected.
14 строк выбрано.

```

Мы имеем возможность привести столбец EMPS к типу таблицы, и Oracle естественным образом выполняет соединение — никаких условий соединения указывать не понадобится. В действительности, т.к. тип EMP не имеет столбца DEPTNO, нет ничего такого, на чем можно было бы выполнить очевидное соединение. Об этом нюансе Oracle заботится самостоятельно.

А каким образом обновлять эти данные? Давайте представим, что требуется выдать отделу 10 премию в сумме \$100. Мы могли бы написать следующий код:

```

EODA@ORA12CR1> update
  2   table( select emps
  3           from dept_and_emp
  4           where deptno = 10
  5           )
  6   set comm = 100
  7   /
3 rows updated.
3 строки обновлено.

```

Именно здесь в игру вступает “иллюзия того, что каждая строка имеет отдельную таблицу”. В показанном ранее предикате SELECT, возможно, не было ясно, что на каждую строку предусмотрена отдельная таблица, особенно из-за отсутствия операций соединения и тому подобного; это было немного похоже на “магию”. Однако в операторе UPDATE видно, что на каждую строку приходится таблица. Для обновления выбрана отдельная таблица — она не имеет имени, а только идентифицирующий ее запрос. Если запрос не выбирает *в точности* одну таблицу, возникает ошибка:

```

EODA@ORA12CR1> update
  2   table( select emps
  3           from dept_and_emp
  4           where deptno = 1
  5         )
  6   set comm = 100
  7   /
update
*
ERROR at line 1:
ORA-22908: reference to NULL table value
ОШИБКА в строке 1:
ORA-22908: ссылка на нулевое значение таблицы

```

```

EODA@ORA12CR1> update
  2   table( select emps
  3           from dept_and_emp
  4           where deptno > 1
  5         )
  6   set comm = 100
  7   /
table( select emps
*
ERROR at line 2:
ORA-01427: single-row subquery returns more than one row
ОШИБКА в строке 2:
ORA-01427: однострочный подзапрос возвращает более одной строки

```

Если возвращается менее одной строки (один экземпляр вложенной таблицы), обновление работать отказывается. Обычно обновление нулевого числа строк допускается, но не в этом случае — здесь возникает ошибка, как и в ситуации, когда при обычном обновлении не указано имя таблицы. Если возвращается более одной строки (более одного экземпляра вложенной таблицы), операция обновления тоже не проходит. Обычно с обновлением множества строк проблем не возникает. Это говорит о том, что Oracle рассматривает каждую строку в таблице DEPT_AND_EMP как указывающую на другую таблицу, а не просто на другой набор строк, как это имеет место в реляционной модели.

В этом и состоит семантическое отличие между вложенной таблицей и родительской/дочерней реляционной таблицей. В модели с вложенной таблицей на каждую строку родительской таблицы приходится одна таблица, а в реляционной модели — один набор строк. Такое отличие иногда делает использование вложенных таблиц несколько громоздким. Например, возьмем нашу модель: она прекрасно представляет данные с точки зрения одиночного отдела, но совсем не годится для получения ответов на вопросы вроде “В каком отделе работает сотрудник по фамилии KING?”, “Сколько у нас работает бухгалтеров?” и т.п. Такие вопросы лучше задавать в реляционной таблице EMP, но в модели с вложенной таблицей доступ к данным EMP осуществляется только через данные DEPT. Мы должны всегда выполнять соединение; мы не можем запрашивать единственно данные таблицы EMP. Точнее, это нельзя сделать поддерживаемым и документированным способом, но можно воспользоваться одним трюком (который рассматривается позже). Если необходимо обновить

каждую строку в EMP\$_NT, пришлось бы выполнить четыре обновления: по одному для каждой строки в DEPT_AND_EMP, чтобы обновить ассоциированные с ними виртуальные таблицы.

Также обратите внимание, что при обновлении данных о сотрудниках отдела 10 мы семантически обновляли столбец EMP\$_ в таблице DEPT_AND_EMP. Мы понимаем, что физически были задействованы две таблицы, но с семантической точки зрения есть только одна таблица. И хотя в таблице отделов данные не обновлялись, строка, содержащая вложенную таблицу, которая действительно обновлялась, была заблокирована от обновления другими сеансами. В традиционной модели таблиц, объединенных отношением “родительская/дочерняя”, подобное бы не произошло.

Это и есть причины, по которым я предпочитаю не применять вложенные таблицы в качестве механизма постоянного хранения. Дочерняя таблица, которая бы не запрашивалась отдельно, встречается *редко*. В предыдущем примере таблица EMP\$_ должна быть отдельной сущностью. Она является автономной, поэтому нуждается в отдельных запросах. Я считаю, что так бывает в большинстве случаев. Я придерживаюсь использования вложенных таблиц через представления на реляционных таблицах.

Итак, теперь, когда вы видели, как обновлять экземпляр вложенной таблицы, вставка и удаление покажутся достаточно прямолинейными. Давайте добавим строку в экземпляр вложенной таблицы для отдела 10 и удалим строку из отдела 20:

```

EODA@ORA12CR1> insert into table
  2 ( select emps from dept_and_emp where deptno = 10 )
  3 values
  4 ( 1234, 'NewEmp', 'CLERK', 7782, sysdate, 1200, null );
1 row created.
1 строка создана.

```

```

EODA@ORA12CR1> delete from table
  2 ( select emps from dept_and_emp where deptno = 20 )
  3 where ename = 'SCOTT';
1 row deleted.
1 строка удалена.

```

```

EODA@ORA12CR1> select d.dname, e.empno, ename, deptno
  2 from dept_and_emp d, table(d.emps) e
  3 where d.deptno in ( 10, 20 );

```

| DNAME | EMPNO | ENAME | DEPTNO |
|------------|-------|--------|--------|
| ACCOUNTING | 7782 | CLARK | 10 |
| ACCOUNTING | 7839 | KING | 10 |
| ACCOUNTING | 7934 | MILLER | 10 |
| ACCOUNTING | 1234 | NewEmp | 10 |
| RESEARCH | 7369 | SMITH | 20 |
| RESEARCH | 7566 | JONES | 20 |
| RESEARCH | 7876 | ADAMS | 20 |
| RESEARCH | 7902 | FORD | 20 |

```

8 rows selected.
8 строк выбрано.

```

Вот так выглядит базовый синтаксис для запроса и модификации вложенных таблиц. Вы обнаружите, что для их применения часто необходимо отменять вложен-

ность, особенно в запросах. Как только вы сможете визуальнo представить концепцию “виртуальная таблица на строку”, работа с вложенными таблицами намного упростится.

Ранее я утверждал, что мы всегда должны производить соединение, и запрашивать только данные таблицы EMP нельзя, но существует один трюк, к которому можно прибегнуть, если это действительно требуется. Описанный далее метод официально не документирован, поэтому используйте его *только* в качестве крайней меры. Он наиболее полезен, когда возникает потребность в массовом обновлении вложенной таблицы (вспомните, что это должно делаться только через таблицу DEPT посредством соединения). Существует недокументированная подсказка NESTED_TABLE_GET_REFS (она кратко упоминается в документации, но полностью не описана), которую разнообразные инструменты (включая устаревшие утилиты EXP и IMP) применяют для работы с вложенными таблицами. Она также позволяет узнать немного больше о физической структуре вложенных таблиц. За счет использования этой подсказки можно получить ряд “магических” результатов. Ниже показан запрос, который применяет EXP (утилиту выгрузки данных) для извлечения данных из нашей вложенной таблицы:

```
EODA@ORA12CR1> SELECT /*+NESTED_TABLE_GET_REFS*/
2 NESTED_TABLE_ID,SYS_NC_ROWINFO$ FROM "EODA"."EMPS_NT";
```

| NESTED_TABLE_ID | SYS_NC_ROWINFO\$(EMPNO, ENAME, JOB, MGR, HIREDATE, |
|----------------------------------|--|
| EF6CDA23E32D315AE043B7D04F0AA620 | EMP_TYPE(7782, 'CLARK', 'MANAGER', 7839, '09-JUN-81', 2450, 100) |
| EF6CDA23E32D315AE043B7D04F0AA620 | EMP_TYPE(7839, 'KING', 'PRESIDENT', NULL, '17-NOV-81', 5000, 100) |
| ... | |

Вы несколько удивитесь, просмотрев описание этой таблицы:

```
EODA@ORA12CR1> desc emps_nt
```

| Name | Null? | Type |
|----------|-------|--------------|
| EMPNO | | NUMBER(4) |
| ENAME | | VARCHAR2(10) |
| JOB | | VARCHAR2(9) |
| MGR | | NUMBER(4) |
| HIREDATE | | DATE |
| SAL | | NUMBER(7,2) |
| COMM | | NUMBER(7,2) |

Столбцы NESTED_TABLE_ID и SYS_NC_ROWINFO\$ даже не показаны. Они являются частью скрытой реализации вложенных таблиц.

В действительности NESTED_TABLE_ID — это внешний ключ к родительской таблице DEPT_AND_EMP.

Таблица DEPT_AND_EMP на самом деле имеет скрытый столбец, который используется для соединения с EMPS_NT. Столбец SYS_NC_ROWINFO\$ является “магическим”; это больше функция, чем столбец. Вложенная таблица здесь является объектной таблицей (она создана из объектного типа), а SYS_NC_ROWINFO\$ представляет

собой внутренний способ, посредством которого Oracle ссылается на строку как на объект, вместо того чтобы ссылаться на каждый скалярный столбец. “За кулисами” Oracle реализует модель с родительской и дочерней таблицей и сгенерированными системой первичным и внешним ключами. Если интересуют дополнительные подробности, можно запросить реальный словарь данных и увидеть все столбцы таблицы DEPT_AND_EMP:

```
EODA@ORA12CR1> select name
2  from sys.col$
3  where obj# = ( select object_id
4    from dba_objects
5    where object_name = 'DEPT_AND_EMP'
6    and owner = 'EODA' )
7  /
```

NAME

DEPTNO

DNAME

EMPS

LOC

SYS_NC0000400005\$

Выбрав столбец SYS_NC0000400005\$ из вложенной таблицы, мы получим примерно такой вывод:

```
EODA@ORA12CR1> select SYS_NC0000400005$ from dept_and_emp;
```

SYS_NC0000400005\$

EF6CDA23E32D315AE043B7D04F0AA620

EF6CDA23E32E315AE043B7D04F0AA620

EF6CDA23E32F315AE043B7D04F0AA620

EF6CDA23E330315AE043B7D04F0AA620

Странно выглядящее имя столбца SYS_NC0000400005\$ — это сгенерированный системой ключ, помещенный в таблицу DEPT_AND_EMP. Углубившись в детали еще больше, мы обнаружим, что Oracle помещает на этот столбец уникальный индекс. К сожалению, Oracle пренебрегает индексацией по столбцу NESTED_TABLE_ID в EMPS_NT. По правде говоря, этот столбец нуждается в индексации, т.к. мы всегда производим соединение из DEPT_AND_EMP в EMPS_NT. Это важный момент, который следует запомнить: в случае применения вложенных таблиц со всеми стандартными параметрами, как делалось только что, *всегда индексируйте* в них столбец NESTED_TABLE_ID!

Однако здесь я немного отклонился от темы — речь шла о том, каким образом трактовать вложенную таблицу так, как если бы она была реальной таблицей. Это делает подсказка NESTED_TABLE_GET_REFS. Использовать ее можно следующим образом:

```
EODA@ORA12CR1> select /*+ nested_table_get_refs */ empno, ename
2  from emps_nt where ename like '%A%';
```

```

EMPNO ENAME
-----
7782 CLARK
7876 ADAMS
7499 ALLEN
7521 WARD
7654 MARTIN
7698 BLAKE
7900 JAMES
7 rows selected.
7 строк выбрано.

```

```

EODA@ORA12CR1> update /* nested_table_get_refs */ emps_nt set ename =
initcap(ename);
14 rows updated.
14 строк обновлено.

```

```

EODA@ORA12CR1> select /* nested_table_get_refs */ empno, ename
2 from emps_nt where ename like '%a%';

```

```

EMPNO ENAME
-----
7782 Clark
7876 Adams
7521 Ward
7654 Martin
7698 Blake
7900 James
6 rows selected.
6 строк выбрано.

```

Опять-таки, помните, что подсказка `NESTED_TABLE_GET_REFS` не является полностью документированным и поддерживаемым средством. Она предлагает специфическую функциональность для работы утилит `EXP` и `IMP`. Эти утилиты представляют собой единственную среду, в которой `NESTED_TABLE_GET_REFS` будет гарантированно работать. Применяйте подсказку `NESTED_TABLE_GET_REFS` на свой страх и риск, но воздерживайтесь от ее включения в производственный код. Фактически если вы обнаружите, что *нуждаетесь в ней*, то по определению вы вообще не намеревались использовать вложенную таблицу! Это неподходящая для вас конструкция. Применяйте подсказку `NESTED_TABLE_GET_REFS` только для однократных исправлений данных или чтобы ради любопытства посмотреть, как устроена вложенная таблица. Поддерживаемый способ получения отчета по данным во вложенной таблице предполагает отмену вложенности:

```

EODA@ORA12CR1> select d.deptno, d.dname, emp.*
2 from dept_and_emp D, table(d.emps) emp
3 /

```

Именно этот прием должен использоваться в запросах и производственном коде.

Хранение вложенных таблиц

Мы уже рассмотрели несколько аспектов, касающихся хранения структуры вложенной таблицы. В этом разделе будет подробно описана структура, которую Oracle создает по умолчанию, и методы управления ею. Воспользовавшись тем же оператором CREATE, что и ранее, мы знаем, что Oracle в действительности создаст структуру, подобную показанной на рис. 10.11:

```

EODA@ORA12CR1> create table dept_and_emp
  2  (deptno number(2) primary key,
  3   dname   varchar2(14),
  4   loc     varchar2(13),
  5   emps     emp_tab_type
  6  )
  7  nested table emps store as emps_nt;
Table created.
Таблица создана.

EODA@ORA12CR1> alter table emps_nt add constraint
  2  emps_empno_unique unique(empno)
  3  /
Table altered.
Таблица изменена.

```

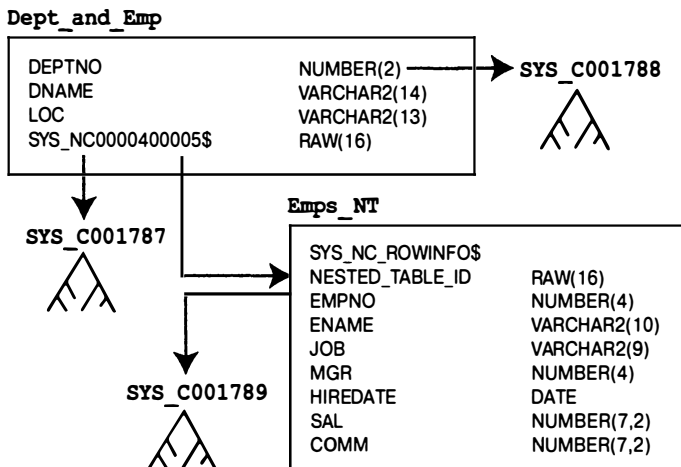


Рис. 10.11. Физическая реализация вложенной таблицы

Этот код создал две реальные таблицы. Таблица, которую мы затребовали, на месте, но в ней содержится дополнительный скрытый столбец (по умолчанию в таблице будет предусмотрен один дополнительный скрытый столбец для *каждого* столбца типа вложенной таблицы). Кроме того, на этом скрытом столбце устанавливается ограничение *уникальности*. База данных Oracle также создала вложенную таблицу EMPNS_NT. Эта таблица имеет два скрытых столбца, один из которых, SYS_NC_ROWINFO\$, является виртуальным столбцом, возвращающим все скалярные элементы в виде объекта.

Другой скрытый столбец представляет собой внешний ключ по имени NESTED_TABLE_ID, который может соединяться с родительской таблицей. Обратите внимание, что индекс на этом столбце *отсутствует*. И, наконец, Oracle добавляет индекс на столбце DEPTNO в таблице DEPT_AND_EMP, чтобы обеспечить наличие первичного ключа. Итак, мы попросили создать таблицу, а получили намного больше, чем ожидали. Полученная структура очень похожа на ту, которую вы могли бы создать для таблиц, связанных отношением “родительская/дочерняя”, но в таком случае в качестве внешнего ключа таблицы EMPS_NT вы применяли бы существующий первичный ключ на столбце DEPTNO, а не генерировали суррогатный ключ RAW(16).

Дамп DBMS_METADATA.GET_DDL нашего примера вложенной таблицы выглядит следующим образом:

```
EODA@ORA12CR1> begin
  2   dbms_metadata.set_transform_param
  3   ( DBMS_METADATA.SESSION_TRANSFORM, 'STORAGE', false );
  4 end;
  5 /
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.
EODA@ORA12CR1> select dbms_metadata.get_ddl( 'TABLE', 'DEPT_AND_EMP' )
  6 from dual;
DBMS_METADATA.GET_DDL('TABLE','DEPT_AND_EMP')
-----
CREATE TABLE "EODA"."DEPT_AND_EMP"
(
  "DEPTNO" NUMBER(2,0),
  "DNAME" VARCHAR2(14),
  "LOC" VARCHAR2(13),
  "EMPS" "EODA"."EMP_TAB_TYPE" ,
  PRIMARY KEY ("DEPTNO")
USING INDEX PCTFREE 10 INITRANS 2 MAXTRANS 255
TABLESPACE "USERS" ENABLE
) SEGMENT CREATION IMMEDIATE
PCTFREE 10 PCTUSED 40 INITRANS 1 MAXTRANS 255
NOCOMPRESS LOGGING
TABLESPACE "USERS"
NESTED TABLE "EMPS" STORE AS "EMPS_NT"
(( CONSTRAINT "EMPS_EMPNO_UNIQUE" UNIQUE ("EMPNO")
  USING INDEX PCTFREE 10 INITRANS 2 MAXTRANS 255 COMPUTE STATISTICS
  TABLESPACE "USERS" ENABLE)
PCTFREE 10 PCTUSED 40 INITRANS 1 MAXTRANS 255 LOGGING
TABLESPACE "USERS" ) RETURN AS VALUE
```

Единственной новой конструкцией здесь является RETURN AS VALUE. Она используется для описания того, каким образом вложенная таблица возвращается клиентскому приложению. По умолчанию Oracle будет возвращать вложенную таблицу клиенту по значению; с каждой строкой будут передаваться действительные данные. Можно также применить конструкцию RETURN AS LOCATOR, которая приводит к тому, что клиент получит указатель на данные, а не сами данные. Данные будут переданы клиенту, если — и только если — он разыменует этот указатель.

Таким образом, если вы полагаете, что клиент в основном не будет просматривать строки вложенной таблицы для каждой строки родительской таблицы, можете возвращать локатор вместо значений, сэкономив на круговых обменах по сети. Например, если есть клиентское приложение, которое отображает списки отделов, а при двойном щелчке пользователя на отделе показывает информацию о сотрудниках этого отдела, то лучше использовать локатор. Причина в том, что подробная информация обычно не просматривается — это исключение, а не правило.

Итак, что еще можно сделать с вложенной таблицей? Прежде всего, столбец `NESTED_TABLE_ID` должен быть проиндексирован. Поскольку мы всегда производим доступ во вложенную таблицу в направлении *от* родительской *к* дочерней, то действительно нуждаемся в таком индексе. Проиндексировать этот столбец можно с помощью оператора `CREATE INDEX`, но более удачное решение предусматривает применение индекс-таблицы для хранения вложенной таблицы. Вложенная таблица является еще одним замечательным примером, когда индекс-таблица подходит как нельзя лучше. В индекс-таблице дочерние строки будут физически размещаться рядом сгруппированными по столбцу `NESTED_TABLE_ID` (а это значит, что для их извлечения будет требоваться меньший объем физического ввода-вывода). Кроме того, устраняется необходимость в избыточном индексе на столбце `RAW(16)`. Далее, поскольку столбец `NESTED_TABLE_ID` будет ведущим столбцом в первичном ключе индекс-таблицы, мы должны также включить сжатие ключей индекса, чтобы подавить избыточные значения `NESTED_TABLE_ID`, которые будут появляться в противном случае. Вдобавок можно также включить в команду `CREATE TABLE` ограничение `UNIQUE` и `NOT NULL` на столбце `EMPNO`. Таким образом, если взять предыдущий оператор `CREATE TABLE` и слегка его модифицировать, то получится набор объектов:

```

EODA@ORA12CR1> CREATE TABLE "EODA"."DEPT_AND_EMP"
 2  ("DEPTNO" NUMBER(2, 0),
 3  "DNAME"  VARCHAR2(14),
 4  "LOC"    VARCHAR2(13),
 5  "EMPS"   "EMP_TAB_TYPE")
 6  PCTFREE 10 PCTUSED 40 INITRANS 1 MAXTRANS 255 LOGGING
 7  TABLESPACE "USERS"
 8  NESTED TABLE "EMPS"
 9  STORE AS "EMPS_NT"
10  ((empno NOT NULL, unique (empno), primary key(nested_table_id,empno))
11  organization index compress 1 )
12  RETURN AS VALUE;
```

Table created.

Таблица создана.

Вместо обычной таблицы `EMPS_NT` мы имеем индекс-таблицу `EMPS_NT`, на что указывает наложенная поверх нее структура индекса на рис. 10.12.

Когда таблица `EMPS_NT` представляет собой индекс-таблицу, использующую сжатие, она будет занимать меньше места, чем первоначальная стандартная вложенная таблица, и будет иметь индекс, в котором мы крайне нуждаемся.

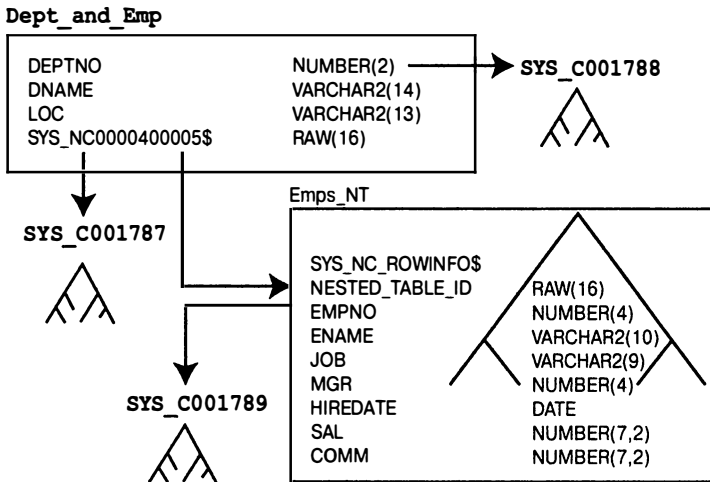


Рис. 10.12. Вложенная таблица, реализованная в виде индекс-таблицы

Заключительные соображения по поводу вложенных таблиц

По описанным ниже причинам я не применяю вложенные таблицы в качестве механизма постоянного хранения.

- Излишние накладные расходы, связанные с хранением добавляемых столбцов RAW (16). Этот дополнительный столбец будет иметь как родительская, так и дочерняя таблица. В родительской таблице столбец RAW (16) добавляется для каждого присутствующего в ней столбца типа вложенной таблицы. Из-за того, что родительская таблица обычно уже располагает первичным ключом (DEPTNO в приведенных ранее примерах), в дочерних таблицах имеет смысл использовать его, а не ключ, сгенерированный системой.
- Ненужные накладные расходы, связанные с дополнительным ограничением уникальности на родительской таблице, в то время как она обычно уже имеет такое ограничение.
- Вложенную таблицу отдельно использовать непросто, не прибегая к применению неподдерживаемых конструкций (NESTED_TABLE_GET_REFS). Можно отменить вложенность таблицы для запросов, но не массовых обновлений. В реальности я еще не сталкивался с таблицей, которая бы не запрашивалась “сама по себе”.

Я интенсивно использую вложенные таблицы в коде и в представлениях. Я считаю, что именно здесь они пребывают в своей стихии. Для механизма хранения лично я предпочитаю создавать родительские и дочерние таблицы. После создания родительской и дочерней таблиц можно построить представление, которое позволит им выглядеть так, как если бы мы располагали настоящей вложенной таблицей. То есть, мы можем получить все преимущества конструкции вложенной таблицы без присущих ей накладных расходов.

Если вы все-таки решите применять вложенную таблицу как механизм хранения, то лучше сделайте ее индекс-таблицей, чтобы избежать накладных расходов,

связанных с индексом на столбце `NESTED_TABLE_ID` и самой вложенной таблицей. Воспользуйтесь советами из предыдущего раздела, посвященного индекс-таблицам, для настройки сегментов переполнения и других параметров. Если вы не захотите применять индекс-таблицу, позаботьтесь об индексировании столбца `NESTED_TABLE_ID` во вложенной таблице, чтобы избежать полного ее сканирования при поиске дочерних строк.

Временные таблицы

Временные таблицы используются для сохранения промежуточных результирующих наборов на протяжении либо сеанса, либо транзакции. Данные, находящиеся во временной таблице, видимы только текущему сеансу: сеансы не могут просматривать данные друг друга, даже если текущий сеанс фиксирует изменения. Проблема с конкуренцией множества пользователей в отношении временных таблиц тоже не возникает, т.к. один сеанс никогда не сможет заблокировать другой сеанс из-за работы с временной таблицей. Даже если мы “заблокируем” временную таблицу, это не мешает другим сеансам пользоваться своими временными таблицами.

На заметку! Как было показано в главе 9, временные таблицы генерируют значительно меньший объем информации redo, чем обычные таблицы. Однако, поскольку временные таблицы генерируют информацию undo для содержащихся в них данных, они будут генерировать определенный объем информации redo. Наибольшее количество данных redo будут генерировать операторы `UPDATE` и `DELETE`, а наименьшее — `INSERT` и `SELECT`. В главе 9 также упоминалось о том, что в версии Oracle 12c появилась возможность конфигурировать временные таблицы для генерации близкого к нулю объема данных redo; это делается путем установки параметра `TEMP_UNDO_ENABLED` в `TRUE`.

Пространство для хранения временных таблиц выделяется во временном табличном пространстве пользователя, который в текущий момент находится в системе. Если же доступ к временным таблицам производится из процедуры, выполняемой с правами владельца, то будет применяться временное табличное пространство владельца этой процедуры. Глобальная временная таблица — это на самом деле всего лишь шаблон для самой такой таблицы. Действие по созданию временной таблицы не подразумевает распределение пространства хранения; никакого начального экстенда не выделяется, как было бы в случае традиционной таблицы (если только не активизировано средство отложенного выделения сегментов). Вместо этого, когда сеанс во время выполнения впервые помещает данные во временную таблицу, для этого сеанса создается временный сегмент. Поскольку каждый сеанс получает собственный временный сегмент (а не просто экстенд существующего сегмента), каждый пользователь может выделять пространство под свои временные таблицы в собственном табличном пространстве. Например, для пользователя `USER1` может быть указано табличное пространство `TEMP1`, поэтому его временные таблицы будут размещаться именно в этом табличном пространстве. Для пользователя `USER2` может быть задано табличное пространство `TEMP2`, так что его временные таблицы будут находиться в нем.

Временные таблицы Oracle похожи на временные таблицы в других реляционных базах данных, но с важным отличием: они определяются *статически*. Вы создаете

их один раз для базы данных, а не один раз для каждой хранимой процедуры в этой базе данных. Временные таблицы существуют всегда — они будут присутствовать в словаре данных как объекты, но будут всегда выглядеть пустыми до тех пор, пока ваш сеанс не поместит в них данные. Тот факт, что они определены статически, позволяет создавать представления, которые ссылаются на них, хранимые процедуры, которые используют статический SQL для ссылки на них, и т.д.

Временные таблицы могут быть основаны на *сеансе* (данные сохраняются в таблице между фиксациями, но не между отключением и повторным подключением). Временные таблицы также могут быть основаны на *транзакции* (данные исчезают после фиксации). Ниже приведен пример того и другого поведения. В качестве шаблона применяется таблица SCOTT.EMP:

```
EODA@ORA12CR1> create global temporary table temp_table_session
2 on commit preserve rows
3 as
4 select * from scott.emp where 1=0
5 /
```

Table created.

Таблица создана.

Конструкция ON COMMIT PRESERVE ROWS делает временную таблицу основанной на сеансе. Строки будут оставаться в этой таблице до тех пор, пока сеанс не отключится или пока они не будут физически удалены с помощью DELETE или TRUNCATE. Эти строки доступны только текущему сеансу; никакой другой сеанс не будет их видеть даже после выполнения фиксации.

```
EODA@ORA12CR1> create global temporary table temp_table_transaction
2 on commit delete rows
3 as
4 select * from scott.emp where 1=0
5 /
```

Table created.

Таблица создана.

Конструкция ON COMMIT DELETE ROWS делает временную таблицу основанной на транзакции. Когда сеанс производит фиксацию, строки исчезают. Строки пропадают за счет возвращения обратно временных экстенгов, выделенных для таблицы — никаких накладных расходов с автоматической очисткой временных таблиц не связано. Теперь давайте посмотрим на отличия между этими двумя типами таблиц:

```
EODA@ORA12CR1> insert into temp_table_session select * from scott.emp;
14 rows created.
14 строк создано.
```

```
EODA@ORA12CR1> insert into temp_table_transaction select * from scott.emp;
14 rows created.
14 строк создано.
```

В каждую из двух временных таблиц помещено по 14 строк, и мы можем в этом убедиться:

```
EODA@ORA12CR1> select session_cnt, transaction_cnt
2 from ( select count(*) session_cnt from temp_table_session ),
3      ( select count(*) transaction_cnt from temp_table_transaction );
```

```

SESSION_CNT TRANSACTION_CNT
-----
14          14
EODA@ORA12CR1> commit;

```

Так как произведена фиксация, мы увидим строки во временной таблице, основанной на сеансе, но не строки во временной таблице, основанной на транзакции:

```

EODA@ORA12CR1> select session_cnt, transaction_cnt
2   from ( select count(*) session_cnt from temp_table_session ),
3         ( select count(*) transaction_cnt from temp_table_transaction
);
SESSION_CNT TRANSACTION_CNT
-----
14          0
EODA@ORA12CR1> disconnect
Disconnected from Oracle Database 12c Enterprise Edition Release
12.1.0.1.0 - 64bit Production
With the Partitioning, OLAP, Advanced Analytics and Real Application
Testing options
EODA@ORA12CR1> connect eoda
Enter password:
Connected.
Подключено.

```

Поскольку запущен новый сеанс, мы не увидим никаких строк в обеих таблицах:

```

EODA@ORA12CR1> select session_cnt, transaction_cnt
2   from ( select count(*) session_cnt from temp_table_session ),
3         ( select count(*) transaction_cnt from temp_table_transaction
);
SESSION_CNT TRANSACTION_CNT
-----
0          0

```

Чтобы проверить, была ли таблица создана как временная, а также продолжительность хранения данных в ней (в течение сеанса или транзакции), можно запросить столбцы TEMPORARY и DURATION представления USER_TABLES. Стандартным значением в DURATION является SYS\$TRANSACTION (означает ON COMMIT DELETE ROWS). В рассматриваемом примере эти значения выглядят следующим образом:

```

EODA@ORA12CR1> select table_name, temporary, duration from user_tables;
TABLE_NAME          T DURATION
-----
TEMP_TABLE_TRANSACTION  Y SYS$TRANSACTION
TEMP_TABLE_SESSION     Y SYS$SESSION

```

При наличии опыта работы с временными таблицами в SQL Server и/или Sybase важно принять во внимание, что вместо выполнения оператора SELECT X, Y, Z INTO #TEMP FROM SOME_TABLE для динамического создания и заполнения временной таблицы придется предпринимать следующие действия.

- Создать все глобальные временные таблицы один раз как часть процесса установки приложения подобно созданию постоянных таблиц.
- Применять в своих процедурах простой оператор `INSERT INTO TEMP (X, Y, Z) SELECT X, Y, Z FROM SOME_TABLE`.

Речь идет о том, что ваша цель здесь — вовсе не создание таблиц в хранимых процедурах во время выполнения. Это не подходящий способ использования временных таблиц в Oracle. Операция DDL является затратной в смысле ресурсов; во время выполнения ее желательно избегать. Временные таблицы для приложения должны создаваться в процессе его установки и *никогда* во время выполнения. Ниже перечислены сложности, с которыми вы столкнетесь, если попытаетесь динамически создавать глобальные временные таблицы (или просто таблицы вообще) во время выполнения в коде PL/SQL.

- Вы будете инициировать операции DDL во время выполнения. Они являются исключительно затратными и вовлекают сотни рекурсивных SQL-операторов. С операциями DDL также связан значительный объем сериализации (они выполняются последовательно друг за другом).
- Для работы с такими таблицами придется применять динамический SQL в коде PL/SQL. Вы утратите все преимущества статического компилируемого SQL. Это огромная потеря.
- У вас не будет возможности одновременно запускать две копии хранимой процедуры. Поскольку оба экземпляра хранимой процедуры будут пытаться уничтожить и создавать одну и ту же временную таблицу, между ними возникнет конфликт. (В такой ситуации можно было бы динамически генерировать уникальное имя для временной таблицы каждый раз, когда процедура создается, но это привносит дополнительную сложность и потенциальную головную боль при поиске и устранении проблем.)
- В конечном итоге вы столкнетесь с тем, что однажды код не сможет корректно удалить временные таблицы. Из-за непредвиденной ошибки (например, сбоя электропитания) процедура может не завершиться. После возобновления нормального электропитания временная таблица продолжит существовать. Вам придется время от времени вручную удалять такие таблицы.

Короче говоря, не существует веских причин для того, чтобы создавать таблицы в коде PL/SQL во время выполнения, а есть только причины никогда этого *не* делать.

Временные таблицы могут иметь многие атрибуты постоянных таблиц. Они могут располагать триггерами, проверочными ограничениями целостности, индексами и т.д. Ниже перечислены возможности постоянных таблиц, которые временные таблицы не поддерживают.

- Они не могут иметь ограничений ссылочной целостности. Они не могут быть ни целевым столбцом внешнего ключа, ни содержать определенный в них внешний ключ.
- Они не могут иметь столбцы типа `NESTED TABLE`. В Oracle9i и предшествующих версиях они не могли также содержать столбцы типа `VARRAY`; это ограничение устранено, начиная с версии Oracle 10g.

- Они не могут быть индекс-таблицами.
- Они не могут находиться в кластере любого типа.
- Они не могут быть секционированы.
- Они не позволяют генерировать статистические данные с помощью табличной команды ANALYZE.

Один из недостатков временной таблицы в любой базе данных связан с тем, что оптимизатор обычно не располагает реальными статистическими данными о ней. При использовании *оптимизатора по стоимости* (cost-based optimizer — CBO) наличие действительных статистических данных жизненно важно для успешной его работы. Если статистические данные отсутствуют, оптимизатор будет делать предположения относительно распределения данных, их объема и селективности индекса. Если эти его предположения оказываются некорректными, то планы, которые оптимизатор CBO генерирует для запросов, интенсивно работающих с временными таблицами, могут оказаться далеко не оптимальными. Во многих случаях правильным решением будет вообще не применять временную таблицу, а использовать вместо нее представление `INLINE VIEW` (пример представления `INLINE VIEW` можно найти в последнем выполненном операторе `SELECT` — он содержит два таких представления). Это позволит Oracle иметь доступ ко всем необходимым статистическим данным для таблицы и составлять оптимальный план.

По моим наблюдениям, люди часто применяют временные таблицы из-за того, что в средах других баз данных они привыкли к тому, что соединение слишком большого количества таблиц в единственном запросе является плохой практикой. При разработке в среде Oracle от этой привычки следует избавиться. Вместо того чтобы пытаться перехитрить оптимизатор и разбивать то, что должно быть одним запросом, на три или четыре запроса, которые сохраняют свои частичные результаты во временных таблицах, и затем объединять эти временные таблицы, вы должны просто написать одиночный запрос, решающий исходную задачу. Ссылаться на множество таблиц в одном запросе вполне допустимо; поддержка со стороны временных таблиц для этой цели не нужна.

Однако бывают случаи, когда использование временной таблицы в процессе представляет собой корректный подход. Например, однажды я разрабатывал для карманного компьютера Palm приложение синхронизации данных книги дат Palm Pilot с календарной информацией, хранящейся в Oracle. В Palm находился список записей, которые были изменены с момента последней синхронизации. Я должен был взять эти записи, сравнить их с актуальными данными в базе, обновить записи базы данных и затем сгенерировать список изменений, предназначенный для применения к данным в Palm. Это великолепный пример ситуации, когда временная таблица оказывается очень полезной. Я использовал временную таблицу для хранения изменений из Palm в базе данных. Затем я запустил хранимую процедуру, которая сравнивала сгенерированные Palm изменения с актуальными (и очень большими) постоянными таблицами, чтобы выяснить, какие изменения должны быть внесены в данные Oracle, а какие — в данные Palm. Я должен был сделать несколько проходов по этим данным. Сначала я искал все записи, которые были модифицированы только в Palm, и вносил соответствующие изменения в базу данных Oracle. Затем я искал все записи, которые изменились на стороне Palm и в моей базе данных с мо-

мента последней синхронизации, и корректировал их. Далее я искал все записи, которые были модифицированы только в базе данных, и помещал изменения во временную таблицу. Наконец, приложение синхронизации Palm извлекало изменения из временной таблицы и применяло их к самому устройству Palm. После закрытия подключения временные данные исчезали.

Однако я столкнулся с проблемой, которая заключалась в том, что поскольку постоянные таблицы анализировались, использовался оптимизатор СВО. Временная таблица не имела статистических данных (анализировать временную таблицу можно, но статистика при этом не собирается), поэтому оптимизатор СВО должен был делать много предположений об этой таблице. Мне, как разработчику, было известно ожидаемое среднее количество строк, распределение данных, селективность индексов и т.п. Мне нужен был способ информирования оптимизатора об имеющихся *лучших* предположениях. Это делается посредством генерации статистических данных для временной таблицы. Итак, мы вплотную подошли к теме, касающейся того, каким образом генерируется статистика для временной таблицы.

На заметку! По причине значительных усовершенствований в сборе статистики по временным таблицам, появившихся в версии Oracle 12c, я решил разбить эту тему на два раздела: "Сбор статистики до версии Oracle 12c" и "Сбор статистики, начиная с версии Oracle 12c".

Сбор статистики до версии Oracle 12c

Существуют три способа передачи оптимизатору статистических данных по глобальным временным таблицам. Один из них предполагает выполнение динамической выборки (в Oracle9i Release 2 и последующих версиях), а другой — применение пакета DBMS_STATS, который позволяет делать это двумя путями. Давайте сначала взглянем на динамическую выборку.

Динамическая выборка (dynamic sampling) — это возможность оптимизатора при полном разборе запроса сканировать сегменты в базе данных (производить их выборку) с целью сбора статистических данных, полезных для оптимизации этого конкретного запроса. Она похожа на миниатюрную команду сбора статистики во время полного разбора. В Oracle 10g и последующих версиях динамическая выборка будет работать сразу же после установки, т.к. стандартный уровень динамической выборки был увеличен с 1 до 2. На уровне 2 перед построением плана выполнения запроса оптимизатор будет осуществлять динамическую выборку каждого не прошедшего анализ объекта, на который имеется ссылка в обрабатываемом запросе. В версии Oracle9i Release 2 уровень динамической выборки 1 приводил к тому, что она использовалась менее часто. Чтобы обеспечить в Oracle9i Release 2 такое же поведение динамической выборки, как по умолчанию в Oracle 10g, можно применить команду ALTER SESSION|SYSTEM или подсказку dynamic_sampling:

```
ops$tkyte@ORA9IR2> create global temporary table gtt
2 as
3 select * from scott.emp where l=0;
Table created.
Таблица создана.
```

```
ops$tkyte@ORA9IR2> insert into gtt select * from scott.emp;
14 rows created.
14 строк создано.
```

```
ops$tkyte@ORA9IR2> set autotrace traceonly explain
ops$tkyte@ORA9IR2> select /*+ first_rows */ * from gtt;
```

Execution Plan

```
-----
 0      SELECT STATEMENT Optimizer=HINT: FIRST_ROWS (Cost=17 Card=8168 Bytes...
 1      0      TABLE ACCESS (FULL) OF 'GTT' (Cost=17 Card=8168 Bytes=710616)
```

```
ops$tkyte@ORA9IR2> select /*+ first_rows dynamic_sampling(gtt 2) */ * from gtt;
```

Execution Plan

```
-----
 0      SELECT STATEMENT Optimizer=HINT: FIRST_ROWS (Cost=17 Card=14 Bytes=1218)
 1      0      TABLE ACCESS (FULL) OF 'GTT' (Cost=17 Card=14 Bytes=1218)
```

```
ops$tkyte@ORA9IR2> set autotrace off
```

Мы установили уровень динамической выборки в 2 для таблицы GTT в этом запросе. Сам по себе оптимизатор предположил бы, что из таблицы GTT должны возвращаться 8168 строк.

На заметку! Стандартное значение 8168 — в действительности функция от стандартного размера блока. В базе данных с размером блоков 4 Кбайт количество ожидаемых строк будет меньше, а при блоках 16 Кбайт — больше.

За счет использования динамической выборки оценка предполагаемой кардинальности окажется намного ближе к реальности (что приведет получению в целом лучших планов выполнения запросов). Установка уровня 2 привела к тому, что оптимизатор быстро просканировал таблицу и получил более реалистичные оценки настоящего размера таблицы GTT. В Oracle 10g и последующих версиях вы должны обнаружить, что проблема с этим оператором уменьшилась, т.к. по умолчанию будет происходить динамическая выборка:

```
EODA@ORA11GR2> create global temporary table gtt
2 as
3 select * from scott.emp where 1=0;
Table created.
Таблица создана.
```

```
EODA@ORA11GR2> insert into gtt select * from scott.emp;
14 rows created.
14 строк создано.
```

```
EODA@ORA11GR2> set autotrace traceonly explain
EODA@ORA11GR2> select * from gtt;
```

| Id | Operation | Name | Rows | Bytes | Cost (%CPU) | Time |
|----|-------------------|------|------|-------|-------------|----------|
| 0 | SELECT STATEMENT | | 14 | 1218 | 2 (0) | 00:00:01 |
| 1 | TABLE ACCESS FULL | GTT | 14 | 1218 | 2 (0) | 00:00:01 |

Note

- dynamic sampling used for this statement (level=2)

Примечание

- для этого оператора использовалась динамическая выборка (уровень 2)

EODA@ORA11GR2> set autotrace off

Мы получили правильную кардинальность без специального запроса для этого. Тем не менее, динамическая выборка не обходится бесплатно: ее выполнение во время разбора запроса связано с определенными затратами. Если собрать необходимые статистические данные заранее, то таких затрат можно избежать во время полного разбора. В этом может помочь пакет DMBS_STATS.

Существуют три метода применения пакета DMBS_STATS для сбора репрезентативных статистических данных. Первый метод предусматривает использование DMBS_STATS с вызовом GATHER_SCHEMA_STATS или GATHER_DATABASE_STATS. Эти процедуры принимают булевский параметр GATHER_TEMP, который имеет стандартное значение FALSE. Когда он установлен в TRUE, для любой глобальной временной таблицы ON COMMIT PRESERVER ROWS будет собираться и сохраняться статистика (этот прием не работает для таблиц ON COMMIT DELETE). Взгляните на следующий код (обратите внимание, что он выполнялся в пустой схеме; созданы только те объекты, которые вы видите):

```
EODA@ORA11GR2> create table emp as select * from scott.emp;
```

Table created.

Таблица создана.

```
EODA@ORA11GR2> create global temporary table gtt1 ( x number )
```

```
2 on commit preserve rows;
```

Table created.

Таблица создана.

```
EODA@ORA11GR2> create global temporary table gtt2 ( x number )
```

```
2 on commit delete rows;
```

Table created.

Таблица создана.

```
EODA@ORA11GR2> insert into gtt1 select user_id from all_users;
```

49 rows created.

49 строк создано.

```
EODA@ORA11GR2> insert into gtt2 select user_id from all_users;
```

49 rows created.

49 строк создано.

```
EODA@ORA11GR2> exec dbms_stats.gather_schema_stats( user );
```

PL/SQL procedure successfully completed.

Процедура PL/SQL успешно завершена.

```
EODA@ORA11GR2> select table_name, last_analyzed, num_rows from user_tables;
```

| TABLE_NAME | LAST_ANALYZED | NUM_ROWS |
|------------|---------------|----------|
| EMP | 17-JAN-14 | 14 |
| GTT2 | | |
| GTT1 | | |

Как видите, в этом случае была проанализирована только таблица EMP; две глобальные временные таблицы были проигнорированы. Такое поведение можно изменить, вызвав процедуру GATHER_SCHEMA_STATS с параметром GATHER_TEMP=>TRUE:

```
EODA@ORA11GR2> insert into gtt2 select user_id from all_users;
49 rows created.
49 строк создано.
```

```
EODA@ORA11GR2> exec dbms_stats.gather_schema_stats( user, gather_temp=>TRUE );
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.
```

```
EODA@ORA11GR2> select table_name, last_analyzed, num_rows from user_tables;
```

| TABLE_NAME | LAST_ANALYZED | NUM_ROWS |
|------------|---------------|----------|
| EMP | 17-JAN-14 | 14 |
| GTT1 | 17-JAN-14 | 49 |
| GTT2 | 17-JAN-14 | 0 |

Обратите внимание на то, что у таблицы ON COMMIT PRESERVE ROWS имеются точные статистические данные, а у таблицы ON COMMIT DELETE ROWS — нет. Процедура пакета DBMS_STATS выполняет фиксацию, и вся информация в этой таблице очищается. Однако учтите, что таблица GTT теперь имеет статистические данные, что в действительности *очень плохо*, т.к. эти данные совершенно некорректны! Сомнительно, что таблица будет содержать 0 строк во время выполнения. Таким образом, если вы применяете этот подход, помните о двух моментах.

- Обеспечьте заполнение глобальных временных таблиц репрезентативными данными в сеансе, который собирает статистику. В противном случае для DBMS_STAT они будут выглядеть пустыми.
- Если существуют глобальные временные таблицы ON COMMIT DELETE ROWS, то такой подход не должен использоваться, потому что определенно будут собраны некорректные статистические сведения.

Второй прием, который работает с глобальными временными таблицами ON COMMIT PRESERVE ROWS, предусматривает применение процедуры GATHER_TABLE_STATS прямо на таблице. Глобальную временную таблицу можно заполнить данными, как это делалось только что, и затем запустить на ней GATHER_TABLE_STATS. Однако следует отметить, что этот прием, как и предыдущий, *не будет работать* для глобальных временных таблиц ON COMMIT DELETE ROWS из-за возникновения тех же самых проблем.

Последний прием с использованием пакета DBMS_STATS заключается в ручном заполнении словаря данных репрезентативными статистическими сведениями для временных таблиц. Например, если среднее количество строк во временной таблице будет равно 500, средний размер строк — 100 байтов, а количество блоков — 7, можно запустить приведенный ниже код:

```
EODA@ORA11GR2> create global temporary table t ( x int, y varchar2(100) )
2 on commit preserve rows;
Table created.
Таблица создана.
```

```
EODA@ORA11GR2> begin
2 dbms_stats.set_table_stats( ownname => USER,
3 tabname => 'T',
4 numrows => 500,
5 numblks => 7,
6 avgflen => 100 );
7 end;
8 /
```

PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.

```
EODA@ORA11GR2> select table_name, num_rows, blocks, avg_row_len
2 from user_tables
3 where table_name = 'T';
```

| TABLE_NAME | NUM_ROWS | BLOCKS | AVG_ROW_LEN |
|------------|----------|--------|-------------|
| T | 500 | 7 | 100 |

Теперь оптимизатор не будет применять свои предположения, а воспользуется *нашими* оценками для этой информации.

Сбор статистики, начиная с версии Oracle 12c

Начиная с версии Oracle 12c, сбор и использование статистики о глобальных временных таблицах чрезвычайно улучшились. Ниже перечислены основные изменения.

- По умолчанию статистические данные уровня сеанса генерируются при сборе статистики для временных таблиц.
- Разделяемую статистику по-прежнему можно собирать (во многом подобно тому, как это делалось в версии Oracle 11g), но вы должны сначала установить параметр GLOBAL_TEMP_TABLE_STATS (процедурой DBMS_STATS.SET_TABLE_PREFS) в SHARED.
- Для временных таблиц, определенных как ON COMMIT DELETE ROWS, некоторые процедуры пакета DBMS_STATS (такие как GATHER_TABLE_STATS) больше не выдают неявный оператор COMMIT; следовательно, для такого типа временных таблиц допускается генерировать репрезентативные статистические сведения.
- Для временных таблиц, определенных как ON COMMIT PRESERVE ROWS, статистические данные уровня сеанса автоматически генерируются для операций в прямом режиме (наподобие CREATE TABLE AS SELECT и операторов INSERT прямого режима). Это устраняет необходимость в обращении к DBMS_STATS с целью генерации статистических данных для указанных специальных операций.

Далее мы более подробно рассмотрим все эти изменения, начав со статистики сеанса.

Статистика сеанса

До выхода версии Oracle 12c статистические данные, сгенерированные для временной таблицы, совместно использовались всеми сеансами, работающими с этой временной таблицей. Это могло приводить к получению не дотягивающих до идеала планов выполнения, особенно если разные сеансы генерировали несоизмеримые объемы данных либо имели варьирующиеся модели данных.

Начиная с Oracle 12c, генерируемые статистические данные для временной таблицы специфичны для сеанса, который производит эту генерацию. Такая реализация предоставляет оптимизатору Oracle более точную информацию для создания плана выполнения, подстроенного под данные, которые сгенерированы в сеансе. Продемонстрируем это на небольшом примере, создав для начала временную таблицу:

```
EODA@ORA12CR1> create global temporary table gt(x number) on commit
  1 preserve rows;
  2 Table created.
  3 Таблица создана.
```

Затем вставим в нее некоторые данные:

```
EODA@ORA12CR1> insert into gt select user_id from all_users;
  1 51 rows created.
  2 51 строк создано.
```

Теперь сгенерируем статистику для этой таблицы:

```
EODA@ORA12CR1> exec dbms_stats.gather_table_stats( user, 'GT' );
  1 PL/SQL procedure successfully completed.
  2 Процедура PL/SQL успешно завершена.
```

Проверить существование статистических данных уровня сеанса можно, запросив представление USER_TAB_STATISTICS:

```
EODA@ORA12CR1> select table_name, num_rows, last_analyzed, scope
  1 2 from user_tab_statistics
  2 3 where table_name like 'GT';
```

| TABLE_NAME | NUM_ROWS | LAST_ANALYZED | SCOPE |
|------------|----------|---------------|---------|
| GT | | | SHARED |
| GT | 51 | 18-JAN-14 | SESSION |

Мы можем продолжить проверку осведомленности оптимизатора о закрытой статистике сеанса посредством AUTOTRACE:

```
EODA@ORA12CR1> set autotrace on;
EODA@ORA12CR1> select count(*) from gt;
```

Ближе к концу вывода находится примечание оптимизатора:

Note

```
-----
- Global temporary table session private statistics used
```

Примечание

```
-----
- Использовалась закрытые статистические данные сеанса для глобальной
  временной таблицы
```

Имейте в виду, что статистические данные уровня сеанса действительны только на протяжении сеанса. Если вы отключитесь и затем подключитесь, статистические данные исчезнут:

```
EODA@ORA12CR1> disconnect
EODA@ORA12CR1> connect eoda
Enter password:
```

Возвратившись к запросу, проверяющему существование статистических данных уровня сеанса, легко заметить, что статистика сеанса теперь отсутствует:

```
EODA@ORA12CR1> select table_name, num_rows, last_analyzed, scope
2 from user_tab_statistics
3 where table_name like 'GT';
```

| TABLE_NAME | NUM_ROWS | LAST_ANALYZED | SCOPE |
|------------|----------|---------------|--------|
| GT | | | SHARED |

На заметку! Если при запрашивании временной таблицы статистические данные уровня сеанса существуют, оптимизатор будет их использовать. Если статистические данные уровня сеанса не существуют, тогда оптимизатор проверит существование разделяемой статистики, и если она на месте, то будет работать с ней. Если отсутствует вообще любая статистика, то оптимизатор будет использовать динамическую статистику (до выхода версии Oracle 12c она называлась *динамической выборкой*).

Разделяемая статистика

Как было показано в предыдущем разделе, при генерации статистических данных для временной таблицы эта статистика видна только сеансу, который ее сгенерировал; так принято по умолчанию в Oracle 12c. Если требуется, чтобы множество сеансов совместно использовали одни и те же статистические данные для какой-то временной таблицы, вы должны сначала с помощью процедуры DBMS_STATS.SET_TABLE_PREFS установить параметр GLOBAL_TEMP_TABLE_STATS в SHARED (стандартным значением этого параметра является SESSION). Давайте в целях демонстрации создадим временную таблицу и вставим в нее данные:

```
EODA@ORA12CR1> create global temporary table gt(x number) on commit
❖preserve rows;
Table created.
Таблица создана.

EODA@ORA12CR1> insert into gt select user_id from all_users;
51 rows created.
51 строк создано.
```

Теперь установим параметр GLOBAL_TEMP_TABLE_STATS в SHARED:

```
EODA@ORA12CR1> exec dbms_stats.set_table_prefs(user, -
> 'GT', 'GLOBAL_TEMP_TABLE_STATS', 'SHARED');
```

Далее сгенерируем статистику для этой временной таблицы:

```
EODA@ORA12CR1> exec dbms_stats.gather_table_stats( user, 'GT' );
```

Выполнив следующий запрос, можно удостовериться в том, что разделяемая статистика была сгенерирована:

```
EODA@ORA12CR1> select table_name, num_rows, last_analyzed, scope
  2 from user_tab_statistics
  3 where table_name like 'GT';
```

| TABLE_NAME | NUM_ROWS | LAST_ANALYZED | SCOPE |
|------------|----------|---------------|--------|
| GT | 51 | 18-JAN-14 | SHARED |

Разделяемая статистика для глобальной временной таблицы сохраняется до тех пор, пока не будет явно удалена. Удалить ее можно так:

```
EODA@ORA12CR1> exec dbms_stats.delete_table_stats( user, 'GT' );
```

Убедиться в том, что разделяемая статистика удалена, можно посредством показанного ниже запроса:

```
EODA@ORA12CR1> select table_name, num_rows, last_analyzed, scope
  2 from user_tab_statistics
  3 where table_name like 'GT';
```

| TABLE_NAME | NUM_ROWS | LAST_ANALYZED | SCOPE |
|------------|----------|---------------|--------|
| GT | | | SHARED |

Статистика для временных таблиц ON COMMIT DELETE ROWS

Ранее упоминалось о том, что при запуске таких процедур, как GATHER_TABLE_STATS, выдается неявный оператор COMMIT. Следовательно, когда генерируются статистические данные для временных таблиц, определенных как ON COMMIT DELETE ROWS, собранная статистика отразит наличие в таблицах нуля строк (в этом случае она бесполезна, потому что вам нужна статистика по данным внутри временных таблиц до их удаления оператором COMMIT).

Начиная с версии Oracle 12c, некоторые процедуры в пакете DBMS_STATS (вроде GATHER_TABLE_STATS) больше не выдают неявный оператор COMMIT после сбора статистики для временных таблиц, определенных как ON COMMIT DELETE ROWS. Это означает, что теперь возможно собирать репрезентативные статистические данные для временных таблиц такого типа. Рассмотрим простой пример. Создадим временную таблицу с конструкцией ON COMMIT DELETE ROWS:

```
EODA@ORA12CR1> create global temporary table gt(x number) on commit delete rows;
Table created.
Таблица создана.
```

Вставим в эту таблицу некоторые данные:

```
EODA@ORA12CR1> insert into gt select user_id from all_users;
51 rows created.
51 строк создано.
```

Сгенерируем статистику для схемы:

```
EODA@ORA12CR1> exec dbms_stats.gather_table_stats( user, 'GT' );
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.
```

Быстрый подсчет показывает, что строки в таблице GT по-прежнему существуют:

```
EODA@ORA12CR1> select count(*) from gt;
```

```
  COUNT(*)
-----
        51
```

Запросив представление USER_TAB_STATISTICS, можно проверить существование статистических данных на уровне сеанса:

```
EODA@ORA12CR1> select table_name, num_rows, last_analyzed, scope
  2  from user_tab_statistics
  3  where table_name like 'GT';
```

| TABLE_NAME | NUM_ROWS | LAST_ANALYZED | SCOPE |
|------------|----------|---------------|---------|
| GT | | | SHARED |
| GT | 51 | 18-JAN-14 | SESSION |

Это позволяет генерировать полезную статистику для временных таблиц, из которых строки должны удаляться после каждой транзакции.

На заметку! В версии Oracle 12c следующие процедуры пакета DBMS_STATS больше не выдают оператор COMMIT в процессе сбора статистики для временных таблиц, определенных как ON COMMIT DELETE ROWS: GATHER_TABLE_STATS, DELETE_TABLE_STATS, DELETE_COLUMN_STATS, DELETE_INDEX_STATS, SET_TABLE_STATS, SET_COLUMN_STATS, SET_INDEX_STATS, GET_TABLE_STATS, GET_COLUMN_STATS, GET_INDEX_STATS. Перечисленные процедуры выдают неявный оператор COMMIT для временных таблиц, которые определены как ON COMMIT PRESERVE ROWS.

Автоматический сбор статистики для операций загрузки в прямом режиме

Начиная с версии Oracle 12c, при выполнении операций прямого режима в отношении временной таблицы (определенной как ON COMMIT PRESERVE ROWS) статистика уровня сеанса по умолчанию собирается для временной таблицы, в которую производится загрузка. Двумя типичными операциями загрузки в прямом режиме являются CREATE TABLE AS SELECT (CTAS) и операторы INSERT прямого режима (т.е. INSERT с подсказкой /*+ append */).

Продemonстрируем это на простом примере. Создадим таблицу CTAS:

```
EODA@ORA12CR1> create global temporary table gt on commit preserve rows
  2  as select * from all_users;
```

Table created.

Таблица создана.

Следующий запрос позволяет проверить, была ли сгенерирована статистика уровня сеанса:

```
EODA@ORA12CR1> select table_name, num_rows, last_analyzed, scope
  2  from user_tab_statistics
  3  where table_name like 'GT';
```

| TABLE_NAME | NUM_ROWS | LAST_ANALYZED | SCOPE |
|------------|----------|---------------|---------|
| GT | | | SHARED |
| GT | 51 | 18-JAN-14 | SESSION |

При загрузке в прямом режиме временной таблицы, определенной как `ON COMMIT PRESERVE ROWS`, устраняется необходимость в обращении к пакету `DBMS_STATS` для генерации статистических данных.

Заключительные соображения по поводу временных таблиц

Временные таблицы могут быть полезны в приложении, где необходимо временно сохранять набор строк, которые требуются при обработке других таблиц, на протяжении либо сеанса, либо транзакции. Они не предназначены для применения в качестве средства разбиения одного большого запроса на несколько меньших результирующих наборов, которые затем объединялись бы вместе (пожалуй, самый популярный способ использования временных таблиц в других базах данных). На самом деле в большинстве случаев вы обнаружите, что одиночный запрос, разделенный на меньшие по размеру запросы, работающие с временными таблицами, выполняется в Oracle гораздо медленнее, чем он выполнялся бы без разбиения. Я неоднократно сталкивался с таким поведением, когда получал возможность переписывать последовательность операторов `INSERT`, вставляющих данные во временные таблицы, как операторов `SELECT` в форме одного большого запроса: единственный результирующий запрос выполнялся намного быстрее исходного многошагового процесса.

Временные таблицы генерируют минимальный объем информации redo, но все же генерируют. До выхода Oracle 12c не существовало способа отключить это. Информация redo генерируется для данных отката, и в большинстве ситуаций ее будет пренебрежимо мало. Если для временных таблиц выдаются только операторы `INSERT` и `SELECT`, то объем генерируемой информации redo окажется незначительным. Крупные объемы redo будут генерироваться только в случае выполнения для временных таблиц операторов `DELETE` и `UPDATE`.

На заметку! Начиная с версии Oracle 12c, базу данных можно проинструктировать относительно записи информации undo во временное табличное пространство и тем самым устранить генерацию redo почти полностью. Это делается установкой параметра `TEMP_UNDO_ENABLED` в `TRUE` (за подробными сведениями обращайтесь в главу 9).

Соблюдая определенную осторожность, для временной таблицы можно сгенерировать статистику, применяемую оптимизатором СВО; однако лучшие статистические предположения для временной таблицы можно получить с использованием пакета `DBMS_STATS` или динамически собрать оптимизатором во время полного разбора с помощью динамической выборки. В Oracle 12c появилась возможность генерировать статистические данные, специфичные для сеанса. Это предоставит оптимизатору более точную информацию для построения планов выполнения, которые будут оптимальными в отношении данных, загруженных в отдельно взятом сеансе.

Объектные таблицы

Мы уже видели частичный пример объектной таблицы, когда исследовали вложенные таблицы. *Объектная таблица* — это таблица, которая создается на основе типа, а не коллекции столбцов.

Обычно оператор CREATE TABLE выглядит следующим образом:

```
create table t ( x int, y date, z varchar2(25) );
```

Оператор CREATE TABLE для создания объектной таблицы выглядит иначе:

```
create table t of Some_Type;
```

Атрибуты (столбцы) таблицы T выводятся из определения типа SOME_TYPE. Давайте рассмотрим краткий пример с парой типов и затем проанализируем результирующие структуры данных:

```
EODA@ORA12CR1> create or replace type address_type
  2 as object
  3 ( city      varchar2(30),
  4   street    varchar2(30),
  5   state     varchar2(2),
  6   zip       number
  7 )
  8 /
```

Type created.

Тип создан.

```
EODA@ORA12CR1> create or replace type person_type
  2 as object
  3 ( name      varchar2(30),
  4   dob       date,
  5   home_address address_type,
  6   work_address address_type
  7 )
  8 /
```

Type created.

Тип создан.

```
EODA@ORA12CR1> create table people of person_type
  2 /
```

Table created.

Таблица создана.

```
EODA@ORA12CR1> desc people
```

| Name | Null? | Type |
|--------------|-------|--------------|
| NAME | | VARCHAR2(30) |
| DOB | | DATE |
| HOME_ADDRESS | | ADDRESS_TYPE |
| WORK_ADDRESS | | ADDRESS_TYPE |

В сущности это все. Мы создали несколько определений типов, после чего можно создать таблицу нужного типа. У таблицы оказалось четыре столбца, представляющих четыре атрибута созданного типа PERSON_TYPE. Теперь можно выполнять DML-операторы в отношении этой объектной таблицы для создания и запрашивания данных:

```
EODA@ORA12CR1> insert into people values ( 'Tom', '15-mar-1965',
  2 address_type( 'Denver', '123 Main Street', 'Co', '12345' ),
  3 address_type( 'Redwood', '1 Oracle Way', 'Ca', '23456' ) );
1 row created.
1 строка создана.
```

```
EODA@ORA12CR1> select name, dob, p.home_address Home, p.work_address work
2   from people p;
```

```
Tom                15-MAR-65
ADDRESS_TYPE('Denver', '123 Main Street', 'Co', 12345)
ADDRESS_TYPE('Redwood', '1 Oracle Way', 'Ca', 23456)
```

```
EODA@ORA12CR1> select name, p.home_address.city from people p;
```

```
NAME                HOME_ADDRESS.CITY
-----
Tom                Denver
```

Мы начинаем видеть объектный синтаксис, необходимый для работы с объектными типами. Например, в операторе INSERT пришлось выполнить приведение для столбцов HOME_ADDRESS и WORK_ADDRESS — скалярные значения были приведены к типу ADDRESS_TYPE. По-другому можно было бы сказать, что для этой строки мы создали экземпляр объекта ADDRESS_TYPE с применением его стандартного конструктора.

Теперь, если судить о внешней стороне таблицы, то в ней есть четыре столбца. После того как мы наблюдали “магию”, скрытую во вложенных таблицах, можно предположить, что здесь происходит что-то еще. Все объектно-реляционные данные Oracle хранит в обычных реляционных таблицах — к концу дня все они заполнены строками и столбцами. Заглянув в реальный словарь данных, можно увидеть, как в действительности выглядит наша таблица:

```
EODA@ORA12CR1> select name, segcollength
2   from sys.col$
3   where obj# = ( select object_id
4                  from user_objects
5                  where object_name = 'PEOPLE' )
6   /
```

| NAME | SEGCOLLENGTH |
|------------------|--------------|
| SYS_NC_OID\$ | 16 |
| SYS_NC_ROWINFO\$ | 1 |
| NAME | 30 |
| DOB | 7 |
| HOME_ADDRESS | 1 |
| SYS_NC00006\$ | 30 |
| SYS_NC00007\$ | 30 |
| SYS_NC00008\$ | 2 |
| SYS_NC00009\$ | 22 |
| WORK_ADDRESS | 1 |
| SYS_NC00011\$ | 30 |
| SYS_NC00012\$ | 30 |
| SYS_NC00013\$ | 2 |
| SYS_NC00014\$ | 22 |

```
14 rows selected.
14 строк выбрано.
```


Это существенно отличается от того, что сообщала команда DESCRIBE. Очевидно, что таблица имеет не четыре, а четырнадцать столбцов, которые описаны ниже.

- SYS_NC_OID\$. Это сгенерированный системой объектный идентификатор таблицы. Представляет собой уникальный столбец типа RAW(16). Он имеет ограничение уникальности и на нем также создан соответствующий уникальный индекс.
- SYS_NC_ROWINFO\$. Это та же самая “магическая” функция, которую мы наблюдали во вложенных таблицах. Выбор данного столбца из таблицы приводит к возвращению всей строки как единственного столбца:

```
EODA@ORA12CR1> select sys_nc_rowinfo$ from people;
```

```
SYS_NC_ROWINFO$(NAME, DOB, HOME_ADDRESS(CITY, STREET, STATE, ZIP),  
WORK_ADDRESS(CITY, STREET, STATE,
```

```
-----  
PERSON_TYPE('Tom', '15-MAR-65', ADDRESS_TYPE('Denver', '123 Main Street',  
'Co', 12345), ADDRESS_TYPE('Redwood', '1 Oracle Way', 'Ca', 23456))
```

- NAME, DOB. Это скалярные атрибуты нашей объектной таблицы. Как и ожидалось, они хранятся в виде обычных столбцов.
- HOME_ADDRESS, WORK_ADDRESS. Это также “магические” функции. Они возвращают коллекцию столбцов, которые представляют одиночный объект. Они не потребляют реальное пространство за исключением указания NULL или NOT NULL для сущности.
- SYS_NCnnnnnn\$. Это скалярные реализации наших встроенных объектных типов. Поскольку тип PERSON_TYPE имеет встроенный тип ADDRESS_TYPE, требуется предусмотреть место для сохранения такого объекта в столбцах соответствующего типа. Сгенерированные системой имена являются обязательными, т.к. имя столбца должно быть уникальным, а один и тот же объектный тип может использоваться более одного раза, что и было сделано. Если бы эти имена не генерировались, столбец ZIP встретился бы дважды.

Таким образом, как и во вложенной таблице, здесь много чего происходит. Был добавлен псевдо-первичный ключ размером 16 байтов, появились виртуальные столбцы и создан индекс. Стандартное поведение, касающееся назначаемого объекту идентификатора, можно изменить, и вскоре мы посмотрим как именно. Сначала давайте взглянем на полный SQL-код, который генерирует нашу таблицу. Код был сформирован с помощью Data Pump, т.к. необходимо увидеть все зависимые объекты, включая весь SQL-код, который предназначен для воссоздания этого экземпляра объекта. Для этого выполнены следующие действия:

```
$ expdp eoda directory=tk tables='PEOPLE' dumpfile=p.dmp logfile=p.log  
Export: Release 12.1.0.1.0 - Production on Sat Jan 18 17:10:11 2014  
Copyright (c) 1982, 2013, Oracle and/or its affiliates. All rights reserved.  
...  
$ impdp eoda directory=tk dumpfile=p.dmp logfile=pi.log sqlfile=people.sql  
Import: Release 12.1.0.1.0 - Production on Sat Jan 18 17:11:54 2014  
...
```

Connected to: Oracle Database 12c Enterprise Edition Release 12.1.0.1.0 -
 64bit Production
 With the Partitioning, OLAP, Advanced Analytics and Real Application
 Testing options
 Master table "EODA"."SYS_SQL_FILE_FULL_01" successfully loaded/unloaded
 Starting "EODA"."SYS_SQL_FILE_FULL_01":
 eoda/***** directory=tk dumpfile=p.dmp logfile=pi.log sqlfile=people.sql

Ниже показано содержимое результирующего файла people.sql:

```
-- new object type path: TABLE_EXPORT/TABLE/TABLE
-- путь к новому объектному типу: TABLE_EXPORT/TABLE/TABLE
CREATE TABLE "EODA"."PEOPLE" OF "EODA"."PERSON_TYPE"
OID 'F0484A73A93A7093E043B7D04F0A821B'
OIDINDEX ( PCTFREE 10 INITRANS 2 MAXTRANS 255
  STORAGE(INITIAL 65536 NEXT 1048576 MINEXTENTS 1 MAXEXTENTS 2147483645
    PCTINCREASE 0 FREELISTS 1 FREELIST GROUPS 1
    BUFFER_POOL DEFAULT FLASH_CACHE DEFAULT CELL_FLASH_CACHE DEFAULT)
  TABLESPACE "USERS" )
PCTFREE 10 PCTUSED 40 INITRANS 1 MAXTRANS 255
NOCOMPRESS LOGGING
  STORAGE(INITIAL 65536 NEXT 1048576 MINEXTENTS 1 MAXEXTENTS 2147483645
    PCTINCREASE 0 FREELISTS 1 FREELIST GROUPS 1
    BUFFER_POOL DEFAULT FLASH_CACHE DEFAULT CELL_FLASH_CACHE DEFAULT)
  TABLESPACE "USERS" ;
```

Приведенный код дает чуть больше информации о том, что здесь происходит на самом деле. Теперь мы ясно видим конструкцию `OIDINDEX` и ссылку на столбец `OID`, на которой следует шестнадцатеричное число.

Синтаксис `OID '<большое шестнадцатеричное число>'` в документации Oracle не описан. Он всего лишь гарантирует, что во время последовательного выполнения `expdp` и `impdp` лежащий в основе тип `PERSON_TYPE` является действительно тем же самым. Это предотвратит ошибку, которая может произойти при выполнении следующих шагов.

1. Создание таблицы `PEOPLE`.
2. Экспортирование этой таблицы.
3. Удаление этой таблицы и лежащего в основе типа `PERSON_TYPE`.
4. Создание нового типа `PERSON_TYPE` с другими атрибутами.
5. Импорт данных старой таблицы `PEOPLE`.

Очевидно, что такие экспортированные данные не могут быть импортированы в новую структуру — они не подойдут. Это предотвращается благодаря конструкции `OID '<большое шестнадцатеричное число>'`.

Если вы помните, я ранее упоминал, что поведение объектного идентификатора, назначаемого экземпляру объекта, можно изменить. Вместо псевдо-первичного ключа, генерируемого системой, можно применять настоящий ключ объекта. На первый взгляд это может показаться обреченным на провал — столбец `SYS_NC_OID$` по-прежнему будет присутствовать в определении таблицы в `SYS.COL$` и на самом деле покажется, что он потребляет намного больше пространства, чем столбец, генерируемый системой. Однако здесь вновь в игру вступает “магия”.

Столбец SYS_NC_OID\$ объектной таблицы, который основан на *первичном ключе* и не сгенерирован *системой*, является виртуальным столбцом и не занимает реального пространства на диске.

Ниже приведен пример, который демонстрирует происходящее в словаре данных и показывает, что никакого физического пространства столбец SYS_NC_OID\$ не использует. Мы начнем с анализа таблицы с объектным идентификатором, генерируемым системой:

```
EODA@ORA12CR1> create table people of person_type
2 /
Table created.
Таблица создана.
```

```
EODA@ORA12CR1> select name, type#, segcollength
2   from sys.col$
3   where obj# = ( select object_id
4                   from user_objects
5                   where object_name = 'PEOPLE' )
6   and name like 'SYS\NC\_%' escape '\
7 /
```

| NAME | TYPE# | SEGCOLLENGTH |
|------------------|-------|--------------|
| SYS_NC_OID\$ | 23 | 16 |
| SYS_NC_ROWINFO\$ | 121 | 1 |

```
EODA@ORA12CR1> insert into people(name)
2   select rownum from all_objects;
72069 rows created.
72069 строк создано.
```

```
EODA@ORA12CR1> exec dbms_stats.gather_table_stats( user, 'PEOPLE' );
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.
```

```
EODA@ORA12CR1> select table_name, avg_row_len from user_object_tables;
```

| TABLE_NAME | AVG_ROW_LEN |
|------------|-------------|
| PEOPLE | 24 |

Здесь мы видим, что средняя длина строк составляет 24 байта: 16 байтов занимает столбец SYS_NC_OID\$ и 8 байтов — столбец NAME. Теперь давайте сделаем то же самое, но применим в качестве объектного идентификатора первичный ключ на столбце NAME:

```
EODA@ORA12CR1> CREATE TABLE "PEOPLE"
2   OF "PERSON_TYPE"
3   ( constraint people_pk primary key(name) )
4   object identifier is PRIMARY KEY
5 /
Table created.
Таблица создана.
```

```

EODA@ORA12CR1> select name, type#, segcollength
2   from sys.col$
3   where obj# = ( select object_id
4                   from user_objects
5                   where object_name = 'PEOPLE' )
6   and name like 'SYS\_NC\_%' escape '\ '
7 /

```

| NAME | TYPE# | SEGCOLLENGTH |
|------------------|-------|--------------|
| SYS_NC_OID\$ | 23 | 81 |
| SYS_NC_ROWINFO\$ | 121 | 1 |

Согласно этому выводу, вместо небольшого 16-байтового столбца мы имеем крупный столбец размером 81 байт! В реальности никакие данные в нем не хранятся. Столбец будет иметь значение NULL. Система сгенерирует уникальный идентификатор на основе объектной таблицы, лежащего в ее основе типа и значения в самой строке. В этом можно удостовериться следующим образом:

```

EODA@ORA12CR1> insert into people (name) values ( 'Hello World!' );
1 row created.
1 строка создана.

```

```

EODA@ORA12CR1> select sys_nc_oid$ from people p;

```

```

SYS_NC_OID$

```

```

F04931FE974478A7E043B7D04F0A0820000000172601000100010029000000000000C070
01E0100002A00078401FE000000140C48656C6C6F20576F726C64210000000000000000
00000000000000000000

```

```

EODA@ORA12CR1> select utl_raw.cast_to_raw( 'Hello World!' ) data from dual;

```

```

DATA

```

```

48656C6C6F20576F726C6421

```

```

EODA@ORA12CR1> select utl_raw.cast_to_varchar2(sys_nc_oid$) data from people;

```

```

DATA

```

```

<ненужные данные...>Hello World!

```

Если мы выберем столбец SYS_NC_OID\$ и проинспектируем шестнадцатеричный дамп вставленной строки, то увидим, что сами данные строки встроены в объектный идентификатор. Преобразовав объектный идентификатор в значение типа VARCHAR2, мы можем получить визуальное подтверждение этого. Означает ли это, что наши данные сохранены дважды, порождая огромные накладные расходы? Нет, они просто учитываются той “магической” функцией, которая вступает в действие при извлечении столбца SYS_NC_OID\$. Данные синтезируются Oracle во время выборки из таблицы.

Теперь я изложу свое мнение. Я часто называю объектно-реляционные компоненты (вложенные и объектные таблицы) *синтаксическим сахаром*. Их всегда можно

транслировать в старые добрые реляционные строки и столбцы. Лично я предпочитаю не использовать их как физические механизмы хранения. Слишком много в них происходит “магических” действий с неясными побочными эффектами. Вы получаете скрытые столбцы, дополнительные индексы, неожиданные псевдостолбцы и т.п. *Это вовсе не означает, что применение объектно-реляционных компонентов является бесполезной тратой времени.* Напротив, я постоянно использую их в коде PL/SQL. Я применяю их с объектными представлениями. Я могу получить преимущества конструкции вложенных таблиц (меньший объем данных, передаваемых по сети, для отношения “главный-подчиненный”, концептуально более простая работа с данными и т.д.) безо всяких проблем, касающихся физического хранения. Причина в том, что с помощью объектных представлений я могу синтезировать свои объекты из реляционных данных. Это решает большинство проблем объектных и вложенных таблиц: я самостоятельно выбираю способ физического хранения, сам настраиваю условия соединения, а таблицы естественным образом доступны как реляционные (чего требуют многие сторонние инструменты и приложения). Те, кому нужно объектное представление реляционных данных, могут получить его, а те, кого интересует реляционное представление, могут работать с ним. Поскольку объектные таблицы в действительности являются замаскированными реляционными таблицами, мы делаем то же самое, что Oracle выполняет “за кулисами”, только более эффективно, т.к. мы не обязаны делать это обобщенным образом, как поступает Oracle. Например, имея определенные ранее типы, можно легко написать следующий код:

```

EODA@ORA12CR1> create table people_tab
 2  (   name          varchar2(30) primary key,
 3      dob           date,
 4      home_city     varchar2(30),
 5      home_street   varchar2(30),
 6      home_state    varchar2(2),
 7      home_zip      number,
 8      work_city     varchar2(30),
 9      work_street   varchar2(30),
10      work_state    varchar2(2),
11      work_zip      number
12  )
13  /

```

Table created.

Таблица создана.

```

EODA@ORA12CR1> create view people of person_type
 2  with object identifier (name)
 3  as
 4  select name, dob,
 5      address_type(home_city,home_street,home_state,home_zip) home_adress,
 6      address_type(work_city,work_street,work_state,work_zip) work_adress
 7      from people_tab
 8  /

```

View created.

Представление создано.

```

EODA@ORA12CR1> insert into people values ( 'Tom', '15-mar-1965',
 2  address_type( 'Denver', '123 Main Street', 'Co', '12345' ),

```

```
3 address_type( 'Redwood', '1 Oracle Way', 'Ca', '23456' );
1 row created.
1 строка создана.
```

Несмотря на то что я добился практически того же самого результата, мне точно известно, что сохраняется, как сохраняется и где сохраняется. Для более сложных объектов можно закодировать в объектных представлениях триггеры `INSTEAD OF`, чтобы разрешить внесение изменений посредством представления.

Заключительные соображения по поводу объектных таблиц

Объектные таблицы используются для реализации объектно-реляционной модели в Oracle. Одиночная объектная таблица обычно будет приводить к созданию множества физических объектов в базе данных и добавлению в схему дополнительных столбцов, предназначенных для управления ими. С объектными таблицами связана определенная “магия”. Объектные представления позволяют получить преимущества от применения синтаксиса и семантики объектов, одновременно обладая полным контролем над физическим хранением данных и позволяя реляционный доступ к лежащим в основе данным. Поступая подобным образом, вы можете извлечь лучшее из реляционного и объектно-реляционного миров.

Резюме

После прочтения этой главы вы должны были прийти к заключению, что не все таблицы создаются одинаково. В Oracle предлагается широкое разнообразие типов таблиц, которые вы можете эксплуатировать. В главе были рассмотрены многие самые заметные аспекты таблиц в целом и проанализированы различные типы таблиц, предоставляемые Oracle.

Сначала была описана терминология и параметры хранения, которые ассоциированы с таблицами. Было показано, насколько полезными являются списки свободных блоков в многопользовательской среде, где таблица часто подвергается вставкам и обновлениям со стороны множества пользователей одновременно, и как с помощью табличных пространств ASSM можно сделать так, чтобы даже не думать об этом. Вы узнали, что собой представляют параметры `PCTFREE` и `PCTUSED`, и ознакомились с некоторыми рекомендациями по их корректной установке.

После этого мы занялись исследованием различных типов таблиц, начав с традиционных. Традиционная таблица пока что остается наиболее часто используемым в приложениях Oracle табличным типом, а также типом, который выбирается по умолчанию. Затем мы перешли к описанию индекс-таблиц, которые позволяют хранить табличные данные в индексной структуре. Вы увидели, что они применимы во многих сценариях, таких как справочные таблицы и инвертированные списки, где традиционные таблицы оказались бы просто избыточной копией данных. Далее было показано, насколько полезными могут быть индекс-таблицы при их использовании вместе с таблицами других типов, особенно в случае вложенных таблиц.

Мы взглянули на кластерные объекты, которых в Oracle есть три вида: индексные кластеры, хеш-кластеры и отсортированные хеш-кластеры. Цель кластеров двояка.

- Предоставить возможность хранения данных из многих таблиц вместе в одном и том же блоке (блоках) базы данных.

- Предоставить возможность принудительного сохранения похожих данных физически вместе на основе некоторого кластерного ключа. Это позволяет, например, хранить вместе все данные для отдела 10 (из множества таблиц).

Такие возможности позволяют получать доступ к связанным данным очень быстро и с минимальным объемом ввода-вывода. Мы обсудили главные отличия между индексными кластерами и хеш-кластерами, а также ситуации, в которых каждый из них подходит (или не подходит).

Следующими рассматривались вложенные таблицы. Вы ознакомились с синтаксисом, семантикой и применением вложенных таблиц. Вы увидели, что на самом деле они представляют собой генерируемые и обслуживаемые системой пары родительских и дочерних таблиц, и узнали, каким образом Oracle делает это физически. Было показано, что в качестве вложенных таблиц можно использовать разные типы таблиц, а не только традиционные, которые применяются по умолчанию. Как обнаружилось, почти во всех случаях для вложенных таблиц лучше использовать не традиционные, а индекс-таблицы.

Затем было дано описание временных таблиц: вы узнали, как они создаются, где получают пространство для хранения, а также то, что во время выполнения они не создают никаких проблем, связанных с параллелизмом. Были рассмотрены отличия между временными таблицами уровня сеанса и временными таблицами уровня транзакции, а также предложен подходящий метод применения временных таблиц в базе данных Oracle.

Глава завершилась исследованием объектных таблиц. Вы узнали, что в них, как и во вложенных таблицах, происходит много скрытых действий. Мы обсудили, как с помощью объектных представлений на реляционных таблицах можно получить функциональность объектной таблицы и в то же время обеспечить простой доступ к лежащим в основе реляционным данным.

глава 11

Индексы

Индексация является важнейшим аспектом при проектировании и разработке приложения. Если индексов слишком много, от этого страдает производительность операций модификации (вставок, обновлений, слияний и удалений). Когда индексов слишком мало, ухудшается производительность операций DML (в том числе операций выборки, вставки, обновления и удаления). Нахождение правильного сочетания индексов критически важно для производительности приложения.

Я нередко обнаруживаю, что во время разработки приложений индексы добавляются с опозданием. Считаю это ошибочным подходом. Если вы понимаете, каким образом будут использоваться данные, с самого начала процесса, то должны быть в состоянии получить репрезентативный набор индексов, которые будут применяться в приложении. Однако очень часто все выглядит так, что приложение сначала создают, а потом смотрят, где необходимы индексы. Это говорит о том, что вы не потрудились понять, как будут использоваться данные и сколько строк в конечном итоге потребуется. Вы будете постоянно добавлять индексы в приложение по мере роста объема данных с течением времени (т.е. выполнять *реактивную настройку*). Вы получите индексы, которые являются избыточными и никогда не задействуются; это приводит к напрасному расходованию не только дискового пространства, но и вычислительных ресурсов. Несколько часов, потраченных в самом начале на обдумывание того, когда и как индексировать данные, сэкономят многие часы, которые придется посвятить настройке в процессе эксплуатации (обратите внимание: *экономят* многие часы, а не *могут* сэкономить).

Главная цель настоящей главы заключается в том, чтобы предоставить обзор индексов, доступных для применения в Oracle, и обсудить, когда и где их можно использовать. Эта глава отличается от остальных глав в книге стилем и форматом. Индексация — огромная тема, и о ней можно написать целую книгу, отчасти потому, что индексация связывает вместе роли разработчика и администратора базы данных. Разработчик должен быть осведомленным об индексах, о том, как их применять в приложениях, когда их использовать (а когда нет) и т.д. Администратора базы данных заботит рост индекса, расходование хранилища под индекс, другие его физические свойства и общая производительность базы данных. Мы рассмотрим индексы главным образом с точки зрения их практического применения в приложениях. В первой половине главы предлагаются базовые сведения, которые вам понадобятся для принятия обоснованных решений о том, когда индексировать и какой тип индекса использовать. Во второй половине главы вы найдете ответы на некоторые часто задаваемые вопросы относительно индексов.

Разнообразные примеры, приведенные в настоящей главе, требуют разных выпусков Oracle. Когда в конкретном примере необходимы средства из редакции Enterprise или Personal, а не Standard, на это будет указано специально.

Обзор индексов Oracle

В Oracle предоставляется много разных типов индексов. Ниже приведен их краткий перечень.

- **Индексы со структурой В-дерева (B*Tree index).** Это то, что я называю обычными индексами. Они являются самыми распространенными индексами в Oracle и большинстве других баз данных. Подобный по конструкции двоичному дереву, индекс со структурой В-дерева предоставляет быстрый доступ по ключу к индивидуальной строке или диапазону строк, обычно требуя нескольких операций чтения для нахождения нужной строки. Однако важно отметить, что В в термине В-дерево означает не binary (двоичное), а balanced (сбалансированное). Индекс со структурой В-дерева — это вообще не двоичное дерево, если взглянуть на то, как он физически хранится на диске. Различают несколько подтипов индекса со структурой В-дерева.
 - **Индекс-таблицы (Index Organized Table — IOT).** Это таблицы, хранящиеся в структуре В-дерева. В то время как строки данных в традиционной таблице хранятся в неорганизованной манере (данные отправляются туда, где есть свободное место), данные в индекс-таблице сохраняются и сортируются по первичному ключу. С точки зрения приложения индекс-таблицы ведут себя как “обычные” таблицы: для доступа к ним применяется язык SQL, как происходит в нормальной ситуации. Индекс-таблицы особенно удобны в информационно-поисковых, пространственных и OLAP-приложениях. Мы довольно подробно обсуждали индекс-таблицы в главе 10.
 - **Кластерные индексы со структурой В-дерева (B*Tree cluster index).** Это небольшая вариация обычных индексов со структурой В-дерева. Они используются для индексации кластерных ключей (см. раздел “Кластеризованные индекс-таблицы” в главе 10) и здесь повторно не рассматриваются. Вместо наличия ключа, который указывает на строку, как в обычных индексах со структурой В-дерева, в таких индексах есть кластерный ключ, указывающий на блок, который содержит строки, относящиеся к этому кластерному ключу.
 - **Индексы, упорядоченные по убыванию (descending index).** Эти индексы позволяют данным быть отсортированными в порядке от больших значений к меньшим (по убыванию) вместо порядка от меньших значений к большим (по возрастанию), принятого в индексной структуре. Мы посмотрим, зачем они нужны и как работают.
 - **Индексы по реверсированным ключам (reverse key index).** Это индексы со структурой В-дерева, в которых порядок следования байтов ключа изменен на противоположный. Индексы по реверсированным ключам могут применяться для обеспечения более равномерного распределения записей в индексах, которые наполняются увеличивающимися значениями. Например, если для генерации первичного ключа используется последовательность, то будут по-

лучены значения, подобные 987500, 987501, 987502 и т.д. Значения последовательности являются монотонными, поэтому в случае применения обычного индекса со структурой В-дерева они имеют тенденцию поступать в один и тот же правосторонний блок, в итоге увеличивая конкуренцию за данный блок. В индексе с реверсированным ключом вместо этого Oracle будет логически индексировать значения 205789, 105789, 005789 и т.д. Порядок байтов данных перед помещением в индекс изменяется на противоположный, так что значения, которые до этого находились бы в индексе рядом, окажутся на большом расстоянии друг от друга. Реверсирование байтов индекса рассредоточивает значения, вставляемые в индекс, по множеству блоков.

- **Битовые индексы (bitmap index).** Обычно в В-дереве существует отношение “один к одному” между записью индекса и строкой: запись индекса указывает на строку. В битовых индексах запись индекса использует битовую карту для указания на множество строк одновременно. Они подходят для данных с высокой повторяемостью (данных с небольшим числом отличающихся значений по сравнению с общим количеством строк в таблице), которые доступны в основном только для чтения. Рассмотрим столбец, который принимает три возможных значения — Y, N и NULL — в таблице с миллионом строк. Этот столбец может быть подходящим кандидатом для построения битового индекса, если, к примеру, нужно часто определять количество строк, имеющих в этом столбце значение Y. Речь вовсе не о том, что битовый индекс по столбцу с 1000 разных значений в той же самой таблице будет недопустимым — он вполне имеет право на существование. Битовые индексы не должны применяться в базе данных OLTP для решения проблем, связанных с параллелизмом (мы обсудим это позже). Обратите внимание, что битовые индексы требуют редакции Enterprise или Personal системы Oracle.
- **Битовые индексы соединений (bitmap join index).** Эти индексы предоставляют средство денормализации данных в индексной структуре, а не в таблице. Например, рассмотрим простые таблицы EMP и DEPT. У кого-то может возникнуть вопрос: “Сколько людей работает в подразделениях, расположенных в Бостоне?”. Таблица EMP имеет внешний ключ для DEPT, и чтобы подсчитать сотрудников в подразделениях, у которых значение LOC равно Boston (Бостон), обычно приходится соединять таблицы, чтобы получить столбец LOC, соединенный с записями EMP. Используя битовый индекс соединения, можно вместо этого проиндексировать столбец LOC по таблице EMP. То же самое предупреждение, касающееся применения обычного битового индекса в OLTP-системах, относится и к битовому индексу соединения.
- **Индексы на основе функций (function-based index).** Это индексы со структурой В-дерева или битовые индексы, которые хранят вычисленный результат функции на столбце (столбцах) строки, а не сами данные столбца. Их можно трактовать как индексы на виртуальном (или производном) столбце — другими словами, на столбце, который не хранится физически в таблице. Такие индексы могут использоваться для ускорения запросов в форме `SELECT * FROM T WHERE FUNCTION(DATABASE_COLUMN)=SOME_VALUE`, поскольку значение `FUNCTION(DATABASE_COLUMN)` уже вычислено и сохранено в индексе.

- **Индексы предметной области (application domain index).** Это индексы, которые вы строите и сохраняете самостоятельно — либо в Oracle, либо даже за пределами Oracle. Вы сообщаете оптимизатору, насколько селективным является индекс, и насколько дорогостоящим в смысле ресурсов будет его привлечение, а оптимизатор на основе этой информации решает, применять его или нет. Примером индекса предметной области можно считать текстовый индекс Oracle; он создается с помощью тех же инструментов, которые используются при построении собственного индекса. Следует отметить, что индекс, созданный подобным образом, не нуждается в применении традиционной индексной структуры. Скажем, текстовый индекс Oracle для реализации своей концепции индекса использует набор таблиц.

Как видите, на выбор предлагается много типов индексов. В последующих разделах будут представлены некоторые технические детали относительно того, как работает и когда должен применяться каждый из них. Я хотел бы снова подчеркнуть, что здесь не будут раскрываться определенные темы, имеющие отношение к работе администратора базы данных. Например, мы не будем обсуждать механизмы оперативной перестройки индексов; вместо этого мы сосредоточим внимание на практических деталях, связанных с приложениями.

Индексы со структурой В-дерева

Индексы со структурой В-дерева, или то, что я называю обычными индексами, являются наиболее часто используемым типом индексных структур в базе данных. Они похожи в реализации на двоичное дерево поиска. Их цель заключается в том, чтобы минимизировать время, затрачиваемое Oracle на поиск данных. Говоря упрощенно, если есть индекс на числовом столбце, то структура концептуально может выглядеть так, как показано на рис. 11.1.

На заметку! На уровне блоков существуют оптимизация и сжатие данных, которые приводят к тому, что реальная структура блоков будет отличаться от представленной на рис. 11.1. Кроме того, на рис. 11.1 изображен неunikальный индекс (т.е. он допускает наличие дублированных ключей). Например, если вы хотите найти значение 11, то с таким значением в индексе обнаружатся две разных записи.

Блоки самого нижнего уровня в дереве, называемые *листовыми узлами* (leaf node) или *листовыми блоками* (leaf block), содержат все индексированные ключи и идентификаторы строк (rowid), которые указывают на индексируемые строки. Промежуточные блоки, находящиеся над листовыми, называются *блоками ветвления* (branch block). Они применяются для навигации по структуре. Например, если мы хотим найти в индексе значение 42, то должны начать с вершины дерева и двигаться влево. Мы должны просмотреть этот блок и выяснить, что необходимо двигаться к блоку в диапазоне 42..50. Этот блок должен быть листовым и привести к строкам, содержащим число 42.

Интересно отметить, что листовые узлы индекса в действительности представляют собой двухсвязные списки.

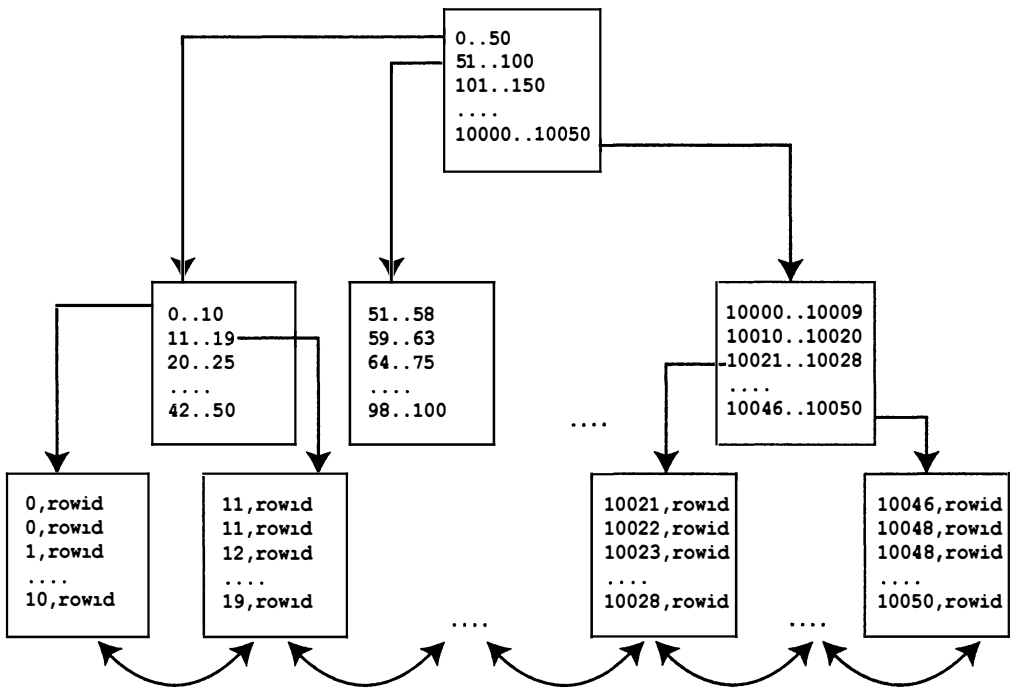


Рис. 11.1. Типичное устройство индекса со структурой В-дерева

Как только мы находим, где начинаются листовые узлы (т.е. находим первое значение), выполнение упорядоченного просмотра значений (также известного под названием *сканирование диапазона по индексу* (index range scan)) становится очень простым. Навигация по структуре больше не понадобится; мы просто по мере необходимости перемещаемся вперед и назад по листовым узлам. Это позволяет довольно просто обеспечить удовлетворение предиката вроде следующего:

where x between 20 and 30

База данных Oracle находит первый листовой блок индекса, содержащий наименьшее значение ключа, которое равно 20 или больше, а затем просто движется горизонтально через связанный список листовых узлов до тех пор, пока в итоге не попадет на значение больше 30.

В индексе со структурой В-дерева на самом деле не бывает таких вещей, как неуникальные записи. В неуникальном индексе Oracle просто сохраняет идентификатор строки, добавляя его к ключу в качестве дополнительного столбца с такой длиной в байтах, чтобы сделать ключ уникальным. Например, индекс наподобие `CREATE INDEX I ON T (X, Y)` — это концептуально `CREATE UNIQUE INDEX I ON T (X, Y, ROWID)`. В уникальном индексе, определяемом вами, Oracle не добавляет идентификатор строки к ключу индекса. В неуникальном индексе вы обнаружите, что данные отсортированы сначала по значениям ключа индекса (в порядке ключа индекса), а затем по возрастанию идентификаторов строк. В уникальном индексе данные отсортированы только по значениям ключа индекса.

Одно из свойств В-дерева связано с тем, что листовые блоки должны находиться на одном и том же уровне в дереве. Этот уровень также известен как *высота* индекса — в том смысле, что при любом обходе от корневого блока индекса к его листовому блоку должно быть посещено одинаковое количество блоков. То есть, чтобы попасть в листовой блок с целью извлечения первой строки для запроса в форме `SELECT INDEXED_COL FROM T WHERE INDEXED_COL=:X`, понадобится одно и то же число операций ввода-вывода независимо от используемого значения :X. Другими словами, индекс *сбалансирован по высоте*. Большинство индексов со структурой В-дерева будут иметь высоту 2 или 3 даже для миллионов записей. Это значит, что в общем случае для нахождения ключа в индексе потребуются две или три операции ввода-вывода, т.е. совсем неплохой результат.

На заметку! Для обозначения количества блоков, вовлеченных в обход индекса от корневого до листового блока, в Oracle применяются два термина со слегка отличающимся смыслом. Первый из них — `HEIGHT` — это число блоков, через которые нужно пройти на пути от корневого блока к листовому. Значение `HEIGHT` можно отыскать в представлении `INDEX_STATS` после того, как индекс был проанализирован посредством команды `ANALYZE INDEX <имя> VALIDATE STRUCTURE`. Второй термин — `BLEVEL` — представляет собой количество уровней ветвления и отличается от `HEIGHT` на единицу (он не учитывает листовые блоки). Значение `BLEVEL` можно найти в обычных таблицах словаря, таких как `USER_INDEXES`, после сбора статистики.

Например, предположим, что есть таблица с 10 000 000 строк (подробные сведения о создании таблицы `BIG_TABLE` приведены в разделе “Настройка среды” в начале книги), которая имеет индекс по первичному ключу на числовом столбце:

```

EODA@ORA12CR1> select index_name, blevel, num_rows
2   from user_indexes
3  where table_name = 'BIG_TABLE';

```

| INDEX_NAME | BLEVEL | NUM_ROWS |
|--------------|--------|----------|
| ----- | ----- | ----- |
| BIG_TABLE_PK | 2 | 9848991 |

Значение `BLEVEL` равно 2, т.е. `HEIGHT` равно 3, поэтому поиск листа потребует двух операций ввода-вывода (а его выдача — третьей операции ввода-вывода). Таким образом, мы должны ожидать трех операций ввода-вывода для извлечения любого значения ключа из этого индекса:

```

EODA@ORA12CR1> set autotrace on
EODA@ORA12CR1> select id from big_table where id = 42;

```

Execution Plan

| ----- | | | | | | | |
|-------|-------------------|--------------|------|-------|-------------|----------|--|
| Id | Operation | Name | Rows | Bytes | Cost (%CPU) | Time | |
| ----- | | | | | | | |
| 0 | SELECT STATEMENT | | 1 | 6 | 2 (0) | 00:00:01 | |
| * 1 | INDEX UNIQUE SCAN | BIG_TABLE_PK | 1 | 6 | 2 (0) | 00:00:01 | |

Statistics

```
...      3 consistent gets
...      1 rows processed
```

```
EODA@ORA12CR1> select id from big_table where id = 12345;
```

Statistics

```
...      3 consistent gets
...      1 rows processed
```

```
EODA@ORA12CR1> select id from big_table where id = 1234567;
```

Statistics

```
...      3 consistent gets
...      1 rows processed
```

В-дерево представляется собой великолепный механизм индексации общего назначения, который хорошо работает с большими и малыми таблицами, а также характеризуется незначительным (если вообще каким-либо) снижением производительности при росте размера таблицы.

Сжатие ключей индекса

Одной из интересных вещей, которые можно делать с индексом со структурой В-дерева, является его *сжатие*. Это не сжатие, применяемое в файлах ZIP; скорее, это сжатие, которое устраняет избыточность в сцепленных (многостолбцовых) индексах.

Тема сжатых индексов уже затрагивалась в разделе “Индекс-таблицы” главы 10, и здесь мы снова вернемся к ней. Базовая концепция, лежащая в основе индекса со сжатым ключом, состоит в том, что каждая запись индекса разбивается на две части: *префиксный* и *суффиксный* компоненты. Префиксный компонент строится по головным столбцам сцепленного индекса и имеет много повторяющихся значений. Суффиксный компонент строится по хвостовым столбцам индексного ключа и является уникальным компонентом записи индекса внутри префикса.

В качестве примера мы создадим таблицу и сцепленный индекс, после чего изменим пространство, занимаемое им без сжатия, с использованием команды ANALYZE INDEX.

На заметку! Существует распространенное заблуждение, что команда ANALYZE применяться не должна — ее заменяет пакет DBMS_STATS. Это неправда. На самом деле команда ANALYZE не должна использоваться для сбора статистики, но другие возможности ANALYZE по-прежнему востребованы. Команда ANALYZE должна применяться для выполнения таких операций, как проверка структуры индекса (как будет сделано позже) или вывода списка расщепленных строк в таблице. Пакет DBMS_STATS должен использоваться исключительно для сбора статистики по объектам.

Далее мы пересоздадим индекс со сжатием ключа, сжимая различное количество записей ключа, чтобы увидеть разницу. Начнем со следующей таблицы и индекса:

```

EODA@ORA12CR1> create table t
2 as
3 select * from all_objects
4 where rownum <= 50000;

```

Table created.

Таблица создана.

```

EODA@ORA12CR1> create index t_idx on
2 t(owner,object_type,object_name);

```

Index created.

Индекс создан.

```

EODA@ORA12CR1> analyze index t_idx validate structure;

```

Index analyzed.

Индекс проанализирован.

Затем создадим таблицу IDX_STATS, в которой будет храниться информация INDEX_STATS, и пометим строки в этой таблице как noncompressed (несжатые):

```

EODA@ORA12CR1> create table idx_stats
2 as
3 select 'noncompressed' what, a.*
4 from index_stats a;

```

Table created.

Таблица создана.

Теперь можно заметить, что компонент OWNER повторяется многократно, а это значит, что единственный блок этого индекса будет иметь десятки записей (рис. 11.2).

```

Sys.Package.Dbms_Alert
Sys.Package.Dbms_Application_Info
Sys.Package.Dbms_Aq
Sys.Package.Dbms_Aqadm
Sys.Package.Dbms_Aqadm_Sys
Sys.Package.Dbms_Aqadm_Syscalls
Sys.Package.Dbms_Aqin
Sys.Package.Dbms_Aqjms
. . .

```

Рис. 11.2. Блок индекса с повторяющимся столбцом OWNER

Повторяющийся столбец OWNER можно факторизовать, получив в результате блок, который выглядит, как показано на рис. 11.3.

```

Sys
Package.Dbms_Alert
Package.Dbms_Application_Info
Package.Dbms_Aq
Package.Dbms_Aqadm
Package.Dbms_Aqadm_Sys
Package.Dbms_Aqadm_Syscalls
Package.Dbms_Aqin
Package.Dbms_Aqjms
. . .

```

Рис. 11.3. Блок индекса с факторизованным столбцом OWNER

На рис. 11.3 имя владельца (OWNER) встречается один раз на листовой блок, а не один раз на повторяющуюся запись. Мы запустим приведенный ниже сценарий, передав ему число 1, чтобы воссоздать ситуацию, когда индекс использует сжатие только на головном столбце:

```
drop index t_idx;
create index t_idx on
  t(owner,object_type,object_name)
  compress &1;
analyze index t_idx validate structure;
insert into idx_stats
  select 'compress &1', a.*
  from index_stats a;
```

В целях сравнения мы запустим этот сценарий не только с одним столбцом, но также с двумя и тремя сжатыми столбцами, и посмотрим, что произойдет. Запросив в конце IDX_STATS, мы должны наблюдать следующее:

```
EODA@ORA12CR1> select what, height, lf_blks, br_blks,
  2      btree_space, opt_cmpr_count, opt_cmpr_pctsave
  3      from idx_stats
  4      /
```

| WHAT | HEIGHT | LF_BKLS | BR_BKLS | BTREE_SPACE | OPT_CMPR_COUNT | OPT_CMPR_PCTSAVE |
|---------------|--------|---------|---------|-------------|----------------|------------------|
| noncompressed | 2 | 227 | 1 | 1823120 | 2 | 28 |
| compress 1 | 2 | 206 | 1 | 1654380 | 2 | 21 |
| compress 2 | 2 | 162 | 1 | 1302732 | 2 | 0 |
| compress 3 | 2 | 268 | 1 | 2149884 | 2 | 39 |

Мы видим, что размер индекса COMPRESS 1 составляет примерно 90% от размера несжатого индекса (сравнивая значения в столбце BTREE_SPACE). Количество листовых блоков заметно уменьшилось. В случае применения COMPRESS 2 экономия еще более впечатляюща. Размер результирующего индекса становится примерно 71% размера исходного индекса. Фактически, используя столбец OPT_CMPR_PCTSAVE, который означает *optimum compression percent saved* (оптимальный процент экономии от сжатия), или ожидаемую экономию от сжатия, мы можем сделать предположение о размере индекса COMPRESS 2:

```
EODA@ORA12CR1> select 1823120*(1-0.28) from dual;
1823120*(1-0.28)
-----
1312646.4
```

На заметку! Команда ANALYZE в отношении несжатого индекса заполняет столбцы OPT_CMPR_PCTSAVE/OPT_CMPR_COUNT и оценивает экономию от сжатия COMPRESS 2 в 28%, чего мы почти точно достигли в действительности.

Но обратите внимание, что происходит со сжатием COMPRESS 3. Полученный в результате индекс больше исходного. Его размер составляет 117% от размера исходного несжатого индекса. Это объясняется тем фактом, что каждый повторяющийся префикс, который мы удаляем, экономит пространство для N своих копий, однако

добавляет 4 байта накладных расходов в листовом блоке как часть схемы сжатия. Добавляя столбец `OBJECT_NAME` к сжатому ключу, мы делаем этот ключ почти уникальным — в данном случае это означает отсутствие дублированных копий, которые можно было бы факторизовать. Следовательно, мы просто *добавляем* 4 байта практически к каждому отдельному ключу индекса и не факторизуем какие-либо повторяющиеся данные. Столбец `OPT_CMPR_COUNT` в `IDX_STATS` предназначен как раз для предоставления оптимального количества сжатых столбцов, а `OPT_CMPR_PCTSAVE` сообщает точно, какой экономии можно ожидать.

Помните, что вы не получаете это сжатие бесплатно. Структура сжатого индекса теперь сложнее, чем обычно бывает. База данных Oracle будет тратить больше времени на обработку данных в этой структуре как при поддержании индекса во время проведения модификаций, так и при поиске в индексе во время выполнения запросов. В действительности за сокращение времени ввода-вывода мы расплачиваемся возросшим временем работы процессора. Благодаря сжатию кеш буферов блоков будет способен уместить больше записей индекса, чем ранее, процент попаданий к кеш может возрасти, а количество физических операций ввода-вывода должно снизиться, но это потребует чуть больших вычислительных возможностей процессора для обработки индекса и также увеличит вероятность конкуренции за блоки. Вспомните обсуждение хеш-кластера, где говорилось, что извлечение миллиона произвольных строк может занять больше процессорного времени, но половину объема ввода-вывода. Как и тогда, мы должны помнить о компромиссе. Если процессор является ограниченным ресурсом, то добавление индексов со сжатым ключом может замедлить обработку. С другой стороны, если существуют ограничения, касающиеся ввода-вывода, то применение таких индексов может увеличить скорость работы.

Индексы по реверсированным ключам

Еще одной возможностью индекса со структурой B-дерева является способность реверсировать ключи. Поначалу может возникнуть вопрос о том, зачем это может понадобиться. Индексы со структурой B-дерева были спроектированы для специфической среды и решения определенной проблемы. Они реализованы так, чтобы снизить конкуренцию за листовые блоки индекса в “правосторонних” индексах, таких как индексы на столбцах, которые заполняются последовательными значениями или отметками времени в среде Oracle RAC.

На заметку! Технология Oracle RAC обсуждалась в главе 2.

RAC — это конфигурация Oracle, в которой множество экземпляров могут монтировать и открывать одну и ту же базу данных. Если двум экземплярам необходимо модифицировать единственный блок данных одновременно, они будут совместно использовать блок, передавая его вперед и назад по внутреннему аппаратному соединению — частному сетевому подключению между двумя (или более) машинами. При наличии индекса по первичному ключу на столбце, заполняемом из последовательности (весьма популярная реализация), во время вставки новых значений каждый процесс будет пытаться модифицировать один блок, который в текущий момент является левым блоком в правой стороне индексной структуры (на рис. 11.1 было показано, что высокие значения уходят внутри индекса вправо, а низкие —

влево). Модификации индексов на столбцах, заполняемых последовательностями, сосредоточены на небольшом множестве листовых блоков. Реверсирование ключей индекса позволяет распределить вставки по всем листовым блокам в индексе, хотя это может сделать его гораздо менее эффективно упакованным.

На заметку! Вы можете также счесть индексы по реверсированным ключам полезными в качестве метода сокращения конкуренции даже в одиночном экземпляре Oracle. Вы будете применять их главным образом для того, чтобы смягчить влияние событий ожидания занятого буфера, возникающих в правой части занятого индекса, как было описано в этом разделе.

Прежде чем мы посмотрим, как измерить влияние индекса по реверсированному ключу, давайте обсудим физическое действие этого индекса. Индекс по реверсированному ключу просто меняет порядок следования байт в каждом столбце своего ключа на противоположный. Если мы возьмем числа 90101, 90102 и 90103 и взглянем на их внутреннее представление с использованием Oracle-функции DUMP, то получим следующий вывод:

```
EODA@ORA12CR1> select 90101, dump(90101,16) from dual
2 union all
3 select 90102, dump(90102,16) from dual
4 union all
5 select 90103, dump(90103,16) from dual
6 /

90101 DUMP(90101,16)
-----
90101 Typ=2 Len=4: c3,a,2,2
90102 Typ=2 Len=4: c3,a,2,3
90103 Typ=2 Len=4: c3,a,2,4
```

Внутреннее представление каждого числа имеет длину 4 байта, причем отличается только последний байт. Эти числа должны оказаться в структуре индекса рядом друг с другом. Но если изменить порядок следования байтов в числах на противоположный, то Oracle вставит такие данные:

```
90101 reversed = 2,2,a,c3
90102 reversed = 3,2,a,c3
90103 reversed = 4,2,a,c3
```

Полученные числа будут располагаться далеко друг от друга. Это сокращает количество экземпляров RAC, обращающихся к одному и тому же блоку (крайнему справа), и уменьшает объем передач блоков между экземплярами RAC. Одним из недостатков индекса по реверсированному ключу является невозможность его применения во всех случаях, в которых может использоваться обычный индекс. Например, при формировании ответа на запрос со следующим предикатом индекс по реверсированному ключу X окажется бесполезным:

```
where x > 5
```

Данные в индексе перед сохранением были отсортированы не по X, а по REVERSE(X), поэтому просмотр диапазона для $X > 5$ не сможет задействовать этот

индекс. С другой стороны, некоторые просмотры диапазона могут быть сделаны с применением индекса по реверсированному ключу. Если имеется сцепленный индекс на (X, Y) , то следующий предикат будет в состоянии использовать его и выполнит просмотр диапазона:

```
where x = 5
```

Это объясняется тем, что сначала реверсированы байты X , а затем байты Y . База данных Oracle не реверсирует байты $(X || Y)$, а взамен сохраняет $(\text{REVERSE}(X) || \text{REVERSE}(Y))$. В итоге все значения для $X=5$ будут сохраняться вместе, так что Oracle сможет просканировать этот индекс и найти их все.

Теперь предположим, что в таблице имеется суррогатный первичный ключ, заполняемый через последовательность, и просмотр диапазона по этому индексу не требуется — т.е. нет необходимости выполнять запросы $\text{MAX}(\text{primary_key})$, $\text{MIN}(\text{primary_key})$, $\text{WHERE primary_key} < 100$ и т.д. В таком случае можно рассмотреть возможность применения индекса по реверсированному ключу для сценариев с интенсивной вставкой строк, даже в единственном экземпляре Oracle. Для демонстрации отличий между вставкой в таблицу, имеющую индекс по реверсированному ключу на первичном ключе, и в таблицу с обычным индексом я подготовил два разных теста: один использует чистую среду PL/SQL, а другой — язык Pro*C. В обоих случаях таблица создавалась с помощью показанного ниже кода DDL (мы избегаем конкуренции за блоки таблицы за счет применения ASSM, так что мы можем изолировать конкуренцию за индексные блоки):

```
create tablespace assm
datafile size 1m autoextend on next 1m
segment space management auto;

create table t tablespace assm
as
select 0 id, owner, object_name, subobject_name,
       object_id, data_object_id, object_type, created,
       last_ddl_time, timestamp, status, temporary,
       generated, secondary
from all_objects a
where 1=0;

alter table t add constraint t_pk primary key (id)
using index (create index t_pk on t(id) &indexType tablespace assm);

create sequence s cache 1000;
```

Здесь `&indexType` заменяется либо ключевым словом `REVERSE` с целью создания индекса по реверсированному ключу, либо ничем, что приводит к использованию “обычного” индекса. Код на PL/SQL, который должен запускаться параллельно одним, двумя, пятью, десятью, пятнадцатью или двадцатью пользователями, выглядит следующим образом:

```
create or replace procedure do_sql
as
begin
  for x in ( select rownum r, OWNER, OBJECT_NAME, SUBOBJECT_NAME,
                  OBJECT_ID, DATA_OBJECT_ID, OBJECT_TYPE, CREATED,
                  LAST_DDL_TIME, TIMESTAMP, STATUS, TEMPORARY,
```

```

GENERATED, SECONDARY from all_objects )
loop
  insert into t
    ( id, OWNER, OBJECT_NAME, SUBOBJECT_NAME,
      OBJECT_ID, DATA_OBJECT_ID, OBJECT_TYPE, CREATED,
      LAST_DDL_TIME, TIMESTAMP, STATUS, TEMPORARY,
      GENERATED, SECONDARY )
  values
    ( s.nextval, x.OWNER, x.OBJECT_NAME, x.SUBOBJECT_NAME,
      x.OBJECT_ID, x.DATA_OBJECT_ID, x.OBJECT_TYPE, x.CREATED,
      x.LAST_DDL_TIME, x.TIMESTAMP, x.STATUS, x.TEMPORARY,
      x.GENERATED, x.SECONDARY );
  if ( mod(x.r,100) = 0 )
  then
    commit;
  end if;
end loop;
commit;
end;
/

```

Так как мы обсуждали в главе 9 оптимизацию времени фиксации PL/SQL, я теперь хочу запустить тест, применяющий другую среду, чтобы эта оптимизация не вводила в заблуждение. С помощью Pro*C я эмулирую процедуру извлечения, трансформации и загрузки (Extraction, Transformation and Loading — ETL) для хранилища данных, которая обрабатывает строки пакетами по 100 за раз между фиксациями:

```

exec sql declare c cursor for select * from all_objects;
exec sql open c;
exec sql whenever notfound do break;
for(;;)
{
  exec sql
  fetch c into :owner:owner_i,
  :object_name:object_name_i, :subobject_name:subobject_name_i,
  :object_id:object_id_i, :data_object_id:data_object_id_i,
  :object_type:object_type_i, :created:created_i,
  :last_ddl_time:last_ddl_time_i, :timestamp:timestamp_i,
  :status:status_i, :temporary:temporary_i,
  :generated:generated_i, :secondary:secondary_i;
  exec sql
  insert into t
    ( id, OWNER, OBJECT_NAME, SUBOBJECT_NAME,
      OBJECT_ID, DATA_OBJECT_ID, OBJECT_TYPE, CREATED,
      LAST_DDL_TIME, TIMESTAMP, STATUS, TEMPORARY,
      GENERATED, SECONDARY )
  values
    ( s.nextval, :owner:owner_i, :object_name:object_name_i,
      :subobject_name:subobject_name_i, :object_id:object_id_i,
      :data_object_id:data_object_id_i, :object_type:object_type_i,
      :created:created_i, :last_ddl_time:last_ddl_time_i,
      :timestamp:timestamp_i, :status:status_i,
      :temporary:temporary_i, :generated:generated_i,
      :secondary:secondary_i );
}

```

```

    if ( ++cnt%100) == 0 )
    {
        exec sql commit;
    }
}
exec sql whenever notfound continue;
exec sql commit;
exec sql close c;

```

Код Pro*C предварительно скомпилирован с параметром PREFETCH, установленным в 100, что делает этот код аналогичным коду PL/SQL в Oracle 10g. Например, сохранив приведенный выше код Pro*C в файле по имени t.pc, команда вызова компилятора Pro*C будет выглядеть так:

```
$ proc iname=t.pc MODE=ORACLE PREFETCH=100
```

На заметку! В Oracle 10g Release 1 и последующих версиях простая конструкция FOR X IN (SELECT * FROM T) в PL/SQL будет молча делать выборку массивом по 100 строк за раз, тогда как в Oracle9i и предшествующих версиях она извлекает по одной строке за раз. Таким образом, если вы хотите воспроизвести этот пример в Oracle9i и предшествующих версиях, придется модифицировать код PL/SQL, чтобы он также производил выборку массивом, используя синтаксис BULK COLLECT.

Оба теста должны извлекать по 100 строк за раз и затем вставлять одну строку в другую таблицу. В табл. 11.1–11.4 подводятся итоги отличий между разными запусками, начиная с однопользовательского теста в табл. 11.1.

Таблица 11.1. Тест производительности применения индексов по реверсированным ключам с помощью PL/SQL и Pro*C: случай с одним пользователем

| | Реверсированный индекс, PL/SQL | Обычный индекс, PL/SQL | Реверсированный индекс, Pro*C | Обычный индекс, Pro*C |
|--|--------------------------------|------------------------|-------------------------------|-----------------------|
| Транзакций в секунду | 41,5 | 41,5 | 42,2 | 45,5 |
| Время ЦП (секунды) | 0,38 | 0,37 | 0,73 | 0,69 |
| Ожидания занятого буфера (количество/секунды) | 0/0 | 0/0 | 0/0 | 0/0 |
| Затраченное время (минуты) | 0,38 | 0,34 | 3,31 | 3,09 |
| Синхронизации журнальных файлов (количество/секунды) | 2/0 | 2/0 | 650/4 | 650/4 |

По первому тесту с одним пользователем мы можем видеть, что при выполнении этой операции PL/SQL заметно эффективнее, чем Pro*C, и эту тенденцию мы продолжим наблюдать по мере роста пользовательской нагрузки. Отчасти причиной того, что Pro*C не так хорошо масштабируется, как PL/SQL, является необходимость ожидания синхронизации журнальных файлов, которая имеет место в Pro*C, в то время как PL/SQL располагает оптимизацией, чтобы избежать этого ожидания.

Из этого однопользовательского теста следует, что индексы по реверсированным ключам потребляют чуть больше ресурсов процессора. Это имеет смысл, поскольку

ку база данных должна выполнять дополнительную работу по изменению порядка следования байтов в ключе на противоположный. Однако по мере увеличения количества пользователей мы обнаружим, что эта логика перестает быть справедливой. Когда возникает конкуренция, накладные расходы, связанные с индексом по реверсированному ключу, полностью исчезают. На самом деле даже в случае теста с двумя пользователями накладные расходы почти перекрываются конкуренцией за правую часть индекса (табл. 11.2).

Таблица 11.2. Тест производительности применения индексов по реверсированным ключам с помощью PL/SQL и Pro*C: случай с двумя пользователями

| | Реверсированный индекс, PL/SQL | Обычный индекс, PL/SQL | Реверсированный индекс, Pro*C | Обычный индекс, Pro*C |
|--|--------------------------------|------------------------|-------------------------------|-----------------------|
| Транзакций в секунду | 55,0 | 55,0 | 59,1 | 53,8 |
| Время ЦП (секунды) | 0,80 | 0,77 | 1,57 | 1,55 |
| Ожидания занятого буфера (количество/секунды) | 823/0 | 615/0 | 649/0 | 1 580/0 |
| Затраченное время (минуты) | 0,79 | 0,75 | 6,99 | 6,90 |
| Синхронизации журнальных файлов (количество/секунды) | 3/0 | 3/0 | 1 229/19 | 1 227/26 |

По результатам этого теста с двумя пользователями можно утверждать, что PL/SQL по-прежнему превосходит Pro*C, но использование индекса по реверсированному ключу демонстрирует некоторые преимущества со стороны PL/SQL и не такие большие — со стороны Pro*C. Эта тенденция сохранится. Для программы Pro*C индекс по реверсированному ключу решает проблему ожидания занятого буфера, вызванную конкуренцией за крайний правый блок индексной структуры; тем не менее, он никак не влияет на ожидание синхронизации журнальных файлов, которое затрагивает программу на Pro*C. В этом была главная причина проведения двух тестов, PL/SQL и Pro*C — увидеть отличия между указанными двумя средами. Возникает вопрос: почему в данном случае индекс по реверсированному ключу обеспечивает очевидные преимущества для PL/SQL, но не для Pro*C? Это связано с событием ожидания синхронизации журнальных файлов. Среда PL/SQL была способна непрерывно выполнять вставку и редко вынуждена ожидать события синхронизации журнальных файлов при фиксации, тогда как среда Pro*C ожидала синхронизации через каждые 100 строк. Следовательно, в этом случае среда PL/SQL была больше подвержена влиянию со стороны событий ожидания занятого буфера по сравнению с Pro*C. Смягчение проблемы ожидания занятого буфера в случае PL/SQL позволяет обработать больше транзакций, поэтому применение индекса по реверсированному ключу оказывает положительное воздействие на PL/SQL. Но в случае Pro*C события ожидания занятого буфера не являются проблемой — они не были главным узким местом в плане производительности, так что устранение событий ожидания никак не повлияло на общую производительность.

Давайте перейдем к тесту с пятью пользователями (табл. 11.3).

Таблица 11.3. Тест производительности применения индексов по реверсированным ключам с помощью PL/SQL и Pro*C: случай с пятью пользователями

| | Реверсирован- ный индекс, PL/SQL | Обычный индекс, PL/SQL | Реверсирован- ный индекс, Pro*C | Обычный индекс, Pro*C |
|---|--|------------------------------|---------------------------------------|-----------------------------|
| Транзакций в секунду | 82,2 | 82,2 | 65,6 | 70,3 |
| Время ЦП (секунды) | 1,93 | 1,91 | 3,97 | 4,36 |
| Ожидания занятого буфера (количество/секунды) | 1 963/1 | 2 644/1 | 2 707/0 | 9 839/1 |
| Затраченное время (минуты) | 5,26 | 5,59 | 22,14 | 22,17 |
| Синхронизации журнальных файлов (количество/секунды) | 6/0 | 6/0 | 3 061/138 | 3 202/128 |

Здесь мы видим все то же самое. На среду PL/SQL,двигающуюся на всех парах, с небольшим числом событий ожидания синхронизации журнальных файлов, оказывали очень сильное влияние ожидания занятого буфера. В случае обычного индекса, когда все пять пользователей пытались вставлять строки в правую часть индексной структуры, среда PL/SQL больше всего страдала от событий ожидания занятого буфера, поэтому она больше выиграла от сокращения количества таких событий.

Взглянув на результаты тестирования с десятью пользователями, показанные в табл. 11.4, мы увидим, что тенденция сохраняется.

Таблица 11.4. Тест производительности применения индексов по реверсированным ключам с помощью PL/SQL и Pro*C: случай с десятью пользователями

| | Реверсирован- ный индекс, PL/SQL | Обычный индекс, PL/SQL | Реверсирован- ный индекс, Pro*C | Обычный индекс, Pro*C |
|---|--|------------------------------|---------------------------------------|-----------------------------|
| Транзакций в секунду | 88,3 | 91,2 | 96,8 | 88,1 |
| Время ЦП (секунды) | 3,83 | 4,06 | 8,57 | 10,01 |
| Ожидания занятого буфера (количество/секунды) | 2 897/28 | 6 831/7 | 5 284/7 | 34 312/25 |
| Затраченное время (минуты) | 25,6 | 26,04 | 116,62 | 124,70 |
| Синхронизации журнальных файлов (количество/секунды) | 11/0 | 119/0 | 6 051/301 | 6 441/352 |

В отсутствие ожидания синхронизации журнальных файлов код PL/SQL очень много выигрывает от исключения ожиданий занятого буфера. Код Pro*C теперь показывает более высокую конкуренцию при ожидании занятого буфера, но из-за того, что часто приходится ожидать синхронизации журнальных файлов, никаких преимуществ не получает. Одним из способов увеличения производительности реализации PL/SQL с обычным индексом было бы введение небольшой задержки. Это снизило бы конкуренцию за правую часть индекса и улучшило общую производительность. В целях экономии места я не привожу результаты тестов с 15 и 20 пользователями, но подтверждаю, что с ростом количества пользователей наблюдаемая тенденция сохраняется.

Совет. Исходный код для примеров оценки производительности индекса по реверсированному ключу доступен для загрузки на веб-сайте издательства. В каталоге ch11 находится несколько файлов demo3* (а также файл t.pc), которые автоматизируют выполнение всего этого тестового комплекта.

Из этой демонстрации можно сделать два вывода. Индекс по реверсированному ключу может содействовать в смягчении последствий ожидания занятого буфера, но в зависимости от других факторов вы получите варьирующийся эффект. Глядя на результаты теста для десяти пользователей в табл. 11.4, несложно заметить, что устранение ожиданий занятого буфера (наиболее длительное из всех событий ожидания в данном случае) влияет на пропускную способность транзакций минимально, но при этом показывает повышенную масштабируемость с более высокой степенью параллелизма. Выполнение того же действия применительно к PL/SQL имеет заметно отличающееся влияние на производительность: мы достигаем умеренного роста пропускной способности за счет устранения этого узкого места.

Индексы, упорядоченные по убыванию

Индексы, упорядоченные по убыванию, появились в Oracle8i с целью расширения функциональности индексов со структурой B-дерева. Они позволяют столбцам храниться в индексе отсортированными по убыванию (от большего значения к меньшему), а не по возрастанию (от меньшего значения к большему). В версиях Oracle, предшествующих Oracle8i, ключевое слово DESC (по убыванию) всегда поддерживалось *синтаксически*, но в основном игнорировалось — оно не оказывало никакого влияния на то, как сохранялись либо использовались данные в индексе. Однако в Oracle8i и последующих версиях ключевое слово DESC изменяет способ создания и применения индекса.

В Oracle давно поддерживалась возможность чтения индекса в обратном порядке, поэтому вас может удивить, почему это средство является важным. Например, создадим таблицу T, как показано ниже:

```
EODA@ORA12CR1> create table t
2 as
3 select *
4   from all_objects
5 /
```

Table created.

Таблица создана.

```
EODA@ORA12CR1> create index t_idx on t(owner,object_type,object_name);
Index created.
```

Индекс создан.

```
EODA@ORA12CR1> begin
2     dbms_stats.gather_table_stats
3       ( user, 'T', method_opt=>'for all indexed columns' );
4 end;
5 /
```

PL/SQL procedure successfully completed.

Процедура PL/SQL успешно завершена.

Теперь запустим следующий запрос к этой таблице:


```

EODA@ORA12CR1> set autotrace traceonly explain
EODA@ORA12CR1> select owner, object_type
2   from t
3   where owner between 'T' and 'Z'
4     and object_type is not null
5   order by owner DESC, object_type DESC;

```

Execution Plan

Plan hash value: 2685572958

| Id | Operation | Name | Rows | Bytes | Cost (%CPU) | Time |
|-----|-----------------------------|-------|------|-------|-------------|----------|
| 0 | SELECT STATEMENT | | 5008 | 50080 | 24 (0) | 00:00:01 |
| * 1 | INDEX RANGE SCAN DESCENDING | T_IDX | 5008 | 50080 | 24 (0) | 00:00:01 |

База данных Oracle будет просто читать индекс в обратном порядке. В этом плане отсутствует финальный шаг сортировки — данные уже сортированы. Однако средство индексов, упорядоченных по убыванию, вступает в игру, когда имеется смесь столбцов, часть которых отсортирована по возрастанию (ASC), а часть — по убыванию (DESC), например:

```

EODA@ORA12CR1> select owner, object_type
2   from t
3   where owner between 'T' and 'Z'
4     and object_type is not null
5   order by owner DESC, object_type ASC;

```

Execution Plan

Plan hash value: 2813023843

| Id | Operation | Name | Rows | Bytes | Cost (%CPU) | Time |
|-----|------------------|-------|------|-------|-------------|----------|
| 0 | SELECT STATEMENT | | 5008 | 50080 | 24 (0) | 00:00:01 |
| 1 | SORT ORDER BY | | 5008 | 50080 | 24 (0) | 00:00:01 |
| * 2 | INDEX RANGE SCAN | T_IDX | 5008 | 50080 | 24 (0) | 00:00:01 |

Predicate Information (identified by operation id):

```

2 - access("OWNER">='T' AND "OWNER"<='Z')
    filter("OBJECT_TYPE" IS NOT NULL)

```

База данных Oracle больше не может использовать индекс на (OWNER, OBJECT_TYPE, OBJECT_NAME) для *сортировки* данных. Она должна была бы читать его в обратном порядке, чтобы получить данные, отсортированные по столбцу OWNER по убыванию, но читать его “вперед” для получения значений OBJECT_TYPE, отсортированных по возрастанию. Вместо этого Oracle собирает вместе все строки и затем сортирует их. Введем индекс DESC:

```

EODA@ORA12CR1> create index desc_t_idx on t(owner desc,object_type asc);
Index created.
Индекс создан.

```

```

EODA@ORA12CR1> select owner, object_type
2   from t
3   where owner between 'T' and 'Z'
4     and object_type is not null
5   order by owner DESC, object_type ASC;

```

Execution Plan

Plan hash value: 2494308350

| Id | Operation | Name | Rows | Bytes | Cost (%CPU) | Time |
|-----|------------------|------------|------|-------|-------------|----------|
| 0 | SELECT STATEMENT | | 5008 | 50080 | 2 (0) | 00:00:01 |
| * 1 | INDEX RANGE SCAN | DESC_T_IDX | 5008 | 50080 | 2 (0) | 00:00:01 |

Predicate Information (identified by operation id):

```

1 - access(SYS_OP_DESCEND("OWNER")>=HEXTORAW('A5FF') AND
      SYS_OP_DESCEND("OWNER")<=HEXTORAW('ABFF'))
      filter(SYS_OP_UNDESCEND(SYS_OP_DESCEND("OWNER"))>='T' AND
      SYS_OP_UNDESCEND(SYS_OP_DESCEND("OWNER"))<='Z' AND
      "OBJECT_TYPE" IS NOT NULL)

```

Мы снова можем читать данные в отсортированном виде, а дополнительный шаг сортировки в конце плана отсутствует.

На заметку! Никогда не поддавайтесь искушению убрать конструкцию ORDER BY из запроса. То, что ваш план выполнения включает индекс, вовсе не означает, что строки будут возвращаться в “каком-то порядке”. Единственный способ извлечь данные из базы в некотором отсортированном порядке предусматривает включение в запрос конструкции ORDER BY. Никакой замены ORDER BY не существует.

Когда должен использоваться индекс со структурой В-дерева?

Не являясь большим поклонником “эмпирических правил” (из каждого правила есть исключения), я не располагаю такими правилами относительно того, когда применять (или не применять) индекс со структурой В-дерева. Чтобы продемонстрировать, почему у меня отсутствуют любые эмпирические правила для этого случая, я представлю два в равной степени справедливых утверждения.

- Используйте В-дерево для индексации столбцов, только если вы собираетесь получать доступ через индекс к очень небольшому проценту строк таблицы.
- Применяйте индекс со структурой В-дерева, только если вы планируете обрабатывать много строк таблицы, и индекс может использоваться вместо таблицы.

Приведенные утверждения кажутся противоречащими друг другу, но в реальности это не так — они просто описывают два совершенно разных случая. С учетом предыдущей рекомендации существуют два способа применения индекса.

- **Как средство доступа к строкам таблицы.** Вы читаете индекс, чтобы получить строку таблицы. Здесь вам необходим доступ к очень небольшому проценту строк в таблице.

- **Как средство ответа на запрос.** Индекс содержит достаточно информации, чтобы ответить на целый запрос — при этом обращаться к таблице вообще не нужно. Индекс будет использоваться в качестве облегченной версии таблицы.

Есть также и другие способы — например, индекс можно применять для извлечения *всех* строк таблицы, включая столбцы, которые не содержатся в самом индексе. Это очевидным образом противоречит обоим представленным правилам. Такой случай может иметь место в интерактивном приложении, где вы извлекаете некоторые строки и отображаете их, затем — еще немного строк и т.д. Вы хотите иметь запрос, оптимизированный по начальному времени ответа, а не по общему времени выполнения.

Первый случай (т.е. использование индекса для доступа к небольшому проценту записей таблицы) говорит, что если у вас есть таблица T (применяется та же самая таблица, что и ранее), и вы имеете план выполнения запроса, который выглядит следующим образом:

```
EODA@ORA12CR1> set autotrace traceonly explain
EODA@ORA12CR1> select owner, status
2   from t
3   where owner = USER;
```

Execution Plan

Plan hash value: 1695850079

| Id | Operation | Name | Rows | Bytes |
|-------------|-------------------------------------|------------|------|-------|
| Cost (%CPU) | ... | | | |
| 0 | SELECT STATEMENT | | 1716 | 17160 |
| 13 (0) | ... | | | |
| 1 | TABLE ACCESS BY INDEX ROWID BATCHED | T | 1716 | 17160 |
| 13 (0) | ... | | | |
| * 2 | INDEX RANGE SCAN | DESC_T_IDX | 288 | |
| 2 (0) | ... | | | |

Predicate Information (identified by operation id):

```
2 - access(SYS_OP_DESCEND("OWNER")=SYS_OP_DESCEND(USER@!))
   filter(SYS_OP_UNDESCEND(SYS_OP_DESCEND("OWNER"))=USER@!)
```

то вы должны обращаться к очень небольшому проценту строк в этой таблице. Проблемой, на которую здесь необходимо обратить внимание, является операция INDEX (RANGE SCAN), за которой следует операция TABLE ACCESS BY INDEX ROWID. Это значит, что Oracle прочтает индекс, после чего для его записей будет выполняться чтение блока базы данных (логическую или физическую операцию ввода-вывода), чтобы получить данные строк. Это не самый эффективный метод, если вы собираетесь получить доступ к большому проценту строк в таблице T через индекс (вскоре мы определим, каким может быть этот большой процент).

Во втором случае (т.е. когда индекс может использоваться *вместо* таблицы) вы можете обработать 100% (или меньше) строк через индекс. Индекс можно применять для создания облегченной версии таблицы.

Следующий запрос демонстрирует эту концепцию:

```
EODA@ORA12CR1> select count(*)
2   from t
3   where owner = user;
```

Execution Plan

Plan hash value: 293504097

| Id | Operation | Name | Rows | Bytes | Cost (%CPU) | Time |
|-----|------------------|-------|------|-------|-------------|----------|
| 0 | SELECT STATEMENT | | 1 | 3 | 10 (0) | 00:00:01 |
| 1 | SORT AGGREGATE | | 1 | 3 | | |
| * 2 | INDEX RANGE SCAN | T_IDX | 1716 | 5148 | 10 (0) | 00:00:01 |

Predicate Information (identified by operation id):

2 - access("OWNER"=USER@!)

Здесь для выдачи ответа на запрос использовался только индекс — не имеет значения, к какому проценту строк мы обращались, т.к. применялся только индекс. В плане выполнения видно, что доступ к таблице вообще не производился; мы просто просматривали саму индексную структуру.

Важно понимать разницу между этими двумя концепциями. Когда нужно выполнить TABLE ACCESS BY INDEX ROWID, мы должны гарантировать, что будем обращаться только к небольшому проценту от общего количества блоков в таблице, что обычно соответствует небольшому проценту строк, или же нам *необходимо* извлечь первые строки насколько возможно быстро (их с нетерпением ждет конечный пользователь). Если мы обращаемся к слишком высокому проценту строк (больше, чем от 1% до 20% общего количества), то доступ к ним через индекс со структурой В-дерева обычно займет больше времени, чем при полном сканировании таблицы.

Со вторым типом запроса, где ответ находится полностью в индексе, мы имеем другую историю. Мы читаем блок индекса и обнаруживаем много строк для обработки, затем переходим к следующему блоку и т.д. — мы никогда не обращаемся к таблице. Существует также *быстрое полное сканирование*, которое можно выполнять на индексах, чтобы в отдельных случаях еще более ускорить эту операцию. Быстрое полное сканирование — это когда база данных читает индексные блоки без определенного порядка; она просто начинает читать их. Индекс больше не используется как собственно индекс, а скорее как таблица. Строки из быстрого полного сканирования не поступают упорядоченными по записям индекса.

В общем случае индекс со структурой В-дерева должен быть построен на столбцах, которые часто применяются в предикатах запроса, и ожидается возвращение небольшой части данных таблицы или же конечный пользователь требует немедленного отклика. В *тонкой* таблице (т.е. в таблице с несколькими столбцами или со столбцами небольшого размера) такая часть может быть очень малой. Запрос, использующий этот индекс, должен ожидать извлечения из таблицы 2–3% строк или менее. В *толстой* таблице (т.е. в таблице с большим количеством столбцов или с очень широкими столбцами) такая часть *может* достигать 20–25% объема таблицы. Этот совет не сразу выглядит осмысленным; он не является интуитивно понятным, но он точен.

Индекс хранится отсортированным по индексному ключу. Доступ к нему осуществляется в порядке сортировки по ключу. Блоки, на которые он указывает, хранятся произвольно в куче. Поэтому, когда читается индекс для доступа к таблице, выполняется множество *разбросанных*, произвольных операций ввода-вывода. Под “разбросанными” понимается то, что индекс будет указывать на необходимость чтения блока 1, блока 1000, блока 205, блока 321, блока 1, блока 1032, блока 1 и т.д. — он не запросит чтение блока 1, затем блока 2, потом блока 3 в последовательной манере. Появится тенденция читать и перечитывать блоки в совершенно случайном порядке. Ввод-вывод одиночного блока может оказаться очень медленным.

В качестве упрощенного примера предположим, что мы читаем такую тонкую таблицу через индекс, и собираемся прочесть 20% строк. Пусть в таблице находится 100 000 строк. Двадцать процентов — это 20 000 строк. Если каждая строка имеет длину около 80 байт в базе данных с размером блоков 8 Кбайт, то на блок приходится около 100 строк. Это значит, что таблица имеет приблизительно 1000 блоков. Дальнейшие вычисления очень просты. Нам нужно прочесть 20 000 строк через индекс; это будет означать, скорее всего, 20 000 операций TABLE ACCESS BY ROWID. При выполнении такого запроса мы обработаем 20 000 блоков. Однако во всей таблице имеется только около 1000 блоков! То есть каждый блок таблицы будет прочитан и обработан в среднем 20 раз. Даже если увеличить размер строки до 800 байт на строку и 10 строк на блок, то в таблице получится 10 000 блоков. Доступ через индекс к 20 000 строк все равно приведет к чтению каждого блока в среднем 2 раза. В этом случае полное сканирование таблицы будет намного более эффективным, чем применение индекса, т.к. оно потребует только однократного затрагивания каждого блока. Любой запрос, использующий этот индекс для доступа к данным, будет не слишком эффективным до тех пор, пока он не станет обращаться к менее чем 5% объема данных для 800-байтового столбца (тогда мы получаем доступ примерно к 5000 блокам) и еще меньше — для 80-байтового столбца (около 0,5% или меньше).

Физическая организация

Физическая организация данных на диске серьезно влияет на эти вычисления, т.к. она по существу определяет то, насколько дорогим (или недорогим) будет индексный доступ. Предположим, что есть таблица, в которой строки имеют первичный ключ, заполняемый последовательностью. По мере добавления данных к таблице строки с последовательными значениями первичного ключа могут в общем случае размещаться рядом друг с другом.

На заметку! Применение таких средств, как метод ASSM либо несколько списков свободных блоков или групп списков свободных блоков, окажет влияние на то, каким образом данные организованы на диске. Эти средства стремятся рассредоточить данные, и такая естественная кластеризация по первичному ключу может не наблюдаться.

Таблица естественным образом кластеризуется по порядку первичного ключа (поскольку данные добавляются более-менее в этом порядке). Разумеется, она не будет строго кластеризована по порядку этого ключа (чтобы достичь этого, мы должны были бы использовать индекс-таблицу); в общем, строки с близкими значениями первичного ключа будут размещаться физически по соседству. Выдав следующий запрос:

```
select * from T where primary_key between :x and :y
```

вы обнаружите, что необходимые строки обычно находятся в тех же самых блоках. В этом случае сканирование диапазона по индексу может оказаться практичным, даже если оно обращается к большому проценту строк — просто потому, что блоки базы данных, которые нужно читать и перечитывать, с высокой вероятностью окажутся кешированными, т.к. данные находятся рядом. С другой стороны, если строки не расположены по соседству, то применение того же индекса может быть губительным для производительности. Подтвердить этот факт поможет небольшая демонстрация. Начнем с таблицы, которая практически упорядочена по первичному ключу:

```
EODA@ORA12CR1> create table colocated ( x int, y varchar2(80) );
Table created.
```

Таблица создана.

```
EODA@ORA12CR1> begin
2   for i in 1 .. 100000
3   loop
4       insert into colocated(x,y)
5       values (i, rpad(dbms_random.random,75,'*') );
6   end loop;
7 end;
8 /
```

PL/SQL procedure successfully completed.

Процедура PL/SQL успешно завершена.

```
EODA@ORA12CR1> alter table colocated
2 add constraint colocated_pk
3 primary key(x);
```

Table altered.

Таблица изменена.

```
EODA@ORA12CR1> begin
2 dbms_stats.gather_table_stats( user, 'COLOCATED' );
3 end;
4 /
```

PL/SQL procedure successfully completed.

Процедура PL/SQL успешно завершена.

Эта таблица соответствует приведенному выше описанию о 100 строках на блок в базе данных с размером блока 8 Кбайт. В этой таблице существует высокая вероятность того, что строки со значениями X, равными 1, 2 и 3, окажутся в одном блоке. Теперь возьмем эту таблицу и намеренно “дезорганизуем” ее. Мы создадим в таблице COLOCATED столбец Y с головным случайным числом и воспользуемся этим фактом для дезорганизации данных, чтобы они определенно больше не были упорядоченными по первичному ключу:

```
EODA@ORA12CR1> create table disorganized
2 as
3 select x,y
4 from colocated
5 order by y;
```

Table created.

Таблица создана.

```
EODA@ORA12CR1> alter table disorganized
2 add constraint disorganized_pk
3 primary key (x);
```

Table altered.

Таблица изменена.

```
EODA@ORA12CR1> begin
```

```
  2 dbms_stats.gather_table_stats( user, 'DISORGANIZED');
```

```
  3 end;
```

```
  4 /
```

PL/SQL procedure successfully completed.

Процедура PL/SQL успешно завершена.

Возможно, эти таблицы одинаковы — мы имеем дело с реляционной базой данных, значит, физическая организация не оказывает влияния на ответы, возвращаемые для запросов (во всяком случае, так утверждается в теории баз данных). На проверку выясняется, что характеристики производительности этих двух таблиц отличаются буквально как день и ночь, хотя возвращаемые ответы идентичны. Выдав в точности такой же запрос и применив тот же самый план выполнения, просмотрим вывод TKPROF (трассировку SQL):

```
select * from colocated where x between 20000 and 40000
```

| call | count | cpu | elapsed | disk | query | current | rows |
|---------|-------|------|---------|------|-------|---------|--------|
| Parse | 5 | 0.00 | 0.00 | 0 | 0 | 0 | 0 |
| Execute | 5 | 0.00 | 0.00 | 0 | 0 | 0 | 0 |
| Fetch | 6675 | 0.06 | 0.21 | 0 | 14495 | 0 | 100005 |
| total | 6685 | 0.06 | 0.21 | 0 | 14495 | 0 | 100005 |

...

Rows (1st) Rows (avg) Rows (max) Row Source Operation

| | | | |
|-------|-------|-------|---|
| 20001 | 20001 | 20001 | TABLE ACCESS BY INDEX ROWID BATCHED COLOCATED... |
| 20001 | 20001 | 20001 | INDEX RANGE SCAN COLOCATED_PK (cr=1374 pr=0 pw=0... |

```
select /*+ index( disorganized disorganized_pk ) */ * from disorganized
where x between 20000 and 40000
```

| call | count | cpu | elapsed | disk | query | current | rows |
|---------|-------|------|---------|------|--------|---------|--------|
| Parse | 5 | 0.00 | 0.00 | 0 | 0 | 0 | 0 |
| Execute | 5 | 0.00 | 0.00 | 0 | 0 | 0 | 0 |
| Fetch | 6675 | 0.12 | 0.41 | 0 | 106830 | 0 | 100005 |
| total | 6685 | 0.12 | 0.41 | 0 | 106830 | 0 | 100005 |

...

Rows (1st)Rows (avg) Rows (max) Row Source Operation

| | | | |
|-------|-------|-------|---|
| 20001 | 20001 | 20001 | TABLE ACCESS BY INDEX ROWID BATCHED DISORGANIZED... |
| 20001 | 20001 | 20001 | INDEX RANGE SCAN DISORGANIZED_PK(cr=1374 pr=0 pw=0... |

На заметку! Все запросы запускались по пять раз для получения обоснованного среднего значения времени выполнения каждого из них (поэтому выводе TKPROF отражена обработка свыше 100 000 строк).

Я считаю это поразительным. Вот к чему может привести разное физическое размещение данных! В табл. 11.5 представлены итоговые результаты.

Таблица 11.5. Исследование влияния физического размещения данных на стоимость индексного доступа

| Таблица | Время ЦП | Логический ввод-вывод |
|----------------------------------|--------------|-----------------------|
| COLOCATED | 0,21 секунды | 14 495 |
| DISORGANIZED | 0,41 секунды | 106 830 |
| Процент по сравнению с COLOCATED | ~50% | 13% |

В моей базе данных, использующей размер блока 8 Кбайт, эти таблицы имеют следующее общее количество блоков каждая:

```

EODA@ORA12CR1> select a.index_name,
2      b.num_rows,
3      b.blocks,
4      a.clustering_factor
5  from user_indexes a, user_tables b
6  where index_name in ('COLOCATED_PK', 'DISORGANIZED_PK' )
7    and a.table_name = b.table_name
8  /

```

| INDEX_NAME | NUM_ROWS | BLOCKS | CLUSTERING_FACTOR |
|-----------------|----------|--------|-------------------|
| COLOCATED_PK | 100000 | 1252 | 1190 |
| DISORGANIZED_PK | 100000 | 1219 | 99929 |

Запрос к дезорганизованной таблице подтверждает простые расчеты, которые мы делали ранее: понадобилось более 20 000 логических операций ввода-вывода (было запрошено всего 100 000 блоков и производилось пять запусков запроса). Каждый блок обрабатывался по 20 раз! С другой стороны, физически совместное размещение данных значительно сокращает объем логического ввода-вывода. Это блестящая иллюстрация того, насколько трудно ориентироваться на эмпирические правила; в одном случае применение индексов работает замечательно, а в другом — нет. Вспомните об этом в следующий раз, когда перенесете данные из рабочей системы в среду разработки, т.к. вы можете получить, по крайней мере, частичный ответ на вопрос: “Почему на этой машине все работает иначе — разве они не идентичны?”. Нет, они не идентичны.

На заметку! Вспомните из главы 6, что увеличенный объем логического ввода-вывода является лишь верхушкой айсберга. Каждая логическая операция ввода-вывода влечет за собой одну или большее число защелок буферного кеша. В многопользовательской/многопроцессорной среде использование процессора вторым запросом, несомненно, продлится во много раз меньше, чем первым, поскольку мы производим раскрутку и ожидаем получение защелок. Второй пример запроса не только выполняет больше работы, но еще и масштабируется хуже первого.

Влияние ARRAYSIZE на логический ввод-вывод

Интересно отметить влияние ARRAYSIZE на выполняемый логический ввод-вывод. Параметр ARRAYSIZE устанавливает количество строк, которые Oracle возвращает клиенту, когда он запрашивает следующую строку. Клиент затем буферизует эти строки и применяет их перед запросом у базы данных следующего набора строк. Значение ARRAYSIZE может оказывать очень существенное влияние на логический ввод-вывод, выполняемый запросом, из-за того факта, что если необходимо получить доступ к одному и тому же блоку снова и снова между обращениями к базе данных (между командами выборки в данном случае), то Oracle придется заново извлекать этот блок из буферного кеша. Следовательно, если вы запросите 100 строк у базы данных в одном обращении, то Oracle может быть в состоянии полностью обработать блок базы данных и не нуждаться в его повторном извлечении. В случае запроса 15 строк за раз Oracle может понадобиться вновь и вновь получать тот же самый блок, чтобы извлечь один и тот же набор строк.

В примере, приведенном ранее в этом разделе, использовался стандартный размер массива извлечения SQL*Plus, равный 15 строк (если вы разделите общее количество извлеченных строк (100 005) на количество команд выборки (6675), то получите результат, очень близкий к 15). Если сравнить выполнение предыдущих запросов при 15 строках на команду выборки с 100 строками на команду выборки, то для таблицы COLOCATED можно наблюдать следующее:

| Rows | Row | Source | Operation |
|--|------------|------------|---|
| Rows (1st) | Rows (avg) | Rows (max) | Row Source Operation |
| 20001 | 20001 | 20001 | TABLE ACCESS BY INDEX ROWID BATCHED COLOCATED |
| (cr=2899 pr=0 pw=0 ... | | | |
| 20001 | 20001 | 20001 | INDEX RANGE SCAN COLOCATED_PK |
| (cr=1374 pr=0 pw=0 ... | | | |
| select * from colocated a100 where x between 20000 and 40000 | | | |
| Rows (1st) | Rows (avg) | Rows (max) | Row Source Operation |
| 20001 | 20001 | 20001 | TABLE ACCESS BY INDEX ROWID BATCHED COLOCATED |
| (cr=684 pr=0 pw=0 ... | | | |
| 20001 | 20001 | 20001 | INDEX RANGE SCAN COLOCATED_PK |
| (cr=245 pr=0 pw=0 ... | | | |

Первый запрос запускался с ARRAYSIZE, равным 15, и значения (cr=nnnn) в столбце Row Source Operation (Исходная строковая операция) показывают, что было выполнено 1 374 логических операции ввода-вывода по индексу и затем 1 525 логических операций ввода-вывода по таблице (т.е. 2 899 - 1 374; значения в Row Source Operation являются накопительными). Когда значение ARRAYSIZE было увеличено с 15 до 100 (посредством команды SET ARRAYSIZE 100), общий объем ввода-вывода по индексу снизился до 245. Это стало непосредственным результатом того, что повторное чтение листовых блоков индекса из буферного кеша производилось через каждые 100 строк, а не 15 строк, как было ранее. Чтобы понять это, предположим, что мы могли хранить 200 строк на листовой блок. По мере сканирования индекса с чтением по 15 строк за раз приходится извлекать первый листовой блок 14 раз, чтобы получить из него все 200 записей. С другой стороны, когда извлекается по 100 строк за раз, тот же самый листовой блок понадобится извлечь из буферного кеша всего два раза, чтобы получить все его записи.

То же самое происходит в случае табличных блоков. Поскольку таблица была отсортирована в том же порядке, что и ключи индекса, мы будем извлекать каждый блок таблицы менее часто, т.к. с помощью каждой команды выборки получим больше строк.

Итак, если это было хорошо для таблицы COLOCATED, то должно быть не менее хорошо и для таблицы DISARGANIZATED, правильно? Не совсем. Результаты для таблицы DISARGANIZATED будут выглядеть так:

```
select /*+ index( a15 disorganized_pk ) */ *
from disorganized a15 where x between 20000 and 40000
```

```
Rows (1st) Rows (avg) Rows (max) Row Source Operation
```

```
-----
      20001      20001      20001 TABLE ACCESS BY INDEX ROWID BATCHED DISORGANIZED
(cr=21365 pr=0 ...
      20001      20001      20001 INDEX RANGE SCAN DISORGANIZED_PK
(cr=1374 pr=0...
```

```
select /*+ index( a100 disorganized_pk ) */ *
from disorganized a100 where x between 20000 and 40000
```

```
Rows (1st) Rows (avg) Rows (max) Row Source Operation
```

```
-----
      20001      20001      20001 TABLE ACCESS BY INDEX ROWID BATCHED DISORGANIZED
(cr=20236 pr=0 ...
      20001      20001      20001 INDEX RANGE SCAN DISORGANIZED_PK
(cr=245 pr=0...
```

Результаты по индексу идентичны, что вполне понятно, поскольку данные сохраняются в индексе одинаково независимо от организации *таблицы*. Количество логических операций ввода-вывода для одного запуска запроса снизилось с 1 374 до 245, как и ранее. Но общий объем логического ввода-вывода, выполненного этим запросом, изменился незначительно: 21 365 операций против 20 236. В чем причина? Количество логических операций ввода-вывода по таблице вообще не изменилось — если вычесть число логических операций ввода-вывода по индексу из общего количества логических операций ввода-вывода, выполненных каждым запросом, то обнаружится, что оба запроса выполнили 19 991 логическую операцию ввода-вывода в отношении таблицы. Это потому, что каждый раз, когда требуется N строк из базы данных — вероятность того, что любые две из этих строк окажутся в одном и том же блоке, очень мала — нет шансов получить множество строк из табличного блока за единственное обращение.

Каждый известный мне профессиональный язык программирования из тех, что способны взаимодействовать с Oracle, реализует свою концепцию извлечения массивов. В PL/SQL можно применять BULK COLLECT или положиться на неявный размер извлечения массива в 100 элементов, который обеспечивается для неявного курсора в циклах for. В Java/JDBC имеется метод упреждающей выборки у объекта подключения или объекта оператора. Интерфейс уровня вызовов Oracle (Oracle Call Interface — OCI; API-интерфейс C) позволяет программно устанавливать размер упреждающей выборки, как это делает Pro*C. Несложно заметить, что это может оказать заметное влияние на объем логического ввода-вывода, выполняемого запросом, а потому заслуживает вашего внимания.

Просто чтобы завершить пример, давайте посмотрим, что произойдет при полном сканировании таблицы DISORGANIZED:

```
select * from disorganized where x between 20000 and 30000
```

| call | count | cpu | elapsed | disk | query | current | rows |
|---------|-------|------|---------|------|-------|---------|-------|
| Parse | 1 | 0.00 | 0.00 | 0 | 0 | 0 | 0 |
| Execute | 1 | 0.00 | 0.00 | 0 | 0 | 0 | 0 |
| Fetch | 668 | 0.01 | 0.03 | 0 | 1858 | 0 | 10001 |
| total | 670 | 0.01 | 0.03 | 0 | 1858 | 0 | 10001 |

```
Rows (1st) Rows (avg) Rows (max) Row Source Operation
```

```
-----
      10001      10001      10001 TABLE ACCESS FULL DISORGANIZED (cr=1858
pr=0 pw=0...
```

Таким образом, в этом конкретном случае полное сканирование очень хорошо подходит из-за способа физического хранения данных на диске. Однако возникает вопрос: почему оптимизатор сразу не выбрал полное сканирование для данного запроса? Вообще говоря, это оставляется на заложенное в него проектное решение, но в первом примере запроса к таблице DISORGANIZED я намеренно снабдил запрос подсказкой и сообщил оптимизатору о необходимости сконструировать план, который использует индекс. Во втором случае я позволил оптимизатору выбрать наилучший общий план.

Фактор кластеризации

Теперь давайте посмотрим на некоторую информацию, применяемую Oracle. Мы специально обратимся к столбцу CLUSTERING_FACTOR (фактор кластеризации), находящемуся в представлении USER_INDEXES. В руководстве *Oracle Database Reference* назначение этого столбца описано следующим образом.

Указывает степень упорядоченности строк в таблице на основе значений индекса.

- Если значение близко к количеству блоков, то таблица очень хорошо упорядочена. В этом случае записи индекса в единственном листовом блоке имеют тенденцию указывать на строки в тех же самых блоках данных.
- Если значение близко к количеству строк, то таблица упорядочена весьма случайно. В этом случае маловероятно, что записи индекса из одного листового блока будут указывать на строки в тех же самых блоках данных.

Фактор кластеризации можно также рассматривать как число, представляющее количество логических операций ввода-вывода в таблице, которые должны быть выполнены для чтения всей таблицы через индекс. То есть CLUSTERING_FACTOR — это признак того, как упорядочена таблица в отношении самого индекса, и если взглянуть на эти индексы, обнаружится следующее:

```
EODA@ORA12CR1> select a.index_name,
2      b.num_rows,
3      b.blocks,
4      a.clustering_factor
5  from user_indexes a, user_tables b
6  where index_name in ('COLOCATED_PK', 'DISORGANIZED_PK' )
7     and a.table_name = b.table_name
8  /
```

| INDEX_NAME | NUM_ROWS | BLOCKS | CLUSTERING_FACTOR |
|-----------------|----------|--------|-------------------|
| COLOCATED_PK | 100000 | 1252 | 1190 |
| DISORGANIZED_PK | 100000 | 1219 | 99929 |

На заметку! В примере этого раздела используется табличное пространство, управляемое ASSM, что объясняет, почему фактор кластеризации для таблицы COLOCATED меньше, чем количество блоков в таблице. В предстоящей таблице COLOCATED имеются неформатированные блоки перед HWM-маркером, которые не содержат данных, а также блоки, применяемые самим методом ASSM для управления пространством, и эти блоки не будут читаться даже при сканировании диапазона по индексу. Более подробные объяснения HWM и ASSM приведены в главе 10.

Итак, база данных сообщает: “Если бы мы читали каждую строку в таблице COLOCATED через индекс COLOCATED_PK от начала до конца, то инициировали бы 1 190 операций ввода-вывода. Однако если бы мы делали то же самое с таблицей DISORGANIZED, то выполнили бы 99 929 операций ввода-вывода”. Причина такой большой разницы связана с тем, что Oracle производит сканирование диапазона по структуре индекса, и если обнаруживается, что следующая строка индекса находится в том же блоке базы данных, что и предыдущая строка, то еще одна операция ввода-вывода для получения блока таблицы из буферного кеша не выполняется. Дескриптор для блока таблицы уже имеется и он просто используется. Тем не менее, если следующая строка находится *не* в том же самом блоке, то Oracle освобождает этот блок и выполняет еще одну операцию ввода-вывода в буферном кеше, чтобы извлечь следующий блок для обработки. Поэтому индекс COLOCATED_PK в процессе сканирования его диапазона обнаруживает, что следующая строка почти всегда располагается в том же блоке, что и предыдущая. Индекс DISORGANIZED_PK обнаруживает противоположную ситуацию. В действительности мы можем увидеть, что это измерение является очень точным. Применяя подсказки, чтобы заставить оптимизатор использовать полное сканирование индекса для чтения всей таблицы, и просто подсчитав количество отличных от NULL значений Y, мы можем выяснить, сколько точно операций ввода-вывода понадобится для чтения всей таблицы через индекс:

```
select count(Y) from
  (select /*+ INDEX(COLOCATED COLOCATED_PK) */ * from colocated)
```

| call | count | cpu | elapsed | disk | query | current | rows |
|---------|-------|------|---------|------|-------|---------|------|
| Parse | 1 | 0.00 | 0.00 | 0 | 0 | 0 | 0 |
| Execute | 1 | 0.00 | 0.00 | 0 | 0 | 0 | 0 |
| Fetch | 2 | 0.03 | 0.03 | 0 | 1399 | 0 | 1 |
| total | 4 | 0.03 | 0.03 | 0 | 1399 | 0 | 1 |

```
...
Rows (1st) Rows (avg) Rows (max) Row Source Operation
-----
          1          1          1 SORT AGGREGATE (cr=1399 pr=0 pw=0 time=34740 us)
100000      100000      100000 TABLE ACCESS BY INDEX ROWID BATCHED
COLOCATED (cr=1399 pr=0 pw=0 time=90620 us cost=1400 size=7600000...
100000      100000      100000 INDEX FULL SCAN COLOCATED_PK (cr=209 pr=0 pw=0...
*****
select count(Y) from
  (select /*+ INDEX(DISORGANIZED DISORGANIZED_PK) */ * from disorganized)
```

| call | count | cpu | elapsed | disk | query | current | rows |
|---------|-------|------|---------|------|--------|---------|------|
| Parse | 1 | 0.00 | 0.00 | 0 | 0 | 0 | 0 |
| Execute | 1 | 0.00 | 0.00 | 0 | 0 | 0 | 0 |
| Fetch | 2 | 0.11 | 0.11 | 0 | 100138 | 0 | 1 |
| total | 4 | 0.11 | 0.11 | 0 | 100138 | 0 | 1 |

```
...
```

Rows (1st) Rows (avg) Rows (max) Row Source Operation

```

      1          1          1 SORT AGGREGATE (cr=100138 pr=0 pw=0 time=111897 us)
100000    100000    100000 TABLE ACCESS BY INDEX ROWID BATCHED DISORGANIZED
(cr=100138 pr=0 pw=0 time=203332 us cost=100158 size=7600000 card=100000)
100000    100000    100000 INDEX FULL SCAN DISORGANIZED_PK (cr=209 pr=0 pw=0

```

В обоих случаях индекс нуждается в выполнении 209 логических операций ввода-вывода (cr=209 в строках столбца Row Source Operation). Если вы вычтете значение 209 из общего числа согласованных чтений и измерите только количество операций ввода-вывода в таблице, то обнаружите, что оно идентично фактору клас-теризации каждого соответствующего индекса. Индекс COLOCATED_PK — это класси-ческий пример “хорошо упорядоченной таблицы”, в то время как DISORGANIZE_PK — классический пример “весьма произвольно упорядоченной таблицы”. Интересно взглянуть, как теперь это повлияет на оптимизатор. Если мы попытаемся извлечь 25 000 строк, то Oracle теперь выберет полное сканирование таблицы для обоих за-просов (извлечение 25% строк через индекс не является оптимальным планом даже для очень хорошо упорядоченной таблицы). Однако снизим выборку до 10% таблич-ных данных (имейте в виду, что 10% — это не пороговое значение, а просто число мень-ше 25%, которое в данном случае вызывает сканирование диапазона по индексу):

```

EODA@ORA12CR1> set autotrace traceonly explain
EODA@ORA12CR1> select * from colocated where x between 20000 and 30000;

```

Execution Plan

Plan hash value: 2792740192

| Id | Operation | Name | Rows | Bytes | Cost | ... |
|-----|-------------------------------------|--------------|-------|-------|------|-----|
| 0 | SELECT STATEMENT | | 10002 | 791K | 142 | ... |
| 1 | TABLE ACCESS BY INDEX ROWID BATCHED | COLOCATED | 10002 | 791K | 142 | ... |
| * 2 | INDEX RANGE SCAN | COLOCATED_PK | 10002 | | 22 | ... |

Predicate Information (identified by operation id):

2 - access("X">=20000 AND "X"<=30000)

```

EODA@ORA12CR1> select * from disorganized where x between 20000 and 30000;

```

Execution Plan

Plan hash value: 2727546897

| Id | Operation | Name | Rows | Bytes | Cost (%CPU) | Time |
|-----|-------------------|--------------|-------|-------|-------------|----------|
| 0 | SELECT STATEMENT | | 10002 | 791K | 333 (1) | 00:00:01 |
| * 1 | TABLE ACCESS FULL | DISORGANIZED | 10002 | 791K | 333 (1) | 00:00:01 |

Predicate Information (identified by operation id):

1 - filter("X"<=30000 AND "X">=20000)

Здесь мы имеем те же самые структуры таблиц и те же индексы, но разные факторы кластеризации. В этом случае оптимизатор выбирает план доступа через индекс для таблицы COLOCATED и план доступа с полным сканированием для таблицы DISORGANIZED.

Ключевой момент этого обсуждения состоит в том, что индексы не всегда являются подходящим методом доступа. Оптимизатор может совершенно резонно отказываться от применения индекса, что демонстрирует предшествующий пример. На использование индекса оптимизатором влияют многие факторы, *включая* физическую организацию данных. Поэтому может возникнуть соблазн перестроить таблицы так, чтобы все индексы имели хороший фактор кластеризации, но в большинстве случаев подобное будет *бесполезной тратой времени*. Это касается случаев, когда производится сканирование диапазона по индексу для значительного процента строк таблицы. Вдобавок вы должны иметь в виду, что в общем случае таблица будет иметь только *один* индекс с хорошим фактором кластеризации! Строки в таблице могут быть отсортированы единственным способом. Если бы в показанном только что примере существовал другой индекс на столбце Y, то он был бы очень плохо кластеризован в таблице COLOCATED, но очень хорошо кластеризован в таблице DISORGANIZED. Если для вас важно иметь физически кластеризованные данные, рассмотрите возможность применения индекс-таблицы, кластера со структурой В-дерева или же хеш-кластера с непрерывной перестройкой таблиц.

Заключительные соображения по поводу индексов со структурой В-дерева

Индексы со структурой В-дерева — наиболее распространенные и хорошо понимаемые индексные структуры в Oracle. Они представляют собой великолепный механизм индексации универсального назначения. Такие индексы обеспечивают хорошо масштабируемое время доступа, возвращая данные из индексов в 1000 строк примерно за то же самое время, что и из индексов в 100 000 строк.

Когда и какие столбцы индексировать — это аспекты, которым следует уделять внимание при проектировании. Индекс не всегда означает быстрый доступ; фактически во многих случаях вы обнаружите, что индексы даже снижают производительность, когда Oracle их использует. Это функция, основанная исключительно на том, насколько большому проценту строк таблицы требуется получать доступ через индекс, и как физически расположены данные. Если вы можете применять индекс для ответа на вопрос, то доступ к значительному проценту строк имеет смысл, поскольку вы избегаете дополнительных разбросанных операций ввода-вывода в таблице. Если вы используете индекс для доступа к таблице, то должны гарантировать обработку только небольшого процента строк таблицы.

Структура и реализация индексов должны продумываться *во время* проектирования приложения, а не потом (что я часто наблюдаю). При тщательном планировании и с учетом способа доступа к данным необходимые индексы будут очевидными практически во всех случаях.

Битовые индексы

Битовые индексы были добавлены в версии Oracle 7.3. В настоящее время они доступны в редакциях Enterprise и Personal, но не Standard. Битовые индексы предназначены для сред хранилищ данных или сред с нерегламентированными запросами, в которых полный набор запросов к данным не полностью известен во время реализации системы. Они *не* предназначены для систем OLTP или систем, где данные часто обновляются многочисленными параллельными сеансами.

Битовые индексы — это структуры, которые хранят указатели на множество строк в единственном элементе индексного ключа, что отличает их от индексов со структурой В-дерева, где существует соответствие между ключами индекса и строками таблицы. В битовом индексе содержится очень малое количество записей, каждая из которых указывает на большое число строк. В обычном индексе со структурой В-дерева одна запись индекса указывает на единственную строку.

Предположим, что мы создаем битовый индекс на столбце JOB в таблице EMP:

```
EODA@ORA12CR1> create BITMAP index job_idx on emp(job);
Index created.
Индекс создан.
```

В табл. 11.6 показано то, что Oracle будет сохранять в этом индексе.

Таблица 11.6. Представление того, как Oracle будет хранить битовый индекс JOB_IDX

| Значение/ строка | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 |
|---------------------|---|---|---|---|---|---|---|---|---|----|----|----|----|----|
| ANALYST | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 1 | 0 |
| CLERK | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 1 |
| MANAGER | 0 | 0 | 0 | 1 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| PRESIDENT | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 |
| SALESMAN | 0 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

В табл. 11.6 показано, что строки 8, 10 и 13 содержат значение ANALYST, тогда как строки 4, 6 и 7 — значение MANAGER. В ней также видно, что нет ни одной строки со значением NULL (битовые индексы хранят записи NULL; их отсутствие в индексе говорит о том, что в таблице нет строк NULL). Если необходимо выяснить количество строк со значением MANAGER, то битовый индекс позволяет сделать это очень быстро. Если требуется найти все строки, у которых JOB равно CLERK или MANAGER, то можно просто объединить их битовые карты из индекса (табл. 11.7).

Таблица 11.7. Представление битового “ИЛИ”

| Значение/ строка | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 |
|---------------------|---|---|---|---|---|---|---|---|---|----|----|----|----|----|
| CLERK | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 1 |
| MANAGER | 0 | 0 | 0 | 1 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| CLERK or MANAGER | 1 | 0 | 0 | 1 | 0 | 1 | 1 | 0 | 0 | 0 | 1 | 1 | 0 | 1 |

Просматривая табл. 11.7, можно быстро понять, что строки 1, 4, 6, 7, 11, 12 и 14 удовлетворяют нашему критерию. Битовая карта, в которой Oracle хранит каждое значение ключа, устроена так, что каждая позиция представляет идентификатор строки внутри лежащей в основе таблицы — на случай, если понадобится действительно извлечь строку для дальнейшей обработки. Запросы вида

```
select count(*) from emp where job = 'CLERK' or job = 'MANAGER';
```

находят ответ непосредственно в битовом индексе. С другой стороны, запросы типа

```
select * from emp where job = 'CLERK' or job = 'MANAGER';
```

нуждаются в обращении к таблице. Здесь Oracle применит функцию для преобразования единичных битов в результирующий идентификатор строки, который будет использоваться для доступа в таблицу.

Когда должен использоваться битовый индекс?

Битовые индексы больше всего подходят для данных с *низкой кардинальностью* (т.е. данных с относительно небольшим количеством возможных дискретных значений, если сравнивать его с мощностью всего множества). Точно указать величину этой кардинальности невозможно; другими словами, трудно определить, что собой на самом деле представляет низкая кардинальность. Для множества из пары тысяч записей низкой кардинальностью может быть 2, но это не так в случае таблицы с двумя строками. В таблице с десятками или сотнями миллионов записей низкой кардинальностью вполне может оказаться и 100 000. Таким образом, низкая кардинальность является относительной к размеру результирующего набора. Это данные, для которых результат деления количества отличающихся значений в наборе строк на количество строк дает малое число (близкое к нулю). Например, столбец GENDER (пол) может принимать значения M (мужской), F (женский) и NULL (не указан). Если у вас есть таблица с 20 000 записями о сотрудниках, то $3 / 20\,000 = 0,00015$. Аналогично, 100 000 уникальных значений из 10 000 000 в результате деления дает 0,01 — опять-таки, очень малое число. Такие столбцы могут быть кандидатами на создание битовых индексов. Вероятно, они *не* будут кандидатами на участие в индексах со структурой B-дерева, т.к. каждое возможное их значение оказалось бы связанным с исключительно большим процентом строк в таблице. Ранее подчеркивалось, что индексы со структурой B-дерева в общем случае должны быть селективными. Битовые индексы не должны быть селективными — наоборот, обычно они должны быть совершенно *неселективными*.

Битовые индексы чрезвычайно удобны в средах, где происходит много нерегламентированных запросов, особенно запросов, которые ссылаются на множество столбцов произвольным образом или вычисляют агрегаты вроде COUNT. Например, пусть имеется крупная таблица с тремя столбцами: GENDER, LOCATION и AGE_GROUP. В этой таблице столбец GENDER принимает значения M или F, LOCATION — значения от 1 до 50, а AGE_GROUP — код, представляющий признаки 18 and under, 19–25, 26–30, 31–40 и 41 and over. Вам необходимо обслуживать огромное количество специальных запросов в следующей форме:

```
select count(*)
  from t
 where gender = 'M'
```



```

and location in ( 1, 10, 30 )
and age_group = '41 and over';

select *
from t
where ( ( gender = 'M' and location = 20 )
      or ( gender = 'F' and location = 22 ) )
and age_group = '18 and under';

select count(*) from t where location in (11,20,30);
select count(*) from t where age_group = '41 and over' and gender = 'F';

```

Вы обнаружите, что обычная схема индексации со структурой В-дерева не подойдет. Если вы хотите применять индекс для получения ответа, то при доступе к данным через индекс понадобятся минимум три и максимум шесть комбинаций возможных индексов со структурой В-дерева. Поскольку может встретиться любой из трех столбцов или любое их подмножество, потребуются большие сцепленные индексы со структурой В-дерева на перечисленных ниже столбцах.

- GENDER, LOCATION, AGE_GROUP. Для запросов, использующих все три столбца, столбец GENDER с LOCATION или только столбец GENDER.
- LOCATION, AGE_GROUP. Для запросов, применяющих столбцы LOCATION и AGE_GROUP либо только столбец LOCATION.
- AGE_GROUP, GENDER. Для запросов, использующих столбец AGE_GROUP с GENDER или только столбец AGE_GROUP.

Чтобы сократить объем искомых данных, могут оказаться обоснованными и другие перестановки, также призванные снизить размер просматриваемой индексной структуры. Это еще игнорирует тот факт, что создание индекса со структурой В-дерева на данных с такой низкой кардинальностью вообще считается плохой идеей.

Именно здесь в игру вступают битовые индексы. С тремя небольшими битовыми индексами, по одному на каждом отдельном столбце, вы сможете эффективно удовлетворить все приведенные выше предикаты. Чтобы найти результирующий набор для любого предиката, который ссылается на произвольное множество из этих трех столбцов, Oracle просто будет применять функции AND, OR и NOT, объединяя битовые карты трех индексов. База данных возьмет результирующую объединенную битовую карту, при необходимости преобразует единицы в идентификаторы строк и обратится к данным (если вы просто выясняете количество строк, соответствующих критерию, то Oracle будет просто подсчитывать биты, равные 1). Давайте рассмотрим пример. Сначала мы сгенерируем тестовые данные, удовлетворяющие указанной кардинальности, проиндексируем их и соберем статистику. Для генерации случайных данных с подходящим распределением используется пакет DBMS_RANDOM.

```

EODA@ORA12CR1> create table t
2 ( gender not null,
3   location not null,
4   age_group not null,
5   data
6 )
7 as
8 select decode( round(dbms_random.value(1,2)),

```

```

9 1, 'M',
10 2, 'F' ) gender,
11 ceil(dbms_random.value(1,50)) location,
12 decode( round(dbms_random.value(1,5)),
13 1,'18 and under',
14 2,'19-25',
15 3,'26-30',
16 4,'31-40',
17 5,'41 and over'),
18 rpad( '*', 20, '*' )
19 from dual connect by level <=100000;

```

Table created.

Таблица создана.

```
EODA@ORA12CR1> create bitmap index gender_idx on t(gender);
```

Index created.

Индекс создан.

```
EODA@ORA12CR1> create bitmap index location_idx on t(location);
```

Index created.

Индекс создан.

```
EODA@ORA12CR1> create bitmap index age_group_idx on t(age_group);
```

Index created.

Индекс создан.

```
EODA@ORA12CR1> exec dbms_stats.gather_table_stats( user, 'T');
```

PL/SQL procedure successfully completed.

Процедура PL/SQL успешно завершена.

Теперь взглянем на планы выполнения показанных ранее специальных запросов:

```
EODA@ORA12CR1> set autotrace traceonly explain
```

```
EODA@ORA12CR1> select count(*)
```

```

2   from t
3   where gender = 'M'
4     and location in ( 1, 10, 30 )
5     and age_group = '41 and over';

```

Execution Plan

Plan hash value: 320981916

| Id | Operation | Name | Rows | Bytes | Cost (%CPU) |
|-----|---------------------------|--------------|------|-------|-------------|
| 0 | SELECT STATEMENT | | 1 | 13 | 9 (0) |
| 1 | SORT AGGREGATE | | 1 | 13 | |
| 2 | BITMAP CONVERSION COUNT | | 608 | 7904 | 9 (0) |
| 3 | BITMAP AND | | | | |
| 4 | BITMAP OR | | | | |
| * 5 | BITMAP INDEX SINGLE VALUE | LOCATION_IDX | | | |

| | | | | | |
|-----|---------------------------|---------------|--|--|--|
| * 6 | BITMAP INDEX SINGLE VALUE | LOCATION_IDX | | | |
| ... | | | | | |
| * 7 | BITMAP INDEX SINGLE VALUE | LOCATION_IDX | | | |
| ... | | | | | |
| * 8 | BITMAP INDEX SINGLE VALUE | AGE_GROUP_IDX | | | |
| ... | | | | | |
| * 9 | BITMAP INDEX SINGLE VALUE | GENDER_IDX | | | |
| ... | | | | | |

Predicate Information (identified by operation id):

- 5 - access("LOCATION"=1)
- 6 - access("LOCATION"=10)
- 7 - access("LOCATION"=30)
- 8 - access("AGE_GROUP"='41 and over')
- 9 - access("GENDER"='M')

Этот пример демонстрирует мощь битовых индексов. База данных Oracle способна увидеть конструкцию `location in (1,10,30)` и знает, что нужно прочитывать индекс на столбце `location` для этих трех значений и с помощью логической операции `OR` объединить “биты” в битовой карте. Затем она берет результирующую битовую карту и посредством логической операции `AND` складывает ее с битовыми картами для `AGE_GROUP='41 AND OVER'` и `GENDER='M'`. После этого Oracle просто подсчитывает единичные биты — и ответ готов.

```
EODA@ORA12CR1> select *
2   from t
3   where ( ( gender = 'M' and location = 20 )
4          or ( gender = 'F' and location = 22 ) )
5          and age_group = '18 and under';
```

Execution Plan

Plan hash value: 705811684

| Id | Operation | Name | Rows | Bytes |
|----------------|-------------------------------------|---------------|------|-------|
| Cost (%CPU)... | | | | |
| 0 | SELECT STATEMENT | | 408 | 13872 |
| 68 (0)... | | | | |
| 1 | TABLE ACCESS BY INDEX ROWID BATCHED | T | 408 | 13872 |
| 68 (0)... | | | | |
| 2 | BITMAP CONVERSION TO ROWIDS | | | |
| ... | | | | |
| 3 | BITMAP AND | | | |
| ... | | | | |
| * 4 | BITMAP INDEX SINGLE VALUE | AGE_GROUP_IDX | | |
| ... | | | | |
| 5 | BITMAP OR | | | |
| ... | | | | |
| 6 | BITMAP AND | | | |
| ... | | | | |

| | | | | |
|------|---------------------------|--------------|--|--|
| * 7 | BITMAP INDEX SINGLE VALUE | LOCATION_IDX | | |
| ... | | | | |
| * 8 | BITMAP INDEX SINGLE VALUE | GENDER_IDX | | |
| ... | | | | |
| 9 | BITMAP AND | | | |
| ... | | | | |
| * 10 | BITMAP INDEX SINGLE VALUE | LOCATION_IDX | | |
| ... | | | | |
| * 11 | BITMAP INDEX SINGLE VALUE | GENDER_IDX | | |
| ... | | | | |

 Predicate Information (identified by operation id):

```

 4 - access("AGE_GROUP"='18 and under')
 7 - access("LOCATION"=22)
 8 - access("GENDER"='F')
10 - access("LOCATION"=20)
11 - access("GENDER"='M')

```

Логика здесь похожа: в плане отражены условия OR, каждое из которых вычисляется путем объединения посредством операции AND с соответствующими битовыми картами и затем объединения с помощью OR результатов вычислений. Применив еще одну операцию AND для удовлетворения условия AGE_GROUP='18 AND UNDER', мы получаем все, что нужно. Поскольку в этот раз запрашиваются действительные строки, Oracle преобразует каждую единицу и ноль битовой карты в идентификатор строки для извлечения исходных данных.

В хранилище данных или крупной системе генерации отчетов, поддерживающей много нерегламентированных SQL-запросов, возможность совместного использования такого количества индексов, какое имеет смысл, на самом деле становится очень полезной. Применение в этом случае обычных индексов со структурой B-дерева даже близко не сравнится по удобству и эффективности, и по мере роста числа столбцов, по которым производится поиск в нерегламентированных запросах, увеличивается также и количество необходимых комбинаций индексов со структурой B-дерева.

Однако бывают ситуации, когда битовые индексы *не* подходят. Они хорошо работают в среде с интенсивным чтением, но очень плохо ведут себя в среде с интенсивной записью. Причина в том, что одиночная запись ключа битового индекса указывает на *множество* строк. Если сеанс модифицирует проиндексированные данные, то все строки, на которые указывает запись индекса, в большинстве случаев блокируются. База данных Oracle не может блокировать индивидуальный бит в записи битового индекса: она блокирует запись целиком. Любые другие модификации, которым необходимо обновить *ту же самую запись битового индекса*, будут заблокированы. Это серьезно подавляет параллелизм, т.к. каждое обновление будет потенциально блокировать сотни строк, предотвращая параллельные обновления их битовых столбцов. Блокироваться будет *не каждая* строка, как можно было бы подумать, но многие из них. Битовые индексы сохраняются порциями, поэтому, обратившись к приведенному ранее примеру с таблицей EMP, мы можем обнаружить, что индексный ключ ANALYST появляется в индексе много раз, каждый раз указывая на сотни строк. Обновление строки, которое модифицирует столбец JOB, потребует

получения монопольного доступа к двум записям индексного ключа: записи индексного ключа для *старого* значения и записи индексного ключа для *нового* значения. Сотни строк, на которые указывают эти две записи, будут недоступными для модификации другими сеансами до тех пор, пока не произойдет фиксации текущего оператора UPDATE.

Битовые индексы соединений

В версии Oracle9i был добавлен еще один тип индекса: битовый индекс соединения. Обычно индекс создается на одиночной таблице с использованием столбцов только из этой таблицы. Битовый индекс соединения выходит за рамки такого правила и позволяет индексировать заданную таблицу с применением столбцов из какой-то другой таблицы. В сущности, это предоставляет возможность денормализации данных в индексной структуре вместо того, чтобы проводить денормализацию в самих таблицах.

Рассмотрим простые таблицы EMP и DEPT. Таблица EMP имеет внешний ключ для DEPT (столбец DEPTNO). Таблица DEPT содержит атрибут DNAME (название отдела). Конечные пользователи будут часто задавать вопросы вроде: “Сколько людей занимается продажами?”, “Кто работает в отделе продаж?”, “Можете ли вы показать мне N наиболее успешных продавцов?”. Обратите внимание, что они не спрашивают: “Сколько людей работает в отделе, значение DEPTNO которого равно 30?”. Они не используют эти ключевые значения, а вместо этого применяют читабельное для человека название отдела. Таким образом, в итоге они запускают запросы следующего вида:

```
select count(*)
from emp, dept
where emp.deptno = dept.deptno
and dept.dname = 'SALES'
/
select emp.*
from emp, dept
where emp.deptno = dept.deptno
and dept.dname = 'SALES'
/
```

Эти запросы почти обязательно должны получать доступ к таблицам DEPT и EMP, используя обычные индексы. Мы могли бы применять индекс на DEPT.DNAME для поиска строки (строк) в таблице SALES и извлечения значения DEPTNO для SALES, а затем использовать индекс на EMP.DEPTNO для нахождения соответствующих строк. Тем не менее, применяя битовый индекс соединения, всего перечисленного можно избежать. Битовый индекс соединения позволяет проиндексировать столбец DEPT.DNAME, но указывать он будет не на таблицу DEPT, а на EMP. Наличие возможности индексировать атрибуты из других таблиц — довольно основополагающая концепция, которая может изменить способ реализации модели данных в системе генерации отчетов. По существу вы получаете мед, а также возможность есть его ложкой. Нормализованная структура данных может оставаться незатронутой, но при этом появляются преимущества денормализации.

Создадим для этого примера индекс:

```

EODA@ORA12CR1> create bitmap index emp_bm_idx
 2 on emp( d.dname )
 3 from emp e, dept d
 4 where e.deptno = d.deptno
 5 /

```

Index created.

Индекс создан.

Обратите внимание, что начало оператора CREATE INDEX выглядит “нормальным” и создает для таблицы индекс EMP_BM_IDX. Но после этого появляются отличия от “нормального” оператора. Мы видим ссылку на столбец в таблице DEPT: D.DNAME. Мы наблюдаем конструкцию FROM, которая делает этот оператор CREATE INDEX напоминающим запрос. Мы имеем условие соединения между несколькими таблицами. Приведенный оператор CREATE INDEX индексирует столбец DEPT.DNAME, но в контексте таблицы EMP. Если мы зададим упомянутые ранее вопросы, то обнаружим, что база данных вообще не обращается к таблице DEPT, и делать ей это не нужно, т.к. столбец DNAME теперь присутствует в индексе, указывая на строки в таблице EMP. В целях иллюстрации мы сделаем таблицы EMP и DEPT большими (чтобы оптимизатор не посчитал их настолько малыми, что выбрал бы полное сканирование вместо использования индексов):

```

EODA@ORA12CR1> begin
 2 dbms_stats.set_table_stats( user, 'EMP',
 3                             numRows => 1000000, numblks => 300000 );
 4 dbms_stats.set_table_stats( user, 'DEPT',
 5                             numRows => 100000, numblks => 30000 );
 6 dbms_stats.delete_index_stats( user, 'EMP_BM_IDX' );
 7 end;
 8 /

```

PL/SQL procedure successfully completed.

Процедура PL/SQL успешно завершена.

На заметку! Вас может удивить, зачем выше была вызвана процедура DELETE_INDEX_STATE. Причина в том, что в Oracle 10g и последующих версиях оператор CREATE INDEX автоматически выполняет сбор статистики (COMPUTE STATISTICS) после создания индекса. Следовательно, в этом случае мы “обманули” базу данных Oracle — ей кажется, что она видит таблицу с 1 000 000 строк и крошечным индексом (а на самом деле таблица содержит всего 14 строк). Статистика индекса была точной, а статистика таблицы — “сфабрикованной”. Понадобилось также “подделать” и статистику индекса — иначе перед индексацией пришлось бы загружать в таблицу 1 000 000 записей.

Теперь запустим наши запросы:

```

EODA@ORA12CR1> set autotrace traceonly explain
EODA@ORA12CR1> select count(*)
 2 from emp, dept
 3 where emp.deptno = dept.deptno
 4 and dept.dname = 'SALES'
 5 /

```

Execution Plan

Plan hash value: 2538954156

| Id | Operation | Name | Rows | Bytes | Cost (%CPU) |
|-----|---------------------------|------------|------|-------|-------------|
| 0 | SELECT STATEMENT | | 1 | 3 | 7 (0) |
| 1 | SORT AGGREGATE | | 1 | 3 | |
| 2 | BITMAP CONVERSION COUNT | | 250K | 732K | 7 (0) |
| * 3 | BITMAP INDEX SINGLE VALUE | EMP_BM_IDX | | | |

Predicate Information (identified by operation id):

3 - access("EMP"."SYS_NC00009\$"='SALES')

Как видите, для ответа на этот конкретный вопрос не понадобилось обращаться ни к EMP, ни к DEPT — полный ответ был получен из самого индекса. Вся информация, нужная для ответа на вопрос, оказалась доступной в структуре индекса.

Более того, мы можем пропустить доступ к таблице DEPT и, применив индекс на таблице EMP, который включает необходимые данные из DEPT, получить непосредственный доступ к требуемым строкам:

```
EODA@ORA12CR1> select emp.*
2   from emp, dept
3  where emp.deptno = dept.deptno
4     and dept.dname = 'SALES'
5  /
```

Execution Plan

Plan hash value: 4261295901

| Id | Operation | Name | Rows | Bytes |
|-----|-------------------------------------|------------|-------|-------|
| 0 | SELECT STATEMENT | | 10000 | 849K |
| 1 | TABLE ACCESS BY INDEX ROWID BATCHED | EMP | 10000 | 849K |
| 2 | BITMAP CONVERSION TO ROWIDS | | | |
| * 3 | BITMAP INDEX SINGLE VALUE | EMP_BM_IDX | | |

Predicate Information (identified by operation id):

3 - access("EMP"."SYS_NC00009\$"='SALES')

С битовыми индексами соединений связано одно предварительное требование. Условие соединения должно производить соединение по первичному или уникаль-

ному ключу в другой таблице. В предыдущем примере DEPT.DEPTNO является первичным ключом таблицы DEPT, и если первичный ключ отсутствует, то возникнет ошибка:

```

EODA@ORA12CR1> create bitmap index emp_bm_idx
  2 on emp( d.dname )
  3 from emp e, dept d
  4 where e.deptno = d.deptno
  5 /
from emp e, dept d
      *
```

ERROR at line 3:
ORA-25954: missing primary key or unique constraint on dimension
ОШИБКА в строке 3:
ORA-25954: отсутствует первичный ключ или уникальное ограничение на измерении

Заключительные соображения по поводу битовых индексов

При наличии сомнений испытывайте битовые индексы (разумеется, в системе, отличной от OLTP). Добавить битовый индекс к таблице (или набор таких индексов) и посмотреть, что это дает, очень просто. К тому же битовые индексы обычно создаются намного быстрее, чем индексы со структурой В-дерева. Экспериментирование — лучший способ увидеть, насколько они подходят в конкретной среде. Меня часто спрашивают: как определяется низкая кардинальность? Четкого ответа на этот вопрос не существует. Иногда это 3 значения из 100 000, а иногда — 10 000 значений из 1 000 000. Низкая кардинальность вовсе не подразумевает однородные количества отличающихся значений. Только путем экспериментирования можно определить, насколько полезен битовый индекс в конкретном приложении. В общем случае, если есть крупная среда с данными, предназначенными главным образом для чтения, и множеством нерегламентированных запросов, то набор битовых индексов может оказаться в точности тем, что необходимо.

Индексы на основе функций

Индексы на основе функций появились в версии Oracle 8.1.5. Теперь они являются средством редакции Standard, хотя в выпусках, предшествовавших Oracle9i Release 2, они были доступны только в редакции Enterprise.

Индексы на основе функций предоставляют возможность индексировать вычисляемые столбцы и затем использовать такие индексы в запросах. Коротко говоря, это средство позволяет организовать нечувствительный к регистру поиск или сортировку, поиск по сложным уравнениям, а также эффективное расширение языка SQL за счет реализации собственных функций и операций для последующего поиска по ним. Вы будете применять индекс на основе функций по многим причинам, главные из которых описаны ниже.

- Эти индексы легко реализовать и немедленно получать значение.
- Эти индексы могут обеспечить повышение скорости работы существующих приложений без изменения их логики или запросов.

В последующих разделах приводятся подходящие случаи реализации индексов на основе функций.

Простой пример индекса на основе функции

Рассмотрим пример. Мы хотим выполнять нечувствительный к регистру поиск по столбцу ENAME таблицы EMP. До появления индексов на основе функций к решению этой задачи нужно было подходить в совершенно иной манере. В таблицу EMP пришлось бы добавить дополнительный столбец под названием, скажем, UPPER_ENAME. Этот столбец должен был сопровождаться триггером базы данных для операций INSERT и UPDATE; триггер бы просто устанавливал NEW.UPPER_NAME := UPPER(:NEW.ENAME). Упомянутый дополнительный столбец должен был быть проиндексированным. Теперь, благодаря индексам на основе функций, необходимость в наличии дополнительного столбца отпадает.

Начнем с создания копии демонстрационной таблицы EMP из схемы SCOTT и помещения в нее некоторых данных:

```
EODA@ORA12CR1> create table emp
2 as
3 select *
4   from scott.emp
5  where 1=0;
Table created.
Таблица создана.

EODA@ORA12CR1> insert into emp
2  (empno,ename,job,mgr,hiredate,sal,comm,deptno)
3  select rownum empno,
4         initcap(substr(object_name,1,10)) ename,
5         substr(object_type,1,9) JOB,
6         rownum MGR,
7         created hiredate,
8         rownum SAL,
9         rownum COMM,
10        (mod(rownum,4)+1)*10 DEPTNO
11  from all_objects
12  where rownum < 10000;
9999 rows created.
9999 строк создано.
```

Далее создадим индекс по значению функции UPPER, вызванной на столбце ENAME, по существу создавая нечувствительный к регистру индекс:

```
EODA@ORA12CR1> create index emp_upper_idx on emp(upper(ename));
Index created.
Индекс создан.
```

Наконец, мы проанализируем таблицу, поскольку, как было отмечено ранее, это необходимо для того, чтобы заставить оптимизатор СВО использовать индексы на основе функций. Начиная с версии Oracle 10g, этот шаг формально не требуется, т.к. оптимизатор СВО применяется по умолчанию, а динамическая выборка соберет нужную информацию, но сбор статистики является более корректным подходом.

```
EODA@ORA12CR1> begin
2   dbms_stats.gather_table_stats
3   (user,'EMP',cascade=>true);
4 end;
5 /
```

PL/SQL procedure successfully completed.

Процедура PL/SQL успешно завершена.

Теперь у нас есть индекс по значению столбца в верхнем регистре. Любое приложение, которое уже выдает нечувствительные к регистру запросы, подобные показанному ниже, будет использовать этот индекс, обеспечивая привносимый им рост производительности:

```
EOA@ORA12CR1> set autotrace traceonly explain
```

```
EOA@ORA12CR1> select *
```

```
2   from emp
```

```
3   where upper(ename) = 'KING';
```

Execution Plan

Plan hash value: 3831183638

| Id | Operation | Name | Rows | Bytes |
|----|-------------------------------------|---------------|------|-------|
| 0 | SELECT STATEMENT | | 2 | 110 |
| 1 | TABLE ACCESS BY INDEX ROWID BATCHED | EMP | 2 | 110 |
| 2 | INDEX RANGE SCAN | EMP_UPPER_IDX | 2 | |

Predicate Information (identified by operation id):

2 - access(UPPER("ENAME")='KING')

Перед тем, как это средство стало доступным, каждая строка в таблице EMP должна была быть просмотрена, переведена в верхний регистр и подвергнута сравнению. По контрасту с этим при наличии индекса на UPPER(ENAME) запрос принимает константу KING для индекса, сканирует по диапазону небольшой объем данных и обращается к таблице по идентификатору строки с целью извлечения данных. Это работает очень быстро.

Такой рост производительности наиболее заметен в случае индексации столбцов на основе пользовательских функций. В Oracle 7.1 была добавлена возможность написания пользовательских функций на языке SQL, так что можно было поступать следующим образом:

```
SQL> select my_function(ename)
```

```
2   from emp
```

```
3   where some_other_function(empno) > 10
```

```
4   /
```

Это было замечательным нововведением, поскольку позволило эффективно расширять язык SQL за счет включения специфических для приложения функций. Однако, к сожалению, производительность предыдущего запроса временами несколько разочаровывала. Пусть таблица EMP содержит 1000 строк. Тогда во время запроса функция SOME_OTHER_FUNCTION должна быть выполнена 1000 раз — по одному разу на строку. Вдобавок, если предположить, что выполнение функции за-

нимает одну сотую секунды, то такой сравнительно простой запрос продлится, по меньшей мере, десять секунд.

Давайте рассмотрим реальный пример, в котором реализуем модифицированную процедуру SOUNDEX на PL/SQL. Дополнительно мы будем применять в нашей процедуре глобальную переменную пакета в качестве счетчика; это позволит выполнять запросы, которые обращаются к функции MY_SOUNDEX, и видеть, сколько раз она была вызвана:

```

EODA@ORA12CR1> create or replace package stats
2 as
3     cnt number default 0;
4 end;
5 /
Package created.
Пакет создан.

EODA@ORA12CR1> create or replace
2 function my_soundex( p_string in varchar2 ) return varchar2
3 deterministic
4 as
5     l_return_string varchar2(6) default substr( p_string, 1, 1 );
6     l_char          varchar2(1);
7     l_last_digit    number default 0;
8
9     type vcArray is table of varchar2(10) index by binary_integer;
10    l_code_table     vcArray;
11
12 begin
13     stats.cnt := stats.cnt+1;
14
15     l_code_table(1) := 'BPFV';
16     l_code_table(2) := 'CSKGJQXZ';
17     l_code_table(3) := 'DT';
18     l_code_table(4) := 'L';
19     l_code_table(5) := 'MN';
20     l_code_table(6) := 'R';
21
22
23     for i in 1 .. length(p_string)
24     loop
25         exit when (length(l_return_string) = 6);
26         l_char := upper(substr( p_string, i, 1 ));
27
28         for j in 1 .. l_code_table.count
29         loop
30             if (instr(l_code_table(j), l_char ) > 0 AND j <> l_last_digit)
31             then
32                 l_return_string := l_return_string || to_char(j,'fm9');
33                 l_last_digit := j;
34             end if;
35         end loop;
36     end loop;
37

```

```

38     return rpad( l_return_string, 6, '0' );
39 end;
40 /

```

Function created.

Функция создана.

Обратите внимание на использование в этой функции ключевого слова DETERMINISTIC. Оно заявляет о том, что предыдущая функция при получении одних и тех же входных данных всегда будет возвращать в точности те же самые выходные данные. Это необходимо для создания индекса, основанного на пользовательской функции. Мы должны сообщить Oracle о том, что функция является детерминированной (DETERMINISTIC) и будет возвращать согласованный результат для одних и тех же входных данных. Мы указываем Oracle, что этой функции следует доверять в том, что она будет возвращать одно и то же значение, от вызова к вызову, при получении тех же самых входных данных. Если это не так, то при доступе к данным через индекс мы получали бы разные ответы по сравнению с полным сканированием таблицы. Параметр DETERMINISTIC подразумевает, например, что мы не можем создать индекс на основе функции DBMS_RANDOM.RANDOM, представляющей собой генератор случайных чисел. Ее результаты недетерминированы; при передаче одних и тех же входных данных мы получим случайное выходное значение. С другой стороны, встроенная SQL-функция UPPER, которая применялась в первом примере, является детерминированной, поэтому можно создавать индекс по значению столбца, приведенному к верхнему регистру.

Теперь, когда мы располагаем функцией MY_SOUNDINDEX, давайте посмотрим, как она работает без индекса. Воспользуемся созданной ранее таблицей EMP, содержащей 10 000 строк:

```

EODA@ORA12CR1> set autotrace on explain
EODA@ORA12CR1> variable cpu number
EODA@ORA12CR1> exec :cpu := dbms_utility.get_cpu_time
PL/SQL procedure successfully completed.
EODA@ORA12CR1> select ename, hiredate
2     from emp
3  where my_soundindex(ename) = my_soundindex('Kings')
4  /

```

```

ENAME      HIREDATE
-----
Ku$_Chunk_ 17-DEC-13
Ku$_Chunk_ 17-DEC-13
Ku$_Chunk_ 17-DEC-13
Ku$_Chunk_ 17-DEC-13

```

Execution Plan

Plan hash value: 3956160932

| Id | Operation | Name | Rows | Bytes | Cost (%CPU) | Time |
|-----|-------------------|------|------|-------|-------------|----------|
| 0 | SELECT STATEMENT | | 100 | 1900 | 24 (9) | 00:00:01 |
| * 1 | TABLE ACCESS FULL | EMP | 100 | 1900 | 24 (9) | 00:00:01 |

Predicate Information (identified by operation id):

```

1 - filter("MY_SOUNDEX"("ENAME")="MY_SOUNDEX"('Kings'))
EODA@ORA12CR1> set autotrace off
EODA@ORA12CR1> begin
2     dbms_output.put_line
3     ( 'cpu time = ' || round((dbms_utility.get_cpu_time-:cpu)/100,2) );
4     dbms_output.put_line( 'function was called: ' || stats.cnt );
5 end;
6 /

cpu time = .2
function was called: 9916
PL/SQL procedure successfully completed.

время ЦП = .2
функция была вызвана: 9916
Процедура PL/SQL успешно завершена.
```

Здесь видно, что этот запрос занял две десятых доли секунды и должен бы выполнить полное сканирование таблицы. Функция MY_SOUNDEX была вызвана почти 10 000 раз (согласно нашему счетчику). Однако хотелось бы, чтобы это происходило реже.

На заметку! В более старых выпусках (до Oracle 10g Release 2) функция вызывалась бы намного больше раз, чем мы наблюдали выше. На самом деле она бы вызывалась примерно 20 000 раз — дважды для каждой строки! В Oracle 10g Release 2 и последующих версиях для сокращения количества раз вызова функции применяется подсказка DETERMINISTIC.

Давайте посмотрим, каким образом индексация функции может ускорить работу. Первое, что мы сделаем — это создадим индекс, как показано ниже:

```

EODA@ORA12CR1> create index emp_soundex_idx on
2 emp( substr(my_soundex(ename),1,6) )
3 /

Index created.
Индекс создан.
```

В этой команде CREATE INDEX интересно обратить внимание на использование функции SUBSTR. Это связано с тем, что мы индексируем функцию, которая возвращает строку. Если бы мы индексировали функцию, возвращающую число или дату, то вызов SUBSTR не понадобился бы. Необходимость в применении SUBSTR объясняется тем, что пользовательская функция возвращает тип VARCHAR2(4000). Такое значение может оказаться слишком большим для индексации — записи индекса должны уместиться примерно в три четверти размера блока. Если мы попытаемся обойтись без SUBSTR, то получим вот что (в табличном пространстве с размером блока 4 Кбайт):

```

EODA@ORA12CR1> create index emp_soundex_idx on
2 emp( my_soundex(ename) ) tablespace ts4k;
emp( my_soundex(ename) ) tablespace ts4k
*
```

```
ERROR at line 2:
ORA-01450: maximum key length (3118) exceeded
ОШИБКА в строке 2:
ORA-01450: превышена максимальная длина ключа (3118)
```

Это не означает, что индекс действительно содержит ключи такой длины, но с точки зрения базы данных он *может* их иметь. Однако база данных понимает вызов функции SUBSTR. Она видит, что SUBSTR были переданы аргументы 1 и 6, и знает, что наибольшее возвращаемое значение не превысит шести символов, поэтому разрешает создать такой индекс. Проблема размера индекса может возникать, особенно в случае сцепленных индексов. Вот пример для табличного пространства с размером блока 8 Кбайт:

```
EODA@ORA12CR1> create index emp_soundex_idx on
  2 emp( my_soundex(ename), my_soundex(job) );
emp( my_soundex(ename), my_soundex(job) )
*
```

```
ERROR at line 2:
ORA-01450: maximum key length (6398) exceeded
ОШИБКА в строке 2:
ORA-01450: превышена максимальная длина ключа (6398)
```

База данных считает, что максимальный размер ключа равен 8000 байтов и снова отказывается принять оператор CREATE. Таким образом, чтобы проиндексировать пользовательскую функцию, возвращающую строку, мы должны ограничить возвращаемый тип в операторе CREATE INDEX. В этом примере, зная, что MY_SOUNDSEX возвращает максимум 6 символов, мы вырезаем первые 6 символов.

Теперь все готово к тестированию производительности таблицы с индексом. Мы хотели бы отследить влияние индекса на операторы INSERT, а также ускорение выполнения операторов SELECT. В сценарии без индекса запросы выполняются более одной секунды, а если запустить SQL_TRACE и TKPROF во время вставок, то обнаружится, что без индекса вставка 9 999 записей занимает примерно 0,3 секунды:

```
insert into emp NO_INDEX
(empno,ename,job,mgr,hiredate,sal,comm,deptno)
select rownum empno,
       initcap(substr(object_name,1,10)) ename,
       substr(object_type,1,9) JOB,
       rownum MGR,
       created hiredate,
       rownum SAL,
       rownum COMM,
       (mod(rownum,4)+1)*10 DEPTNO
from all_objects
where rownum < 10000
```

| call | count | cpu | elapsed | disk | query | current | rows |
|---------|-------|------|---------|------|-------|---------|------|
| Parse | 1 | 0.08 | 0.08 | 0 | 0 | 0 | 0 |
| Execute | 1 | 0.15 | 0.22 | 0 | 2745 | 13763 | 9999 |
| Fetch | 0 | 0.00 | 0.00 | 0 | 0 | 0 | 0 |
| total | 2 | 0.23 | 0.30 | 0 | 2745 | 13763 | 9999 |

Но с индексом эта вставка потребует около 0,57 секунды:

| call | count | cpu | elapsed | disk | query | current | rows |
|---------|-------|------|---------|------|-------|---------|------|
| Parse | 1 | 0.07 | 0.07 | 0 | 0 | 0 | 0 |
| Execute | 1 | 0.39 | 0.49 | 131 | 2853 | 23313 | 9999 |
| Fetch | 0 | 0.00 | 0.00 | 0 | 0 | 0 | 0 |
| total | 2 | 0.46 | 0.57 | 131 | 2853 | 23313 | 9999 |

Разница объясняется накладными расходами по управлению новым индексом на основе функции MY_SOUNDEX — производительность снижается как просто по причине наличия индекса (индекс любого типа будет влиять на производительность вставки), так и потому, что индекс должен 9 999 раз вызвать хранимую процедуру.

Теперь, чтобы протестировать запрос, мы запустим его повторно:

```
EODA@ORA12CR1> exec stats.cnt := 0
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.

EODA@ORA12CR1> exec :cpu := dbms_utility.get_cpu_time
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.

EODA@ORA12CR1> set autotrace on explain
EODA@ORA12CR1> select ename, hiredate
  2   from emp
  3   where substr(my_soundex(ename),1,6) = my_soundex('Kings')
  4   /
```

```
ENAME          HIREDATE
-----
Ku$_Chunk_     17-DEC-13
Ku$_Chunk_     17-DEC-13
Ku$_Chunk_     17-DEC-13
Ku$_Chunk_     17-DEC-13
```

Plan hash value: 1897478402

```
-----
| Id | Operation                                | Name                | Rows | Bytes |
Cost (%CPU)...
-----
|  0 | SELECT STATEMENT                        |                     |  100 |  3300 |
12 ...
|  1 |  TABLE ACCESS BY INDEX ROWID BATCHED | EMP                 |  100 |  3300 |
12 ...
|*  2 |    INDEX RANGE SCAN                    | EMP_SOUNDEX_IDX    |   40 |       |
1 ...
-----
```

Predicate Information (identified by operation id):

```
2 - access(SUBSTR("EODA"."MY_SOUNDEX"("ENAME"),1,6)="MY_
SOUNDEX"('Kings'))
```

```
EODA@ORA12CR1> set autotrace off
```

```

EODA@ORA12CR1> begin
2   dbms_output.put_line
3   ( 'cpu time = ' || round((dbms_utility.get_cpu_time-:cpu)/100,2) );
4   dbms_output.put_line( 'function was called: ' || stats.cnt );
5 end;
6 /
cpu time = .01
function was called: 1

PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.

```

Если мы сравним эти два примера (без индексации и с индексацией), то заметим, что вставка в индексированную таблицу потребовала более чем вдвое больше времени. Однако выборка вместо двух десятых секунды стала выполняться фактически мгновенно. Ниже перечислены важные замечания относительно рассмотренных примеров.

- Вставка 9 999 записей заняла приблизительно в два раза больше времени. Индексация по пользовательской функции обязательно повлияет на производительность вставок и некоторых обновлений. Вы должны понимать, что любой индекс, конечно же, влияет на производительность. Например, я прогнал простой тест без функции MY_SOUNDEX, просто проиндексировав сам столбец ENAME. Это привело к тому, что выполнение вставки длилось около одной секунды — функция PL/SQL не несет ответственности за все накладные расходы. Поскольку большинство приложений вставляют и обновляют одиночные записи, а каждая строка требует менее 1/10 000 секунды на вставку, в типичном приложении вы этого даже не заметите. Из-за того, что мы вставляем строку однократно, то платим цену за выполнение функции на столбце только однажды, а не тысячи раз, когда запрашиваем данные.
- В то время как вставка протекает вдвое медленнее, запрос работает во много раз быстрее. Он выполняет функцию MY_SOUNDEX несколько раз вместо почти 20 000 раз. Разница в производительности запроса измерима и довольно велика. К тому же по мере роста таблицы запросы с полным сканированием будут требовать все больше и больше времени на выполнение. В случае увеличения размера таблицы запросы на основе индексов будут всегда выполняться с приблизительно одинаковыми характеристиками производительности.
- Мы должны использовать в запросе функцию SUBSTR. Это не настолько изящно, как кодирование WHERE MY_SOUNDEX(ename) = MY_SOUNDEX('King'), но проблему легко обойти, что вскоре будет показано.

Итак, производительность вставки снизилась, но запросы выполняются noticeably быстро. Выигрыш за небольшое сокращение производительности вставки/обновления огромен. Кроме того, если вы никогда не обновляете столбцы, участвующие в вызове функции MY_SOUNDEX, то обновления не страдают вообще (MY_SOUNDEX вызывается, только когда столбец ENAME *модифицируется* и его значение изменяется).

Давайте посмотрим, каким образом сделать так, чтобы запрос не обязан был вызывать функцию SUBSTR. Применение функции SUBSTR может быть чревато ошибками — конечные пользователи должны знать о необходимости вызова SUBSTR для

получения подстроки с первого по шестой символ. Если они укажут другой размер, то индекс использоваться не будет. Вдобавок понадобится управлять на сервере количеством байтов для индексации. Это позволит позже изменить реализацию функции MY_SOUNDX на возвращение семи символов вместо шести, если возникнет такое желание. В Oracle Database 11g Release 1 и последующих версиях мы можем довольно легко скрыть вызов функции SUBSTR с помощью виртуального столбца — или же представления в любом выпуске:

```
EODA@ORA12CR1> create or replace view emp_v
2 as
3 select ename, substr(my_soundex(ename),1,6) ename_soundex, hiredate
4 from emp
5 /
```

View created.

Представление создано.

```
EODA@ORA12CR1> exec stats.cnt := 0;
PL/SQL procedure successfully completed.
```

Процедура PL/SQL успешно завершена.

```
EODA@ORA12CR1> exec :cpu := dbms_utility.get_cpu_time
PL/SQL procedure successfully completed.
```

Процедура PL/SQL успешно завершена.

```
EODA@ORA12CR1> select ename, hiredate
2 from emp_v
3 where ename_soundex = my_soundex('Kings')
4 /
```

| ENAME | HIREDATE |
|-------------|-----------|
| Ku\$ Chunk_ | 17-DEC-13 |
| Ku\$ Chunk_ | 17-DEC-13 |
| Ku\$ Chunk_ | 17-DEC-13 |
| Ku\$ Chunk_ | 17-DEC-13 |

```
EODA@ORA12CR1> begin
2 dbms_output.put_line
3 ( 'cpu time = ' || round((dbms_utility.get_cpu_time-:cpu)/100,2) );
4 dbms_output.put_line( 'function was called: ' || stats.cnt );
5 end;
6 /
```

cpu time = .01

function was called: 1

```
PL/SQL procedure successfully completed.
```

Процедура PL/SQL успешно завершена.

Мы имеем дело с той же разновидностью плана выполнения запроса, что и в случае базовой таблицы. Все, что мы предприняли здесь — это скрыли вызов функции SUBSTR(F(X), 1, 6) внутри самого представления. Оптимизатор по-прежнему распознает, что этот виртуальный столбец в действительности является индексированным столбцом, так что он поступает правильно. Мы наблюдаем такое же улучшение производительности и тот же самый план выполнения запроса. Применение этого представления так же хорошо, как использование базовой таблицы — и даже лучше, потому что оно скрывает сложность и позволяет позже изменить длину подстроки, получаемой с помощью SUBSTR.

В Oracle 11g Release 1 и последующих версиях имеется еще один вариант для реализации. Вместо применения представления с “виртуальным столбцом” мы можем использовать настоящий виртуальный столбец. Применение этого средства предусматривает удаление существующего индекса на основе функции:

```
EODA@ORA12CR1> drop index emp_soundex_idx;
Index dropped.
Индекс удален.
```

Затем добавим в таблицу виртуальный столбец и проиндексируем его:

```
EODA@ORA12CR1> alter table emp
2 add
3   ename_soundex as
4   (substr(my_soundex(ename),1,6))
5 /
Table altered.
Таблица изменена.

EODA@ORA12CR1> create index emp_soundex_idx
2 on emp(ename_soundex);
Index created.
Индекс создан.
```

Теперь можно просто запрашивать базовую таблицу — никакой дополнительный уровень представления здесь не задействован:

```
EODA@ORA12CR1> exec stats.cnt := 0;
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.

EODA@ORA12CR1> exec :cpu := dbms_utility.get_cpu_time;
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.

EODA@ORA12CR1> select ename, hiredate
2   from emp
3   where ename_soundex = my_soundex('Kings')
4   /
```

| ENAME | HIREDATE |
|-------------|-----------|
| Ku\$_Chunk_ | 17-DEC-13 |
| Ku\$_Chunk_ | 17-DEC-13 |
| Ku\$_Chunk_ | 17-DEC-13 |
| Ku\$_Chunk_ | 17-DEC-13 |

```
EODA@ORA12CR1> begin
2   dbms_output.put_line
3   ( 'cpu time = ' || round((dbms_utility.get_cpu_time-:cpu)/100,2) );
4   dbms_output.put_line( 'function was called: ' || stats.cnt );
5 end;
6 /
cpu time = 0
function was called: 1
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.
```

Индексация только некоторых строк

В дополнение к прозрачному содействию нашим запросам, которые используют встроенные функции наподобие UPPER, LOWER и т.д., индексы на основе функций могут применяться для выборочной индексации только определенных строк таблицы. Позже вы узнаете, что индексы со структурой В-дерева не содержат записи для ключей со значением NULL. То есть, если имеется индекс I на таблице T (как показано ниже) и есть строка, в которой A и B содержат NULL, то в структуре индекса никаких записей не будет:

```
Create index I on t(a,b);
```

Это пригодится в ситуации, когда нужно индексировать только некоторые строки в таблице.

Рассмотрим крупную таблицу со столбцом NOT NULL по имени PROCESSED_FLAG, который может принимать одно из двух значений, Y или N, со стандартным значением N. Новые строки добавляются со значением N, указывающим на то, что они не обработаны, а после их обработки значение PROCESSED_FLAG обновляется на Y. Мы хотим проиндексировать этот столбец, чтобы иметь возможность быстро извлекать записи со значением N, но в таблице присутствуют миллионы строк, большинство из которых в итоге получают значение Y для столбца PROCESSED_FLAG. Результирующий индекс со структурой В-дерева окажется огромным, и затраты на его обслуживание при обновлении значения с N на Y будут высокими. Данная таблица выглядит как неплохой кандидат на создание битового индекса (в конце концов, кардинальность является низкой). Однако это транзакционная система, и многие пользователи будут вставлять записи одновременно с обработкой строк, имеющих значение N в столбце PROCESSED_FLAG, а ранее мы выяснили, что битовые индексы не очень хорошо подходят при наличии параллельных модификаций. Учитывая необходимость в постоянном обновлении значений N на Y, битовые индексы можно исключить из обсуждения, т.к. они сделают весь процесс полностью последовательным.

Итак, мы хотели бы проиндексировать только те записи, которые нас интересуют (со значением N). Мы посмотрим, как этого добиться с помощью индексов на основе функций, но сначала давайте взглянем на то, что произойдет, если просто воспользоваться обычным индексом со структурой В-дерева. Применив стандартный сценарий BIG_TABLE, описанный в разделе “Настройка среды” в самом начале книги, мы обновим столбец TEMPORARY, переключая значения Y на N, а N на Y:

```
EODA@ORA12CR1> update big_table set temporary =
decode(temporary, 'N', 'Y', 'N');
1000000 rows updated.
1000000 строк обновлено.
```

Теперь проверим соотношение между значениями Y и N:

```
EODA@ORA12CR1> select temporary, cnt,
2         round( (ratio_to_report(cnt) over ()) * 100, 2 ) rtr
3   from (
4 select temporary, count(*) cnt
5   from big_table
6  group by temporary
7   )
8 /
```

| T | CNT | RTR |
|---|--------|-------|
| - | ----- | ----- |
| Y | 998728 | 99.87 |
| N | 1272 | .13 |

Как видите, из 1 000 000 строк таблицы должны быть проиндексированы только около одной пятой процента. Если мы воспользуемся обычным индексом на столбце TEMPORARY (который в этом примере играет роль столбца PROCESS_FLAG), то обнаружим, что индекс содержит 1 000 000 записей, потребляет приблизительно 14 Мбайт дискового пространства и имеет высоту 3:

```
EODA@ORA12CR1> create index processed_flag_idx
  2 on big_table(temporary);
Index created.
Индекс создан.
```

```
EODA@ORA12CR1> analyze index processed_flag_idx
  2 validate structure;
Index analyzed.
Индекс проанализирован.
```

```
EODA@ORA12CR1> select name, btree_space, lf_rows, height from index_
stats;
```

| NAME | BTREE_SPACE | LF_ROWS | HEIGHT |
|--------------------|-------------|---------|--------|
| ----- | ----- | ----- | ----- |
| PROCESSED_FLAG_IDX | 14528892 | 1000000 | 3 |

Любое извлечение данных через этот индекс потребует трех операций ввода-вывода, чтобы добраться до листового блока. Такой индекс является не только широким, но и высоким. Для получения первой необработанной записи придется выполнить, по меньшей мере, четыре операции ввода-вывода (три по индексу и одну по таблице).

Как можно изменить все это? Нужно сделать так, чтобы индекс был меньше, и чтобы его было легче обслуживать (с меньшими накладными расходами во время выполнения при обновлениях). Создадим индекс на основе функции, который позволит просто написать функцию, возвращающую NULL, когда конкретная строка не должна быть включена в индекс, и отличающееся от NULL значение — когда должна. Например, поскольку мы заинтересованы только в строках со значением N, поступим следующим образом:

```
EODA@ORA12CR1> drop index processed_flag_idx;
Index dropped.
Индекс удален.
```

```
EODA@ORA12CR1> create index processed_flag_idx
  2 on big_table( case temporary when 'N' then 'N' end );
Index created.
Индекс создан.
```

```
EODA@ORA12CR1> analyze index processed_flag_idx validate structure;
Index analyzed.
Индекс проанализирован.
```

```
EOGA@ORA12CR1> select name, btree_space, lf_rows, height from index_stats;
```

| NAME | BTREE_SPACE | LF_ROWS | HEIGHT |
|--------------------|-------------|---------|--------|
| ----- | ----- | ----- | ----- |
| PROCESSED_FLAG_IDX | 32016 | 1272 | 2 |

Разница заметна: индекс теперь занимает около 40 Кбайт вместо 14 Мбайт. Его высота также уменьшилась. В случае применения этого индекса мы будем выполнять на одну операцию ввода-вывода меньше по сравнению с предыдущим более высоким индексом.

Реализация выборочной уникальности

Другой удобный прием с индексами на основе функций предусматривает их использование для принудительного применения сложных ограничений определенных типов. Например, предположим, что есть таблица с версионной информацией, подобная таблице проектов. Проекты имеют одно из двух состояний: `ACTIVE` или `INACTIVE`. Нужно обеспечить выполнение правила, которое заключается в том, что активные проекты должны иметь уникальное имя, а неактивные проекты — нет. То есть может существовать только один активный “проект X”, но неактивных проектов X может быть произвольное количество.

Первая реакция разработчика, который видит такое требование, обычно выглядит так: “Мы просто запустим запрос, чтобы посмотреть, существует ли активный проект X, и если нет, то создадим его”. Если вы читали главу 7, то понимаете, что такая элементарная реализация не сможет работать в многопользовательской среде. Если два пользователя попытаются создать новый активный проект X в одно и то же время, то это получится у обоих. Создание проекта X необходимо сериализовать, но единственный способ сделать это заключается в блокировке целой таблицы проектов (не особенно содействующий параллелизму подход) либо использовании индекса на основе функции и предоставлении базе данных возможности делать все самостоятельно.

Опираясь на тот факт, что можно создавать индексы на основе функций, на то, что записи `NULL` не сохраняются в индексах со структурой B-дерева, а также на возможность создания уникального индекса, мы просто поступим так:

```
Create unique index active_projects_must_be_unique
On projects ( case when status = 'ACTIVE' then name end );
```

Это и все. Когда столбец состояния содержит значение `ACTIVE`, столбец `NAME` будет уникальным образом индексирован. Любая попытка создать активный проект с дублированным именем будет обнаружена и параллельный доступ к таблице вообще не пострадает.

Предостережение относительно ошибки ORA-01743

С индексами на основе функций связана одна особенность, о которой следует упомянуть. Если такой индекс создается на основе встроенной функции `TO_DATE`, то в некоторых случаях сделать это не удастся:

```
EOGA@ORA12CR1> create table t ( year varchar2(4) );
Table created.
Таблица создана.
```

```

EODA@ORA12CR1> create index t_idx on t( to_date(year,'YYYY') );
create index t_idx on t( to_date(year,'YYYY') )
*
```

```

ERROR at line 1:
ORA-01743: only pure functions can be indexed
ОШИБКА в строке 1:
ORA-01743: индексироваться могут только чистые функции
```

Выглядит довольно странно, поскольку *иногда* мы можем создать функцию, изменяющую TO_DATE, например:

```

EODA@ORA12CR1> create index t_idx on t( to_date('01'||year,'ММYYYY') );
Index created.
Индекс создан.
```

Сообщение, сопровождающее ошибку, не особенно проливает свет на причины:

```

EODA@ORA12CR1> !oerr ora 1743
01743, 00000, "only pure functions can be indexed"
// *Cause: The indexed function uses SYSDATE or the user environment.
// *Action: PL/SQL functions must be pure (RNDS, RNPS, WNDS, WNPS). SQL
//          expressions must not use SYSDATE, USER, USERENV(), or anything
//          else dependent on the session state. NLS-dependent functions
//          are OK.
// *Причина: Индексная функция использует SYSDATE или пользовательскую среду.
// *Действие: Функции PL/SQL должны быть чистыми (RNDS, RNPS, WNDS, WNPS).
//          Выражения SQL не должны применять SYSDATE, USER, USERENV()
//          или что-то еще, зависящее от состояния сеанса. Функции,
//          зависящие от NLS, разрешены.
```

Мы не используем SYSDATE. Мы не применяем пользовательскую среду (или все-таки применяем?). Мы не используем какие-либо функции PL/SQL и ничего, что задействовало бы состояние сеанса. Секрет кроется в указанном формате: YYYY. Этот формат для одних и тех же входных данных возвращает разные ответы в зависимости от того, для какого месяца производится вызов. Например, для любого времени в мае формат YYYY возвратит May 1, в июне — June 1 и т.д.:

```

EODA@ORA12CR1> select to_char( to_date('2015','YYYY'),
2                          'DD-Mon-YYYY HH24:MI:SS' )
3      from dual;
```

```

TO_CHAR(TO_DATE('200
-----
01-May-2015 00:00:00
```

Оказывается, что функция TO_DATE в случае применения с форматом YYYY является недетерминированной! Вот почему индекс не может быть создан: он будет работать корректно только для того месяца, в котором вы его создаете (либо вставляете или обновляете строку). Таким образом, это связано с пользовательской средой, которая включает саму текущую дату.

Чтобы использовать TO_DATE в индексе на основе функции, вы *должны* применять однозначный и детерминированный формат даты, который не зависит от текущего дня.

Заключительные соображения по поводу индексов на основе функций

Индексы на основе функций легко создавать и использовать, и они предоставляют немедленное значение. Они могут применяться для ускорения существующих приложений без изменения их логики или запросов. При этом могут наблюдаться существенные усовершенствования. Их можно использовать для предварительного вычисления сложных значений, не прибегая к триггерам. Вдобавок оптимизатор может более точно оценивать селективность, если выражения реализованы в индексе на основе функций. Индексы на основе функций можно применять для избирательной индексации только тех строк, которые представляют интерес, как было продемонстрировано ранее в примере с `PROCESSED_FLAG`. С помощью этого приема, в сущности, можно проиндексировать конструкцию `WHERE`. Наконец, индексы на основе функций могут использоваться для реализации определенного рода ограничений целостности: выборочной уникальности (например, поля `X`, `Y` и `Z` должны быть уникальными, когда истинно некоторое условие).

Индексы на основе функций будут влиять на производительность вставок и обновлений. Насколько это обстоятельство существенно для вас — решать вам. Если данные в основном вставляются и очень редко запрашиваются, то такие индексы могут оказаться неподходящим средством. С другой стороны, имейте в виду, что вставка строки обычно производится однажды, а ее запрос — тысячи раз. Снижение производительности при вставке (которое конечные пользователи, возможно, никогда не заметят) может быть тысячекратно перекрыто ускорением запросов. В общем, в этом случае преимущества значительно перевешивают недостатки.

Индексы предметной области

Индексы предметной области — это то, что в Oracle называется *расширяемой индексацией*. Они позволяют создавать собственные индексные структуры, работающие так же, как и индексы, предлагаемые Oracle. Когда кто-то выдает оператор `CREATE INDEX` с указанием вашего типа индекса, Oracle выполнит ваш код для генерации индекса. Если кто-то анализирует индекс для сбора статистики по нему, Oracle запустит ваш код для генерации статистики в формате, который вас интересует. Когда Oracle разбирает запрос и разрабатывает план его выполнения, который может использовать ваш индекс, будет запрошена стоимость выполнения этой функции для оценки эффективности разных планов. Короче говоря, индексы предметной области предоставляют возможность реализовать новый тип индекса, который пока еще не существует в базе данных. Например, если вы разрабатываете программное обеспечение, анализирующее графические изображения, которые хранятся в базе данных, и формируете сведения об этих изображениях, такие как найденные в них цвета, то можете создать собственный индекс на основе *изображений*. По мере добавления изображений в базу данных ваш код вызывается для извлечения цветов из них и сохранения их в каком-нибудь месте (там, где они должны храниться). Во время выполнения запроса, когда пользователь желает получить все изображения с конкретным цветом, Oracle обратится за ответом к вашему индексу.

Лучшим примером индекса предметной области является собственный *текстовый индекс* Oracle. Этот индекс применяется для обеспечения поиска по ключевым

словам в крупных текстовых элементах. Вы можете создать простой текстовый индекс следующим образом:

```
EODA@ORA12CR1> create index myindex on mytable(docs)
  2  indextype is ctxsys.context
  3  /
Index created.
Индекс создан.
```

Затем можно использовать текстовые операции, введенные в язык SQL создателями индекса этого типа:

```
select * from mytable where contains( docs, 'some words' ) > 0;
```

Текстовый индекс сможет даже отвечать на команды, подобные приведенным ниже:

```
EODA@ORA12CR1> begin
  2  dbms_stats.gather_index_stats( user, 'MYINDEX' );
  3  end;
  4  /
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.
```

Такой индекс будет приниматься во внимание оптимизатором во время выполнения, чтобы определить относительную стоимость применения текстового индекса по сравнению с другими индексами или полным сканированием таблицы. Во всем этом интересно то, что разработать подобный индекс мог бы любой. Реализация текстового индекса выполнена без *знаний внутреннего устройства ядра*. Он построен с использованием выделенного, документированного и открытого API-интерфейса. Ядро базы данных Oracle не осведомлено о том, как хранится текстовый индекс (API-интерфейс сохраняет его во многих таблицах базы данных для созданного индекса). Ядру базы данных Oracle ничего не известно об обработке, происходящей во время вставки новой строки. Текст Oracle на самом деле поддерживается приложением, которое построено поверх ядра, но полностью в него интегрировано. Для нас он выглядит похожим на любую другую функцию ядра базы данных Oracle, хотя в действительности это не так.

Лично у меня не возникало потребности в построении нового экзотического типа индексной структуры. Я видел примеры применения этого средства в основном в составе решений от независимых поставщиков, которые предлагали новаторские технологии индексации.

Я считаю, что наиболее интересным аспектом индексов предметной области является то, что они позволяют другим поставлять новые технологии индексации, которые я могу использовать в своих приложениях. Большинство людей никогда не применяют этот конкретный API-интерфейс для создания новых типов индексов, но многие из нас пользуются конечными результатами. Почти любое приложение, над которым я работал, имело какой-то ассоциированный с ним *текст*, *XML-код*, подлежащий обработке, или *графические изображения*, предназначенные для сохранения и категоризации. Набор функциональности Oracle Multimedia, реализованный с применением средства Application Domain Indexing (Индексация предметной области), предоставляет эти возможности. Со временем перечень доступных типов индексов растет. Мы рассмотрим текстовые индексы более подробно в следующей главе.

Невидимые индексы

В Oracle Database 11g и последующих версиях имеется возможность сделать индекс невидимым оптимизатору. Индекс является невидимым только в том смысле, что оптимизатор не будет его использовать при создании плана выполнения. Вы можете либо создать новый индекс как невидимый, либо изменить существующий индекс, превратив его в невидимый. Ниже мы создаем таблицу, загружаем ее данными, генерируем статистику и затем создаем невидимый индекс:

```

EODA@ORA12CR1> create table t(x int);
Table created.
Таблица создана.

EODA@ORA12CR1> insert into t select round(dbms_random.value(1,10000))
from dual
2 connect by level <=10000;

EODA@ORA12CR1> exec dbms_stats.gather_table_stats(user,'T');
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.

EODA@ORA12CR1> create index ti on t(x) invisible;
Index created.
Индекс создан.
```

Теперь мы включим средство AUTOTRACE и запустим запрос, в котором от оптимизатора ожидается использование индекса при генерации им плана выполнения:

```

EODA@ORA12CR1> set autotrace traceonly explain
EODA@ORA12CR1> select * from t where x=5;
```

| Id | Operation | Name | Rows | Bytes | Cost (%CPU) | Time |
|-----|-------------------|------|------|-------|-------------|----------|
| 0 | SELECT STATEMENT | | 2 | 8 | 7 (0) | 00:00:01 |
| * 1 | TABLE ACCESS FULL | T | 2 | 8 | 7 (0) | 00:00:01 |

В предыдущем выводе видно, что индекс оптимизатором не задействован. Переключать видимость индекса оптимизатору во время сеанса можно, устанавливая параметр инициализации OPTIMIZER_USE_INVISIBLE_INDEXES в TRUE (по умолчанию он установлен в FALSE). Например, для текущего подключенного сеанса следующая команда инструктирует оптимизатор о необходимости принять во внимание невидимые индексы при генерации плана выполнения:

```

EODA@ORA12CR1> alter session set optimizer_use_invisible_indexes=true;
```

Повторный запуск предшествующего запроса показывает, что оптимизатор теперь применяет индекс:

```

EODA@ORA12CR1> select * from t where x=5;
```

| Id | Operation | Name | Rows | Bytes | Cost (%CPU) | Time |
|-----|------------------|------|------|-------|-------------|----------|
| 0 | SELECT STATEMENT | | 2 | 8 | 1 (0) | 00:00:01 |
| * 1 | INDEX RANGE SCAN | TI | 2 | 8 | 1 (0) | 00:00:01 |

Если вы хотите, чтобы невидимые индексы принимались во внимание всеми сеансами, тогда измените посредством оператора `ALTER SYSTEM` параметр `OPTIMIZER_USE_INVISIBLE_INDEXES`. Это сделает все невидимые индексы видимыми оптимизатору во время генерации планов выполнения.

Индекс можно сделать постоянно видимым оптимизатору, как показано ниже:

```
EODA@ORA12CR1> alter index ti visible;
Index altered.
Индекс изменен.
```

Имейте в виду, что хотя невидимый индекс может быть невидимым оптимизатору, он по-прежнему может оказывать влияние на производительность следующим образом.

- Невидимые индексы потребляют пространство и ресурсы, когда в лежащей в основе таблице производится вставка, обновление или удаление записей. Это может сказаться на производительности (замедлить выполнение операторов DML).
- База данных Oracle может все еще использовать невидимый индекс, чтобы предотвратить определенные ситуации с блокировкой, когда индекс со структурой B-дерева помещается на столбец внешнего ключа.
- Если вы создали уникальный невидимый индекс, то уникальность столбцов будет обеспечиваться независимо от установки невидимости.

По перечисленным причинам, даже если вы создали индекс как невидимый, он может по-прежнему влиять на поведение SQL-операторов. Было бы ошибочно полагать, что невидимые индексы не оказывают никакого воздействия на приложения, работающие с таблицами, в которых такие индексы существуют. Невидимые индексы являются невидимыми только в том смысле, что оптимизатор не рассматривает их применение при генерации планов выполнения, если только он не проинструктирован о необходимости их использовать.

Итак, в чем состоит польза от невидимых индексов? Такие индексы все равно будут поддерживаться (отсюда и снижение производительности), но не могут использоваться запросами, которые не способны видеть их (поэтому производительность никогда не повышается). Одним примером может служить ситуация, когда вы хотите удалить индекс в производственной системе. Идея заключается в том, что вы можете сделать индекс невидимым и посмотреть, не пострадала ли в результате производительность. В этом случае перед удалением индекса важно также проверить, не размещен ли он на столбце внешнего ключа либо не задействован ли он для принудительного применения ограничения уникальности. Другой пример связан с необходимостью добавления индекса в производственной системе и прогона тестов, чтобы выяснить, улучшилась ли производительность. Вы можете добавить индекс как невидимый и выборочно сделать его видимым на протяжении сеанса, чтобы определить пользу от него. И снова вы должны помнить о том, что хотя индекс является невидимым, он будет занимать пространство и требовать ресурсов для поддержания.

Множество индексов на одной и той же комбинации столбцов

До выхода Oracle 12c нельзя было иметь несколько индексов, определенных в одной таблице с точно такой же комбинацией столбцов. Например:

```

EODA@ORA11GR2> create table t(x int);
Table created.
Таблица создана.

EODA@ORA11GR2> create index ti on t(x);
Index created.
Индекс создан.

EODA@ORA11GR2> create bitmap index tb on t(x) invisible;
ERROR at line 1:
ORA-01408: such column list already indexed
ОШИБКА в строке 1:
ORA-01408: такой список столбцов уже проиндексирован

```

Начиная с версии Oracle 12c, появилась возможность определять множество индексов на том же самом наборе столбцов. Тем не менее, это можно делать, только если индексы физически различны; например, когда один индекс создается как имеющий структуру В-дерева, а другой как битовый. К тому же для одной и той же комбинации столбцов в таблице может существовать только один видимый индекс. Следовательно, показанные выше операторы CREATE INDEX будут работать в базе данных Oracle 12c:

```

EODA@ORA12CR1> create table t(x int);
Table created.
Таблица создана.

EODA@ORA12CR1> create index ti on t(x);
Index created.
Индекс создан.

EODA@ORA12CR1> create bitmap index tb on t(x) invisible;
Index created.
Индекс создан.

```

Зачем могут понадобиться два индекса, определенные на том же самом наборе столбцов? Предположим, что вы изначально построили хранилище данных, имеющее звездообразную схему, в которой на всех столбцах внешних ключей определены индексы со структурой В-дерева, а позже посредством тестирования обнаружили, что битовые индексы будут выполняться эффективнее для типов запросов, применяемых к звездообразной схеме. Вследствие этого вы хотите преобразовать существующие индексы в битовые как можно более плавно. Сначала вы построите битовые индексы как невидимые. Затем, когда все будет готово, вы можете удалить индексы со структурой В-дерева и сделать битовые индексы видимыми.

Индексация расширенных столбцов

С появлением Oracle 12c типы данных VARCHAR2, NVARCHAR2 и RAW могут быть сконфигурированы для хранения вплоть до 32 767 байтов информации (ранее пре-

дельный объем составлял 4 000 байтов для VARCHAR2 и NVARCHAR2 и 2000 байтов для RAW). Поскольку в главе 12 приводятся подробные сведения по включению расширенных типов в базу данных, здесь они повторяться не будут. Основное внимание в этом разделе посвящено исследованию индексации расширенных столбцов.

Давайте начнем с создания таблицы, имеющей расширенный столбец, и затем попробуем создать на этом столбце обычный индекс со структурой B-дерева:

```

EODA@O12CE> create table t(x varchar2(32767));
Table created.
Таблица создана.

```

На заметку! Если вы попытаетесь создать таблицу со столбцом VARCHAR2, превышающим размер 4 000 байтов, в базе данных, которая не была сконфигурирована для расширенных типов данных, то Oracle выдаст сообщение об ошибке ORA-00910: specified length too long for its datatype (ORA-00910: указанная длина слишком большая для типа данных).

Теперь попробуем создать индекс на расширенном столбце:

```

EODA@O12CE> create index ti on t(x);
create index ti on t(x)
*
ERROR at line 1:
ORA-01450: maximum key length (6398) exceeded
ОШИБКА в строке 1:
ORA-01450: превышена максимальная длина ключа (6398)

```

Мы уже видели такую ошибку ранее в этой главе. Ошибка возникла из-за того, что Oracle налагает ограничение на максимальную длину ключа индекса, которая составляет около трех четвертей размера блока (размер блока для базы данных в рассматриваемом примере равен 8 Кбайт). Несмотря на то что записи в этом индексе пока отсутствуют, базе данных Oracle известно о том, что данный ключ индекса может оказаться больше, чем 6 398 байтов, для строки, способной уместить до 32 767 байтов, следовательно, она не позволит создать индекс при таком сценарии.

Этот вовсе не означает, что индексировать расширенные столбцы невозможно, просто вы должны использовать определенные приемы для ограничения длины ключа индекса значением, не превышающим 6 398 байтов. С учетом этого становятся очевидными следующие варианты.

- Создать виртуальный столбец на основе функции SUBSTR или STANDARD_HASH и затем создать индекс на полученном виртуальном столбце.
- Создать индекс на основе функции с применением SUBSTR или STANDARD_HASH.
- Создать табличное пространство с более высоким размером блока; например, блок размером 16 Кбайт позволил бы располагать ключами индекса с длиной приблизительно 12 000 байтов. Надо сказать, что если вы нуждаетесь в 12 000 байтов для ключа индекса, то вполне вероятно делаете что-то не так и должны переосмыслить свои действия. Этот метод исследоваться не будет.

Первым мы рассмотрим решение с виртуальным столбцом.

Решение с виртуальным столбцом

Идея заключается в том, чтобы сначала создать виртуальный столбец, который использует функцию SQL на расширенном столбце для возврата значения с длиной меньше 6 398 байтов. Затем этот виртуальный столбец может быть проиндексирован, что предоставит механизм для обеспечения более высокой производительности при выдаче запросов в отношении расширенных столбцов. Продемонстрируем такой подход на примере. Прежде всего, создадим таблицу с расширенным столбцом:

```
EODA@012CE> create table t(x varchar2(32767));
Table created.
Таблица создана.
```

Затем вставим в эту таблицу тестовые данные:

```
EODA@012CE> insert into t select to_char(level)|| rpad('abc',10000,'xyz')
  2  from dual connect by level < 1001
  3  union
  4  select to_char(level)
  5  from dual connect by level < 1001;
2000 rows created.
2000 строк создано.
```

Предположим, что первые десять символов расширенного столбца являются достаточно избирательными для возвращения небольших порций строк из таблицы. По этой причине создадим виртуальный столбец на основе части строки расширенного столбца:

```
EODA@012CE> alter table t add (xv as (substr(x,1,10)));
Table altered.
Таблица изменена.
```

Далее создадим индекс на виртуальном столбце и соберем статистику:

```
EODA@012CE> create index te on t(xv);
Index created.
Индекс создан.

EODA@012CE> exec dbms_stats.gather_table_stats(user,'T');
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.
```

Теперь при запросе виртуального столбца оптимизатор может воспользоваться индексом в предикатах равенства и диапазона, указываемых внутри конструкции WHERE, например:

```
EODA@012CE> set autotrace traceonly explain
EODA@012CE> select count(*) from t where x = '800';
```

| Id | Operation | Name | Rows | Bytes | Cost (%CPU) | Time |
|-----|-------------------------------------|------|------|-------|-------------|----------|
| 0 | SELECT STATEMENT | | 1 | 5011 | 2 (0) | 00:00:01 |
| 1 | SORT AGGREGATE | | 1 | 5011 | | |
| * 2 | TABLE ACCESS BY INDEX ROWID BATCHED | T | 1 | 5011 | 2 (0) | 00:00:01 |
| * 3 | INDEX RANGE SCAN | TE | 1 | | 1 (0) | 00:00:01 |

Обратите внимание, что хотя индекс определен на виртуальном столбце, оптимизатор по-прежнему может применять его при запрашивании напрямую расширенного столбца X (но не виртуального столбца XV). Оптимизатор также может использовать такой тип индекса в поиске по диапазону:

```
EODA@012CE> select count(*) from t where x >'800' and x<'900';
```

| Id | Operation | Name | Rows | Bytes | Cost (%CPU) | Time |
|-----|-------------------------------------|------|------|-------|-------------|----------|
| 0 | SELECT STATEMENT | | 1 | 5011 | 4 (0) | 00:00:01 |
| 1 | SORT AGGREGATE | | 1 | 5011 | | |
| * 2 | TABLE ACCESS BY INDEX ROWID BATCHED | T | 239 | 1169K | 4 (0) | 00:00:01 |
| * 3 | INDEX RANGE SCAN | TE | 241 | | 2 (0) | 00:00:01 |

Аналогично функции SUBSTR, виртуальный столбец можно основывать также и на функции STANDARD_HASH. Функция STANDARD_HASH может применяться к длинной символьной строке и возвращает в определенной степени уникальное значение RAW, которое намного меньше 6 398 байтов. Давайте взглянем на пару примеров использования виртуального столбца, основанного на функции STANDARD_HASH.

Предполагая, что применяется та же самая таблица с начальными данными, как и в предшествующих примерах с функцией SUBSTR, добавим к ней виртуальный столбец, использующий STANDARD_HASH, создадим индекс и сгенерируем статистику:

```
EODA@012CE> alter table t add (xv as (standard_hash(x)));
```

Table altered.

Таблица изменена.

```
EODA@012CE> create index te on t(xv);
```

Index created.

Индекс создан.

```
EODA@012CE> exec dbms_stats.gather_table_stats(user, 'T');
```

PL/SQL procedure successfully completed.

Процедура PL/SQL успешно завершена.

Функция STANDARD_HASH работает хорошо, когда применяются предикаты равенства в конструкции WHERE. Например:

```
EODA@012CE> set autotrace traceonly explain
```

```
EODA@012CE> select count(*) from t where x='300';
```

| Id | Operation | Name | Rows | Bytes | Cost (%CPU) | Time |
|-----|-------------------------------------|------|------|-------|-------------|----------|
| 0 | SELECT STATEMENT | | 1 | 5025 | 2 (0) | 00:00:01 |
| 1 | SORT AGGREGATE | | 1 | 5025 | | |
| * 2 | TABLE ACCESS BY INDEX ROWID BATCHED | T | 1 | 5025 | 2 (0) | 00:00:01 |
| * 3 | INDEX RANGE SCAN | TE | 1 | | 1 (0) | 00:00:01 |

Индекс на виртуальном столбце, основанном на функции STANDARD_HASH, позволяет эффективно выполнять поиск по равенству, но не работает для поиска по диапазону, т.к. данные, которые хранятся в индексе, базируются на рандомизированном хеш-значении; например:

```
EODA@012CE> select count(*) from t where x > '800' and x < '900';
```

| Id | Operation | Name | Rows | Bytes | Cost (%CPU) | Time |
|-----|-------------------|------|------|-------|-------------|----------|
| 0 | SELECT STATEMENT | | 1 | 5004 | 6 (0) | 00:00:01 |
| 1 | SORT AGGREGATE | | 1 | 5004 | | |
| * 2 | TABLE ACCESS FULL | T | 239 | 1167K | 6 (0) | 00:00:01 |

Решение с индексом на основе функции

Концепция состоит в том, что вы строите индекс и применяете к нему функцию способом, который ограничивает длину ключа индекса и также дает в результате индекс, пригодный к эксплуатации. Здесь для создания таблицы с расширенным столбцом и ее заполнения тестовыми данными используется тот же самый код (как в предыдущем разделе):

```
EODA@012CE> create table t(x varchar2(32767));
```

Table created.

Таблица создана.

```
EODA@012CE> insert into t
```

```
2 select to_char(level)|| rpad('abc',10000,'xyz')
```

```
3 from dual connect by level < 1001
```

```
4 union
```

```
5 select to_char(level)
```

```
6 from dual connect by level < 1001;
```

2000 rows created.

2000 строк создано.

Предположим, что вы ознакомлены с данными и знаете, что первых десяти символов из расширенного столбца обычно достаточно для идентификации строки; вследствие этого вы создаете индекс на подстроке, состоящей из этих десяти символов, и генерируете статистику для таблицы:

```
EODA@012CE> create index te on t(substr(x,1,10));
```

Index created.

Индекс создан.

```
EODA@012CE> exec dbms_stats.gather_table_stats(user,'T');
```

PL/SQL procedure successfully completed.

Процедура PL/SQL успешно завершена.

Оптимизатор может задействовать индекс подобного рода при наличии предикатов равенства и диапазона в конструкции WHERE. Вот пример случая с предикатом равенства:

```
EODA@012CE> set autotrace traceonly explain
```

```
EODA@012CE> select count(*) from t where x = '800';
```

| Id | Operation | Name | Rows | Bytes | Cost (%CPU) | Time |
|-----|-------------------------------------|------|------|-------|-------------|----------|
| 0 | SELECT STATEMENT | | 1 | 16407 | 2 (0) | 00:00:01 |
| 1 | SORT AGGREGATE | | 1 | 16407 | | |
| * 2 | TABLE ACCESS BY INDEX ROWID BATCHED | T | 1 | 16407 | 2 (0) | 00:00:01 |
| * 3 | INDEX RANGE SCAN | TE | 8 | | 1 (0) | 00:00:01 |

В следующем примере применяется предикат диапазона:

```
EODA@012CE> select count(*) from t where x>'200' and x<'400';
```

| Id | Operation | Name | Rows | Bytes | Cost (%CPU) | Time |
|-----|-------------------------------------|------|------|-------|-------------|----------|
| 0 | SELECT STATEMENT | | 1 | 5011 | 6 (0) | 00:00:01 |
| 1 | SORT AGGREGATE | | 1 | 5011 | | |
| * 2 | TABLE ACCESS BY INDEX ROWID BATCHED | T | 477 | 2334K | 6 (0) | 00:00:01 |
| * 3 | INDEX RANGE SCAN | TE | 479 | | 3 (0) | 00:00:01 |

Имея таблицу и начальные данные из предыдущих примеров использования функции SUBSTR, добавим к ней индекс на основе функции STANDARD_HASH:

```
EODA@012CE> create index te on t(standard_hash(x));
```

Теперь удостоверимся в том, что поиск по равенству применяет этот индекс:

```
EODA@012CE> set autotrace traceonly explain
```

```
EODA@012CE> select count(*) from t where x = '800';
```

| Id | Operation | Name | Rows | Bytes | Cost (%CPU) | Time |
|-----|-------------------------------------|------|------|-------|-------------|----------|
| 0 | SELECT STATEMENT | | 1 | 5004 | 4 (0) | 00:00:01 |
| 1 | SORT AGGREGATE | | 1 | 5004 | | |
| * 2 | TABLE ACCESS BY INDEX ROWID BATCHED | T | 1 | 5004 | 4 (0) | 00:00:01 |
| * 3 | INDEX RANGE SCAN | TE | 8 | | 1 (0) | 00:00:01 |

Это делает возможным эффективный поиск по равенству, но не подходит для поиска по диапазону, поскольку данные хранятся в индексе на основе рандомизированного хеш-значения.

Часто задаваемые вопросы и мифы об индексах

Как упоминалось во введении, ко мне поступает огромное количество вопросов относительно Oracle. Я — тот самый Том, ведущий колонку “Ask Tom” (“Спросите у Тома”) в журнале *Oracle Magazine* и на веб-сайте <http://asktom.oracle.com>, где я отвечаю на вопросы людей о базе данных и инструментах Oracle. Согласно моему опыту, тема индексов вызывает наибольшее число вопросов. В этом разделе я отвечу на некоторые из наиболее часто задаваемых вопросов. Одни ответы могут показаться совершенно очевидными, в то время как другие могут удивить. Достаточно сказать, что индексы окружает множество мифов и недоразумений.

Работают ли индексы в представлениях?

Связанный с этим вопрос звучит так: каким образом можно проиндексировать представление? Прежде всего, факт заключается в том, что представление — это всего лишь хранимый запрос. База данных Oracle заменяет текст запроса, который обращается к представлению, определением самого представления. Представления задумывались для удобства конечного пользователя или программиста — оптимизатор имеет дело с запросом к их базовым таблицам. Любые индексы, которые могли использоваться в запросе, написанном в отношении к базовым таблицам, будут учи-

тываться во время применения представления. Чтобы проиндексировать представление, нужно просто проиндексировать его базовые таблицы.

Могут ли значения NULL и индексы работать вместе?

Индексы со структурой В-дерева, за исключением особого случая кластерных индексов со структурой В-дерева, не хранят записи, содержащие NULL, но битовые и кластерные индексы сохраняют их. Этот побочный эффект может оказаться источником путаницы, но в действительности он позволяет извлечь выгоду, когда вы понимаете то, что именно подразумевает отказ от хранения ключей NULL.

Чтобы увидеть влияние от того факта, что значения NULL *не* сохраняются, рассмотрим следующий пример:

```
EODA@ORA12CR1> create table t ( x int, y int );
Table created.
Таблица создана.

EODA@ORA12CR1> create unique index t_idx on t(x,y);
Index created.
Индекс создан.

EODA@ORA12CR1> insert into t values ( 1, 1 );
1 row created.
1 строка создана.

EODA@ORA12CR1> insert into t values ( 1, NULL );
1 row created.
1 строка создана.

EODA@ORA12CR1> insert into t values ( NULL, 1 );
1 row created.
1 строка создана.

EODA@ORA12CR1> insert into t values ( NULL, NULL );
1 row created.
1 строка создана.

EODA@ORA12CR1> analyze index t_idx validate structure;
Index analyzed.
Индекс проанализирован.

EODA@ORA12CR1> select name, lf_rows from index_stats;
```

| NAME | LF_ROWS |
|-------|---------|
| T_IDX | 3 |

Таблица имеет четыре строки, тогда как индекс — только три. Первые три строки, в которых хотя бы один элемент ключа индекса не равен NULL, находятся в индексе. Последняя строка со значениями (NULL, NULL) в индекс не попадает. Одной из ситуаций, в которых возникает путаница, является случай с уникальным индексом, как в приведенном выше примере.

Рассмотрим результат действия следующих трех операторов INSERT:

```
EODA@ORA12CR1> insert into t values ( NULL, NULL );
1 row created.
1 строка создана.
```

```
EODA@ORA12CR1> insert into t values ( NULL, 1 );
insert into t values ( NULL, 1 )
*
ERROR at line 1:
ORA-00001: unique constraint (EODA.T_IDX) violated
ОШИБКА в строке 1:
ORA-00001: нарушено ограничение уникальности (EODA.T_IDX)
```

```
EODA@ORA12CR1> insert into t values ( 1, NULL );
insert into t values ( 1, NULL )
*
ERROR at line 1:
ORA-00001: unique constraint (EODA.T_IDX) violated
ОШИБКА в строке 1:
ORA-00001: нарушено ограничение уникальности (EODA.T_IDX)
```

Новая строка (NULL, NULL) не считается такой же, как старая строка с (NULL, NULL):

```
EODA@ORA12CR1> select x, y, count(*)
2 from t
3 group by x,y
4 having count(*) > 1;
```

| X | Y | COUNT(*) |
|-------|-------|----------|
| ----- | ----- | ----- |
| | | 2 |

Это выглядит неправдоподобным; если учитывать все записи NULL, то наш уникальный индекс перестает быть уникальным. Дело в том, что когда речь заходит об уникальности, то значения (NULL, NULL) в Oracle не считаются идентичными — так требует стандарт SQL. Однако, что касается агрегирования, то значения (NULL, NULL) и (NULL, NULL) трактуются как одинаковые. Указанные две комбинации в выражениях сравнения рассматриваются как разные, но внутри конструкции GROUP BY — как совпадающие. Это то, что следует принимать во внимание: каждое ограничение уникальности должно иметь, по меньшей мере, один столбец NOT NULL, чтобы быть по-настоящему уникальным.

Вопрос, который возникает в связи с индексами и значениями NULL, звучит следующим образом: почему мой запрос не использует индекс? Запрос, о котором идет речь, выглядит примерно так:

```
select * from T where x is null;
```

Такой запрос не может задействовать только что созданный индекс — строка (NULL, NULL) просто отсутствует в индексе, поэтому применение индекса привело бы к возвращению неправильного ответа. Запрос может использовать индекс, только если хотя бы один столбец определен как NOT NULL. Например, следующие команды показывают, что Oracle задействует индекс для предиката X IS NULL, если

этот индекс содержит X в начале своего списка столбцов и, по крайней мере, один из остальных столбцов определен в базовой таблице как NOT NULL:

```
EODA@ORA12CR1> create table t ( x int, y int NOT NULL );
Table created.
Таблица создана.
```

```
EODA@ORA12CR1> create unique index t_idx on t(x,y);
Index created.
Индекс создан.
```

```
EODA@ORA12CR1> insert into t values ( 1, 1 );
1 row created.
1 строка создана.
```

```
EODA@ORA12CR1> insert into t values ( NULL, 1 );
1 row created.
1 строка создана.
```

```
EODA@ORA12CR1> begin
2   dbms_stats.gather_table_stats(user,'T');
3 end;
4 /
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.
```

Запустив запрос к этой таблице, мы обнаружим вот что:

```
EODA@ORA12CR1> set autotrace on
EODA@ORA12CR1> select * from t where x is null;
```

| X | Y |
|-------|-------|
| ----- | ----- |
| | 1 |

Execution Plan
...

| Id | Operation | Name | Rows | Bytes | Cost (%CPU) | Time |
|-----|------------------|-------|------|-------|-------------|----------|
| 0 | SELECT STATEMENT | | 1 | 5 | 1 (0) | 00:00:01 |
| * 1 | INDEX RANGE SCAN | T_IDX | 1 | 5 | 1 (0) | 00:00:01 |

Ранее говорилось, что отказом от хранения полных записей NULL в индексах со структурой В-дерева можно воспользоваться для собственной выгоды, и сейчас будет показано, как это сделать. Скажем, есть таблица со столбцом, который принимает в точности два значения. Значения распределены весьма ассиметрично: 90% или больше строк имеют одно значение, а 10% или меньше — другое. Такой столбец можно эффективно проиндексировать, чтобы добиться быстрого доступа к меньшинству строк. Это может пригодиться, когда для получения меньшинства строк желательно применять индекс, но для оставшегося большинства строк — полное сканирование, и нужно сберечь дисковое пространство.

Решение заключается в использовании значения NULL для большинства строк и любого выбранного значения для меньшинства или, как демонстрировалось ранее, в применении индекса на основе функции для индексации только тех возвращаемых значений функции, которые не равны NULL.

Теперь, когда вы знаете, каким образом В-дерево трактует значения NULL, это можно задействовать в собственных интересах и наложить ограничения уникальности на наборы столбцов, которые все допускают NULL (подготовившись к тому, что содержать во всех ключевых столбцах NULL может более одной строки).

Должны ли быть проиндексированы внешние ключи?

Вопрос о том, должны ли индексироваться внешние ключи, возникает довольно часто. Мы касались этой темы в главе 6, когда обсуждали взаимоблокировки. Там было указано, что неиндексированные внешние ключи являются самой главной причиной взаимоблокировок среди тех, что я встречал. Это связано с тем фактом, что обновление первичного ключа родительской таблицы либо удаление записи из родительской таблицы поместит блокировку на всю дочернюю таблицу (никакие модификации в дочерней таблице не будут разрешены, пока оператор не завершится). В результате блокируется намного больше строк, чем должно, и снижается степень параллелизма. Я часто наблюдаю такую ситуацию, когда пользователи работают с инструментом, генерирующим SQL-код для модификации таблицы. Инструмент формирует оператор UPDATE, который обновляет каждый столбец в таблице независимо от того, было ли модифицировано его значение. В действительности он обновит первичный ключ, хотя его значение никогда не изменяется. Например, средство Oracle Forms будет делать это по умолчанию, если только вы не сообщите ему о необходимости отправлять в базу данных только модифицированные столбцы. В дополнение к проблеме блокировки таблицы, с которой вы можете столкнуться, неиндексированный внешний ключ плох также в перечисленных ниже случаях.

- Когда присутствует конструкция ON DELETE CASCADE и дочерняя таблица не индексирована. Например, таблица EMP является дочерней по отношению к DEPT. Команда DELETE FROM DEPT WHERE DEPTNO = 10 должна каскадным образом распространиться на EMP. Если столбец DEPTNO в EMP не индексирован, то будет производиться полное сканирование EMP. Это полное сканирование, скорее всего, нежелательно, и в случае удаления множества строк в родительской таблице дочерняя таблица будет сканироваться по одному разу на каждую строку, удаляемую из родительской таблицы.
- Когда выполняется запрос от родительской таблицы к дочерней. Рассмотрим снова пример с таблицами EMP/DEPT. Совсем нередко приходится запрашивать EMP в контексте DEPTNO. Если вы часто выполняете приведенный ниже запрос для генерации отчета или чего-то еще, то обнаружите, что отсутствие индекса по внешнему ключу будет замедлять такие запросы:

```
select *
  from dept, emp
 where emp.deptno = dept.deptno
 and dept.dname = :X;
```

Тот же самый аргумент я приводил в пользу индексации по столбцу NESTED_COLUMN_ID вложенной таблицы в главе 10. Скрытый столбец NESTED_COLUMN_ID вложенной таблицы — это не что иное, как внешний ключ.

Итак, когда же *не* нужно индексировать внешний ключ? В общем случае такая индексация не понадобится при удовлетворении следующих условий.

- Вы не удаляете строки из родительской таблицы.
- Вы не обновляете значение уникального/первичного ключа родительской таблицы — ни преднамеренно, ни случайно (через инструмент).
- Вы не выполняете соединение родительской таблицы с дочерней (такое как DEPT к EMP). В более общих чертах можно сказать, что столбцы внешнего ключа не поддерживают важный путь доступа к дочерней таблице, и вы не используете их в предикатах для выборки данных из этой таблицы.

Если все три критерия удовлетворены, можете отказаться от индекса — в нем нет необходимости, и он замедлит выполнение операций DML в дочерней таблице. Если какое-то из этих трех условий не удовлетворено, отдавайте себе отчет о последствиях.

В качестве стороннего замечания: если вы полагаете, что дочерняя таблица блокируется из-за неиндексированного внешнего ключа, и вы хотите удостовериться в этом (или просто предотвратить блокировку в целом), то можете выдать такую команду:

```
ALTER TABLE <имя дочерней таблицы> DISABLE TABLE LOCK;
```

Теперь любая операция UPDATE или DELETE в родительской таблице, которая вызывает табличную блокировку, завершится ошибкой:

```
ERROR at line 1:
ORA-00069: cannot acquire lock -- table locks disabled for <child table name>
ОШИБКА в строке 1:
ORA-00069: не удастся получить блокировку -- табличные блокировки
запрещены для <имя дочерней таблицы>
```

Такой прием удобен при отслеживании порции кода, которая делает то, что, по вашему мнению, делать не должна (команды UPDATE или DELETE в отношении родительского первичного ключа), т.к. конечные пользователи незамедлительно сообщат вам о возникновении этой ошибки.

Почему индекс не используется?

Возможных причин много. Ниже рассматриваются наиболее распространенные из них.

Случай 1

Мы имеем индекс со структурой В-дерева, и наш предикат не использует головную часть индекса. В этом случае мы можем иметь таблицу T с индексом T(x, y). Мы выдаем запрос SELECT * FROM T WHERE Y=5. Оптимизатор решает не применять индекс, поскольку предикат не задействует столбец X — в данном случае ему придется инспектировать каждую запись индекса (вскоре мы обсудим сканирование индекса с пропусками, где это не так). Обычно оптимизатор будет отдавать предпочтение полному сканированию таблицы T. Это вовсе не препятствует использованию индекса. Если запрос выглядит как SELECT X, Y FROM T WHERE Y=5, то оптимизатор заметит, что он не должен обращаться к таблице для получения ни X, ни Y (они

есть в индексе). В таком случае оптимизатор вполне может выбрать быстрое полное сканирование самого индекса, т.к. он обычно намного меньше лежащей в основе таблицы. Также обратите внимание, что такой путь доступа имеется в наличии только у оптимизатора СВО.

Другой случай, когда индекс на $T(x, y)$ может применяться оптимизатором СВО, касается выполнения сканирования индекса с пропусками. Сканирование с пропусками хорошо работает тогда и только тогда, когда головная часть индекса (x в предыдущем примере) имеет очень немного отличающихся друг от друга значений, и оптимизатору это известно. Например, рассмотрим индекс на $(GENDER, EMPNO)$, где $GENDER$ принимает значения M или F , а столбец $EMPNO$ уникален. Запрос вроде

```
select * from t where empno = 5;
```

может рассмотреть использование индекса на таблице T для удовлетворения запроса методом *сканирования с пропусками*, а это значит, что запрос будет концептуально обработан примерно так:

```
select * from t where GENDER='M' and empno = 5
UNION ALL
select * from t where GENDER='F' and empno = 5;
```

Он бегло просмотрит индекс, как будто бы индексов есть два: один для значений M и еще один для значений F . Это легко заметить в плане выполнения запроса. Создадим таблицу со столбцом, имеющим два значения, и проиндексируем его:

```
EODA@ORA12CR1> create table t
2 as
3 select decode(mod(rownum,2), 0, 'M', 'F' ) gender, all_objects.*
4 from all_objects
5 /
```

Table created.

Таблица создана.

```
EODA@ORA12CR1> create index t_idx on t(gender,object_id);
Index created.
```

Индекс создан.

```
EODA@ORA12CR1> exec dbms_stats.gather_table_stats( user, 'T' );
PL/SQL procedure successfully completed.
```

Процедура *PL/SQL* успешно завершена.

Если теперь выполнить запрос, то мы должны увидеть следующее:

```
EODA@ORA12CR1> set autotrace traceonly explain
EODA@ORA12CR1> select * from t t1 where object_id = 42;
```

Execution Plan

...

| Id | Operation | Name | Rows | Bytes | Cost (%CPU) | Time |
|-----|-------------------------------------|-------|------|-------|-------------|----------|
| 0 | SELECT STATEMENT | | 1 | 91 | 4 (0) | 00:00:01 |
| 1 | TABLE ACCESS BY INDEX ROWID BATCHED | T | 1 | 91 | 4 (0) | 00:00:01 |
| * 2 | INDEX SKIP SCAN | T_IDX | 1 | | 3 (0) | 00:00:01 |

Наличие шага INDEX SKIP SCAN говорит о том, что Oracle бегло просматривает индекс в поисках мест, где столбец GENDER меняет значение, и читает вниз по дереву с этого места в поисках OBJECT_ID=42 в каждом рассматриваемом виртуальном индексе. Если заметно увеличить количество отличающихся значений GENDER, мы увидим, что Oracle перестает считать сканирование с пропусками приемлемым планом выполнения запроса:

```
EODA@ORA12CR1> update t set gender = chr(mod(rownum,256));
17944 rows updated.
17944 строк обновлено.

EODA@ORA12CR1> exec dbms_stats.gather_table_stats( user, 'T', cascade=>TRUE );
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.
```

В таком случае пришлось бы проверить 256 мини-индексов, поэтому для нахождения нужной строки выбирается полное сканирование таблицы:

```
EODA@ORA12CR1> set autotrace traceonly explain
EODA@ORA12CR1> select * from t t1 where object_id = 42;
```

Execution Plan

...

| Id | Operation | Name | Rows | Bytes | Cost (%CPU) | Time |
|-----|-------------------|------|------|-------|-------------|----------|
| 0 | SELECT STATEMENT | | 1 | 92 | 274 (1) | 00:00:01 |
| * 1 | TABLE ACCESS FULL | T | 1 | 92 | 274 (1) | 00:00:01 |

Случай 2

Мы применяем запрос SELECT COUNT (*) FROM T (или подобный ему) и располагаем в таблице T индексом со структурой В-дерева. Тем не менее, оптимизатор выбирает полное сканирование таблицы вместо подсчета записей в индексе (намного меньшем). В этом случае индекс, вероятно, построен на наборе столбцов, которые могут принимать значения NULL. Поскольку записи индекса, полностью содержащие NULL, не сохраняются, количество строк в индексе может не совпадать с количеством строк в таблице. Поэтому здесь оптимизатор поступает правильно — если он воспользуется для подсчета строк индексом, то получит ошибочный ответ.

Случай 3

Для проиндексированного столбца мы запускаем следующий запрос и обнаруживаем, что индекс по INDEXED_COLUMN не был задействован:

```
select * from t where f(indexed_column) = value
```

Это объясняется применением функции к этому столбцу. Мы проиндексировали значения столбца INDEXED_COLUMN, а не значения F(INDEXED_COLUMN). Возможность использования индекса здесь сокращена. При необходимости можно провести индексацию по функции.

Случай 4

Мы проиндексировали символьный столбец. Этот столбец содержит только числовые данные. Мы выполняем запрос с применением такого синтаксиса:

```
select * from t where indexed_column = 5
```

Обратите внимание, что 5 в запросе — это константное *число* 5 (не символьная строка). Индекс на столбце INDEXED_COLUMN не используется. Причина в том, что предыдущий запрос соответствует следующему запросу:

```
select * from t where to_number(indexed_column) = 5
```

Мы имеем неявно примененную к столбцу функцию и, как отмечалось в разделе “Случай 3”, это препятствует использованию индекса. Сказанное очень легко проверить на простом примере, в котором мы задействуем встроенный пакет DBMS_XPLAN. Данный пакет доступен только в Oracle9i Release 2 и последующих версиях (в Oracle9i Release 1 вместо него нужно применять средство AUTOTRACE, что позволит увидеть план, но информация о предикатах отображаться не будет — это возможно лишь в Oracle9i Release 2 и последующих версиях):

```
EODA@ORA12CR1> create table t ( x char(1) constraint t_pk primary key,
  2 y date );
Table created.
Таблица создана.
```

```
EODA@ORA12CR1> insert into t values ( '5', sysdate );
1 row created.
1 строка создана.
```

```
EODA@ORA12CR1> explain plan for select * from t where x = 5;
Explained.
Объяснено.
```

```
EODA@ORA12CR1> select * from table(dbms_xplan.display);
```

```
PLAN_TABLE_OUTPUT
```

```
-----
Plan hash value: 1601196873
```

```
-----
| Id | Operation                | Name | Rows | Bytes | Cost (%CPU) | Time      |
-----+-----+-----+-----+-----+-----+-----+
|  0 | SELECT STATEMENT         |      |     1 |    12 |        3 (0) | 00:00:01 |
|* 1 | TABLE ACCESS FULL      | T    |     1 |    12 |        3 (0) | 00:00:01 |
-----
```

```
Predicate Information (identified by operation id):
```

```
-----
1 - filter(TO_NUMBER("X")=5)
```

Как видите, запрос полностью сканирует таблицу. Даже если мы снабдим запрос подсказкой, он будет использовать индекс, но не для операции INDEX UNIQUE SCAN, как можно было ожидать, а для INDEX FULL SCAN:


```
EODA@ORA12CR1> explain plan for select /*+ INDEX(t t_pk) */ * from t
where x = 5;
Explained.
Объяснено.
```

```
EODA@ORA12CR1> select * from table(dbms_xplan.display);
```

```
PLAN_TABLE_OUTPUT
```

```
-----
Plan hash value: 180604526
```

| Id | Operation | Name | Rows | Bytes | Cost (%CPU) | Time |
|-----|-------------------------------------|------|------|-------|-------------|----------|
| 0 | SELECT STATEMENT | | 1 | 12 | 3 (0) | 00:00:01 |
| 1 | TABLE ACCESS BY INDEX ROWID BATCHED | T | 1 | 12 | 3 (0) | 00:00:01 |
| * 2 | INDEX FULL SCAN | T_PK | 1 | | 2 (0) | 00:00:01 |

```
-----
Predicate Information (identified by operation id):
```

```
-----
2 - filter(TO_NUMBER("X")=5)
```

Причина заключена в последней строке вывода: `filter(TO_NUMBER("X")=5)`. Это неявная функция, применяемая к столбцу таблицы. Символьная строка, хранящаяся в X, должна быть преобразована в число перед выполнением сравнения со значением 5. Мы не можем преобразовать 5 в строку, т.к. представлением числа 5 в форме строки управляют настройки NLS (это недетерминировано), поэтому мы преобразуем строку в число, и это предотвращает использование индекса для быстрого нахождения строки. Если мы просто сравним строки со строками:

```
EODA@ORA12CR1> explain plan for select * from t where x = '5';
Explained.
Объяснено.
```

```
EODA@ORA12CR1> select * from table(dbms_xplan.display);
```

```
PLAN_TABLE_OUTPUT
```

```
-----
Plan hash value: 1303508680
```

| Id | Operation | Name | Rows | Bytes | Cost (%CPU) | Time |
|-----|-----------------------------|------|------|-------|-------------|----------|
| 0 | SELECT STATEMENT | | 1 | 12 | 2 (0) | 00:00:01 |
| 1 | TABLE ACCESS BY INDEX ROWID | T | 1 | 12 | 2 (0) | 00:00:01 |
| * 2 | INDEX UNIQUE SCAN | T_PK | 1 | | 1 (0) | 00:00:01 |

```
-----
Predicate Information (identified by operation id):
```

```
-----
2 - access("X"='5')
```

то получим ожидаемую операцию INDEX UNIQUE SCAN и увидим, что функция не применяется. В любом случае вы должны *всегда* избегать неявных преобразований. Сравнивайте яблоки с яблоками, а апельсины с апельсинами. Другой случай, когда подобное часто случается, касается значений даты. Попробуйте выдать следующий запрос:

```
-- Найти все записи за сегодня.
select * from t where trunc(date_col) = trunc(sysdate);
```

Вы обнаружите, что индекс на DATE_COL не используется. Мы можем либо проиндексировать TRUNC (DATE_COL), либо, что проще, запрашивать с применением операций сравнения по диапазону. Ниже демонстрируется использование операций “больше чем” и “меньше чем” с датой. Поняв, что условие

```
TRUNC (DATE_COL) = TRUNC (SYSDATE)
```

представляет собой то же самое, что и

```
select *
  from t
 where date_col >= trunc(sysdate)
       and date_col < trunc(sysdate+1)
```

мы можем переместить все функции в правую часть равенства, позволив применять индекс на столбце DATE_COL (и обеспечить тот же самый результат, что и WHERE TRUNC (DATE_COL) = TRUNC (SYSDATE)).

По возможности вы должны всегда убирать функции от столбцов базы данных, когда они находятся в предикате. Это не только позволит рассматривать использование большего количества индексов, но также сократит объем работ, которые придется делать базе данных. В предыдущем случае, когда мы применяем конструкцию

```
where date_col >= trunc(sysdate)
       and date_col < trunc(sysdate+1)
```

значения TRUNC вычисляются один раз для всего запроса, после чего индекс может быть использован для нахождения уточненных значений. Когда мы применяем TRUNC (DATE_COL) = TRUNC (SYSDATE), значение TRUNC (DATE_COL) должно вычисляться по одному разу для *каждой строки* таблицы (без индексов).

Случай 5

Использование индекса может в действительности замедлять запросы. Я часто это вижу — люди предполагают, что индекс, конечно же, будет всегда ускорять выполнение запроса. Итак, они создают небольшую таблицу, анализируют ее и обнаруживают, что оптимизатор решил не применять индекс. В этом случае оптимизатор поступает совершенно правильно. База данных Oracle (с оптимизатором СВО) использует индекс только тогда, когда в этом есть смысл. Рассмотрим следующий пример:

```
EOA@ORA12CR1> create table t(x int);
Table created.
Таблица создана.
```

```
EOA@ORA12CR1> insert into t select rownum from dual connect by level < 1000000;
999999 rows created.
999999 строк создано.
```

```
EODA@ORA12CR1> create index ti on t(x);
```

```
Index created.
```

Индекс создан.

```
EODA@ORA12CR1> exec dbms_stats.gather_table_stats(user, 'T');
```

```
PL/SQL procedure successfully completed.
```

Процедура PL/SQL успешно завершена.

Если мы запустим запрос, которому необходимо извлечь относительно малый процент строк таблицы:

```
EODA@ORA12CR1> set autotrace on explain
```

```
EODA@ORA12CR1> select count(*) from t where x < 50;
```

```
COUNT(*)
```

```
-----
```

```
49
```

```
Execution Plan
```

```
...
```

| Id | Operation | Name | Rows | Bytes | Cost (%CPU) | Time |
|-----|------------------|------|------|-------|-------------|----------|
| 0 | SELECT STATEMENT | | 1 | 5 | 3 (0) | 00:00:01 |
| 1 | SORT AGGREGATE | | 1 | 5 | | |
| * 2 | INDEX RANGE SCAN | TI | 49 | 245 | 3 (0) | 00:00:01 |

то он благополучно применит индекс. Однако мы обнаружим, что если ожидаемое количество строк, которые должны быть извлечены через индекс, превышает некоторое пороговое значение (варьирующееся в зависимости от разнообразных настроек оптимизатора, физической статистики, версии и т.д.), то мы начнем наблюдать полное сканирование таблицы:

```
EODA@ORA12CR1> select count(*) from t where x < 1000000;
```

```
COUNT(*)
```

```
-----
```

```
999999
```

```
Execution Plan
```

```
-----
```

```
...
```

| Id | Operation | Name | Rows | Bytes | Cost (%CPU) | Time |
|-----|-------------------|------|------|-------|-------------|----------|
| 0 | SELECT STATEMENT | | 1 | 5 | 620 (1) | 00:00:01 |
| 1 | SORT AGGREGATE | | 1 | 5 | | |
| * 2 | TABLE ACCESS FULL | T | 999K | 4882K | 620 (1) | 00:00:01 |

Этот пример показывает, что оптимизатор не *всегда* использует индекс, и фактически делает правильный выбор, пропуская его. Если во время настройки запросов выясняется, что индекс не задействован, но вы считаете, что это *надлежало бы сделать*, не заставляйте Oracle его применять. Прежде чем отклонять предложение оптимизатора СВО, проведите тестирование и удостоверьтесь, что запрос с индексом работает действительно быстрее (оценив время выполнения и количество операций ввода-вывода).

Случай 6

Мы не обновляли статистику для таблиц. Таблицы обычно были малы, но теперь, когда мы просматриваем их, они довольно-таки значительно выросли. В использовании индекса появился смысл, чего первоначально не было. Если мы сгенерируем статистику, она покажет, что индекс станет применяться. Без корректной статистики оптимизатор СВО *не может* принимать правильные решения.

Заключительные соображения по поводу случаев, связанных с индексами

По моему опыту описанные шесть случаев охватывают *основные* причины, из-за которых индексы не применяются. Обычно все сводится к заключению: “Они не могут использоваться, т.к. это приведет к некорректным результатам” или “Они не должны применяться, иначе серьезно пострадает производительность”.

Миф: пространство никогда не используется в индексе повторно

Это миф, который я хотел бы развеять раз и навсегда: пространство в индексах *используется повторно*. Миф начинается примерно так: вы имеете таблицу T, в которой есть столбец X. В какой-то момент вы помещаете в таблицу значение X=5. Позже вы его удаляете. Миф заключается в том, что пространство, которое занимало значение X=5, никогда не будет задействовано повторно, если только вы вновь не вставите в таблицу X=5. Миф утверждает, что после того как слот индекса начал применяться, он остается на том же месте навсегда, и может повторно использоваться только тем же самым значением. Следствием этого является другой миф о том, что свободное пространство никогда не возвращается в структуру индекса, а блок никогда не используется повторно. Опять-таки, это попросту не соответствует действительности.

Первую часть мифа опровергнуть легко. Все, что для этого понадобится — это создать таблицу вроде показанной ниже:

```
EODA@ORA12CR1> create table t ( x int, constraint t_pk primary key(x) );
Table created.
Таблица создана.
```

```
EODA@ORA12CR1> insert into t values (1);
1 row created.
1 строка создана.
```

```
EODA@ORA12CR1> insert into t values (2);
1 row created.
1 строка создана.
```

```
EODA@ORA12CR1> insert into t values (999999999);
1 row created.
1 строка создана.
```

```
EODA@ORA12CR1> analyze index t_pk validate structure;
Index analyzed.
Индекс проанализирован.
```

```
EODA@ORA12CR1> select lf_blks, br_blks, btree_space from index_stats;
```

| LF_BKLS | BR_BKLS | BTREE_SPACE |
|---------|---------|-------------|
| 1 | 0 | 7996 |

Итак, согласно мифу, если удалить из таблицы T строку, в которой X=2, то это пространство никогда не будет задействовано повторно, если только снова не вставить строку с X=2. В настоящий момент этот индекс занимает один листовой блок пространства. Если пространство записей ключа индекса никогда не используется повторно после их удаления, а вставка и удаление будут продолжены без повторения значений, то этот индекс должен расти бешеными темпами. Давайте посмотрим:

```
EODA@ORA12CR1> begin
2      for i in 2 .. 999999
3      loop
4          delete from t where x = i;
5          commit;
6          insert into t values (i+1);
7          commit;
8      end loop;
9  end;
10 /
```

PL/SQL procedure successfully completed.

Процедура PL/SQL успешно завершена.

```
EODA@ORA12CR1> analyze index t_pk validate structure;
```

Index analyzed.

Индекс проанализирован.

```
EODA@ORA12CR1> select lf_blks, br_blks, btree_space from index_stats;
```

| LF_BLKs | BR_BLKs | BTree_Space |
|---------|---------|-------------|
| 1 | 0 | 7996 |

Вывод показывает, что пространство индекса применялось повторно. Тем не менее, как и в большинстве мифов, доля истины здесь есть. Дело в том, что пространство, используемое начальным значением 2, останется в блоке индекса навсегда. Индекс не будет сам себя перераспределять. Это значит, что если я загрузю таблицу значениями от 1 до 500 000, а затем удалю каждую вторую строку (все четные значения), то в индексе на этом столбце окажется 250 000 дыр. Пространство будет применяться повторно, только если я заново вставлю данные, которые уместятся в блок, где имеется дыра. База данных Oracle не будет пытаться уплотнить или сжать индекс. Это может быть сделано явно посредством команды ALTER INDEX REBUILD или ALTER INDEX COALESCE. С другой стороны, если я загрузю таблицу значениями от 1 до 500 000 и затем удалю из нее все строки, значения в которых равны 250 000 или меньше, то обнаружу, что очищенные блоки индекса возвратились обратно в список свободных блоков для индекса. Это пространство может быть полностью использовано повторно.

Если помните, этот миф сопровождается еще одним мифом: пространство индекса никогда не *возвращается*. Миф утверждает, что после того как началось применение блока индекса, он останется на этом месте в структуре индекса навсегда и будет использован повторно, только если вставляются данные, которые в любом случае попали бы в это место внутри индекса. Мы можем показать, что это также ошибочно. Для начала нам понадобится построить таблицу, содержащую около 500 000 строк.

Для этого мы применим сценарий для таблицы `BIG_TABLE`, который был описан в разделе “Настройка среды” в самом начале книги. Имея эту таблицу с соответствующим индексом по первичному ключу, мы измерим количество листовых блоков, имеющиххся в индексе, и количество блоков, содержащихся в списке свободных блоков для индекса. В отличие от таблицы, блок индекса попадает в список свободных только тогда, когда он совершенно пуст. Таким образом, все блоки, которые мы видим в списке свободных блоков, абсолютно пусты и доступны для повторного использования:

```
EODA@ORA12CR1> select count(*) from big_table;
```

```

COUNT(*)
-----
500000
```

На заметку! В следующем коде PL/SQL индекс, по которому генерируется отчет, должен быть построен в табличном пространстве `MSSM`, а не `ASSM`. Если вы попытаетесь запустить этот код в отношении индекса в табличном пространстве `ASSM`, то получите сообщение `ORA-10618: operation not allowed on this segment (ORA-10618: операция не разрешена в этом сегменте)`.

```
EODA@ORA12CR1> declare
2     l_freelist_blocks number;
3 begin
4     dbms_space.free_blocks
5     ( segment_owner => user,
6       segment_name => 'BIG_TABLE_PK',
7       segment_type => 'INDEX',
8       freelist_group_id => 0,
9       free_blks => l_freelist_blocks );
10    dbms_output.put_line( 'blocks on freelist = ' || l_freelist_blocks );
11 end;
12 /
```

blocks on freelist = 0

блоков в списке свободных блоков = 0

PL/SQL procedure successfully completed.

Процедура PL/SQL успешно завершена.

```
EODA@ORA12CR1> select leaf_blocks from user_indexes where index_name =
'BIG_TABLE_PK';
```

```

LEAF_BLOCKS
-----
1043
```

Перед выполнением массового удаления блоки в списке свободных блоков отсутствуют и есть 1043 *листовых* блока индекса, которые хранят данные. Теперь произведем удаление и снова измерим утилизацию пространства:

```
EODA@ORA12CR1> delete from big_table where id <= 250000;
250000 rows deleted.
250000 строк удалено.
```

```
EODA@ORA12CR1> commit;
Commit complete.
Фиксация завершена.
```

```
EODA@ORA12CR1> declare
2     l_freelist_blocks number;
3 begin
4     dbms_space.free_blocks
5     ( segment_owner => user,
6       segment_name => 'BIG_TABLE_PK',
7       segment_type => 'INDEX',
8       freelist_group_id => 0,
9       free_blks => l_freelist_blocks );
10    dbms_output.put_line( 'blocks on freelist = ' || l_freelist_blocks );
11    dbms_stats.gather_index_stats
12    ( user, 'BIG_TABLE_PK' );
13 end;
14 /
blocks on freelist = 520
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.
```

```
EODA@ORA12CR1> select leaf_blocks from user_indexes where index_name =
'BIG_TABLE_PK';
```

```
LEAF_BLOCKS
-----
523
```

Как видите, теперь примерно половина индекса находится в списке свободных блоков (520 блоков) и есть только 523 листовых блока. Сложение 523 и 520 дает исходные 1043 блока. Это значит, что блоки совершенно пусты и готовы к повторному использованию (блоки в списке свободных блоков для индекса обязательно должны быть пустыми в отличие от блоков в списке свободных блоков для традиционной таблицы).

Эта демонстрация подчеркивает два момента.

- Пространство в блоках индекса применяется повторно, как только появляется строка, которая может его занять.
- Когда блок индекса опустошается, он может быть вынесен из структуры индекса и повторно использован в будущем. Видимо, это и послужило источником мифа: в индексной структуре блоки не отображаются как имеющие свободное пространство, что характерно для таблицы. В случае таблицы вы можете видеть блоки в списке свободных блоков, даже если они содержат в себе некоторые данные. В случае индекса в списке свободных блоков вы будете видеть только абсолютно пустые блоки; блоки, которые имеют хотя бы одну запись индекса (остальное пространство в них свободно), в нем не показаны.

Миф: наиболее отличительные столбцы должны быть в индексе первыми

Данное утверждение выглядит вполне осмысленным. Если вы собираетесь создать индекс на столбцах C1 и C2 в таблице T с 100 000 строк и выясняете, что C1 имеет 100 000 отличающихся значений, а C2 — только 25 000 таких значений, то захотите создать индекс на T(C1, C2). Это значит, что столбец C1 должен быть первым — такой подход соответствует здравому смыслу. На самом деле при сравнении векторов данных (предположим, что комбинация C1, C2 представляет собой вектор) не играет роли, какой столбец находится первым. Рассмотрим следующий пример. Мы создадим таблицу на основе ALL_OBJECTS с индексом на столбцах OWNER, OBJECT_TYPE и OBJECT_NAME (от наименее отличительного к наиболее отличительному), а также на столбцах OBJECT_NAME, OBJECT_TYPE и OWNER:

```
EODA@ORA12CR1> create table t as select * from all_objects;
Table created.
Таблица создана.
```

```
EODA@ORA12CR1> create index t_idx_1 on t(owner,object_type,object_name);
Index created.
Индекс создан.
```

```
EODA@ORA12CR1> create index t_idx_2 on t(object_name,object_type,owner);
Index created.
Индекс создан.
```

```
EODA@ORA12CR1> select count(distinct owner), count(distinct object_type),
2 count(distinct object_name ), count(*)
3 from t;
```

```
COUNT(DISTINCTOWNER) COUNT(DISTINCTOBJECT_TYPE) COUNT(DISTINCTOBJECT_NAME) COUNT(*)
-----
34 36 30813 50253
```

Теперь, чтобы показать, что ни один из них не является более эффективным в отношении пространства, мы измерим утилизацию дискового пространства этими индексами:

```
EODA@ORA12CR1> analyze index t_idx_1 validate structure;
Index analyzed.
Индекс проанализирован.
```

```
EODA@ORA12CR1> select btree_space, pct_used, opt_cmpr_count, opt_cmpr_pctsave
2 from index_stats;
```

```
BTREE_SPACE PCT_USED OPT_CMPR_COUNT OPT_CMPR_PCTSAVE
-----
2526832 89 2 28
```

```
EODA@ORA12CR1> analyze index t_idx_2 validate structure;
Index analyzed.
Индекс проанализирован.
```



```
EODA@ORA12CR1> select btree_space, pct_used, opt_cmp_r_count, opt_cmp_r_pctsave
2 from index_stats;
```

```
BTREE_SPACE  PCT_USED  OPT_CMPR_COUNT  OPT_CMPR_PCTSAVE
-----
2510776      90          0          0
```

Индексы занимают почти одинаковые объемы пространства — особой разницы здесь нет. Однако первый индекс является намного более *сжимаемым*, если применяется сжатие ключей, что доказывает значение OPT_CMPR_PCTSAVE. Это аргумент в пользу организации столбцов внутри индекса в порядке от наименее отличительных до наиболее отличительных. Теперь давайте посмотрим, как индексы работают, чтобы узнать, какой из них в целом более эффективен. Для этого мы будем использовать блок PL/SQL с запросами, снабженными подсказками (чтобы явно задействовать то один индекс, то другой):

```
EODA@ORA12CR1> alter session set sql_trace=true;
Session altered.
Сеанс изменен.
```

```
EODA@ORA12CR1> declare
2     cnt int;
3     begin
4     for x in ( select /*+FULL(t)*/ owner, object_type, object_name from t )
5     loop
6         select /*+ INDEX( t t_idx_1 ) */ count(*) into cnt
7         from t
8         where object_name = x.object_name
9         and object_type = x.object_type
10        and owner = x.owner;
11
12        select /*+ INDEX( t t_idx_2 ) */ count(*) into cnt
13        from t
14        where object_name = x.object_name
15        and object_type = x.object_type
16        and owner = x.owner;
17    end loop;
18 end;
19 /
```

PL/SQL procedure successfully completed.

Процедура PL/SQL успешно завершена.

Запросы читают каждую отдельную строку таблицы посредством индекса. Отчет TKPROF дает следующие сведения:

```
SELECT /*+ INDEX( t t_idx_1 ) */ COUNT(*) FROM T
WHERE OBJECT_NAME = :B3 AND OBJECT_TYPE = :B2 AND OWNER = :B1
```

| call | count | cpu | elapsed | disk | query | current | rows |
|---------|--------|------|---------|------|--------|---------|-------|
| Parse | 1 | 0.00 | 0.00 | 0 | 0 | 0 | 0 |
| Execute | 50253 | 0.69 | 0.67 | 0 | 0 | 0 | 0 |
| Fetch | 50253 | 0.46 | 0.49 | 0 | 100850 | 0 | 50253 |
| total | 100507 | 1.15 | 1.16 | 0 | 100850 | 0 | 50253 |

```

Rows (1st) Rows (avg) Rows (max) Row Source Operation
-----
1          1          1 SORT AGGREGATE (cr=2 pr=0 pw=0 time=16 us)
1          1          1 INDEX RANGE SCAN T_IDX_1 (cr=2 pr=0 pw=0 time=12 us)
*****
SELECT /** INDEX( t t_idx_2 ) */ COUNT(*) FROM T
WHERE OBJECT_NAME = :B3 AND OBJECT_TYPE = :B2 AND OWNER = :B1
call      count      cpu    elapsed      disk      query      current      rows
-----
Parse      1          0.00      0.00          0          0          0          0
Execute  50253          0.68      0.66          0          0          0          0
Fetch    50253          0.48      0.48          0     100834          0     50253
-----
total    100507          1.16      1.15          0     100834          0     50253
...
Rows (1st) Rows (avg) Rows (max) Row Source Operation
-----
1          1          1 SORT AGGREGATE (cr=2 pr=0 pw=0 time=13 us)
1          1          1 INDEX RANGE SCAN T_IDX_2 (cr=2 pr=0 pw=0 time=12 us)

```

Запросы обрабатывают в точности одно и то же число строк и очень близкие количества блоков (незначительные отклонения вызваны случайным порядком следования строк в таблице и последующими оптимизациями, произведенными Oracle), требуют равного времени процессора и выполняются примерно за одинаковое время (повторный запуск этих запросов даст немного отличающиеся значения CPU и ELAPSED, но в среднем они будут почти равными). Порядок размещения столбцов в зависимости от степени различия их значений никак не повлиял на эффективность, и, как утверждалось ранее, в случае сжатия ключа индекса наименее селективный столбец имеет смысл указывать первым. Если вы запустите предыдущий пример с конструкцией COMPRESS 2 для индексов, то обнаружите, что первый индекс выполнит около двух третей объема операций ввода-вывода по сравнению со вторым индексом, учитывая природу запроса в этом случае. Тем не менее, факт заключается в том, что решение о помещении столбца C1 перед C2 должно приниматься на основе того, как индекс применяется. При наличии множества запросов вроде показанных ниже имеет больше смысла создать индекс на T (C2, C1):

```

select * from t where c1 = :x and c2 = :y;
select * from t where c2 = :y;

```

Этот единственный индекс может использоваться любым из запросов. Кроме того, в случае применения сжатия индексного ключа (которое мы рассматривали в отношении индекс-таблиц, а дополнительно будем исследовать позже) можно построить индекс меньшего размера, если столбец C2 будет на первом месте. Причина в том, что каждое значение столбца C2 повторяется в индексе в среднем по четыре раза. Если столбцы C1 и C2 имеют среднюю длину 10 байтов, то записи индекса формально займут 2 000 000 байтов (100 000 × 20). Используя сжатие индексного ключа на (C2, C1), мы можем сократить размер индекса до 1 250 000 байтов (100 000 × 12,5), т.к. три из четырех повторений C2 удастся исключить.

В Oracle 5 (да, версии 5) существовал довод в пользу размещения наиболее селективных столбцов в начале индекса. Это было связано с реализованным в Oracle 5 способом сжатия индекса (вовсе не то же самое, что сжатие индексного ключа). С добавлением низкоуровневого блокирования в версии Oracle 6 это средство было изъято. С тех пор утверждение о том, что помещение в начало индекса наиболее отличительных столбцов делает индекс меньше или эффективнее, не соответствует

действительности. На первый взгляд так должно быть, но на самом деле нет. Со сжатием индексных ключей появился аргумент в пользу противоположного решения, поскольку теперь это приводит к уменьшению размера индекса. Однако, как утверждалось ранее, решение должно управляться тем, *каким образом* применяется индекс.

Резюме

В этой главе были описаны различные типы индексов, предлагаемые Oracle. Мы начали с базового индекса со структурой В-дерева и взглянули на его разнообразные подтипы, такие как индексы по реверсированному ключу (предназначенные для Oracle RAC) и индексы, упорядоченные по убыванию, для извлечения данных, которые отсортированы в смешанном порядке по возрастанию и по убыванию. Мы потратили некоторое время на ознакомление с тем, когда должен использоваться индекс и почему в определенных обстоятельствах он может оказаться бесполезным.

Затем мы рассмотрели битовые индексы — великолепный метод индексации данных с низкой и средней кардинальностью в среде хранилища данных (интенсивного по чтению, не OLTP). Были раскрыты ситуации, когда применение битовых индексов оправдано, и объяснены причины, по которым их не следует применять в среде OLTP или любой среде, где множество пользователей должны параллельно обновлять один и тот же столбец.

Далее мы обратились к индексам на основе функций, которые на самом деле являются специальными случаями индексов со структурой В-дерева и битовых индексов. Это индексы на функциях, примененных к столбцу (или столбцам), т.е. мы можем предварительно вычислить и сохранить результаты сложных вычислений и пользовательских функций с целью их последующего исключительно быстрого извлечения. Мы рассмотрели некоторые важные детали реализации, связанные с индексами на основе функций, такие как необходимые настройки уровня системы и сеанса, которые должны быть на месте, чтобы индексы можно было использовать. Мы исследовали примеры применения индексов на основе функций — как встроенных функций Oracle, так и пользовательских функций. Наконец, было приведено несколько предостережений, касающихся индексов на основе функций.

После этого мы взглянули на высокоспециализированный тип индекса, который называется индексом предметной области. Вместо того чтобы углубляться в детали построения одного из таких индексов с чистого листа (что подразумевает длинную и сложную последовательность событий), мы рассмотрели пример, который уже реализован: текстовый индекс. Затем мы обсудили пару тем, относящихся к версии Oracle 12c: индексация расширенных столбцов и создание нескольких индексов на тех же самых комбинациях столбцов. Индексация расширенных столбцов требует применения либо виртуального столбца и связанного с ним индекса, либо индекса на основе функции. При индексации тех же самых комбинаций столбцов должны использоваться разные физические типы индексов, и только один индекс может быть выбран в качестве видимого.

Глава завершилась ответами на самые часто задаваемые вопросы об индексах, а также ознакомлением с рядом мифов, касающихся индексов. Здесь были раскрыты темы, начиная с простого вопроса вроде того, работают ли индексы с представлениями, до более сложного и тонкого утверждения о том, что пространство в индексе никогда не используется повторно. На эти вопросы были даны исчерпывающие ответы, а мифы были развеяны главным образом с помощью примеров, в которых демонстрировались соответствующие концепции.

ГЛАВА 12

Типы данных

Выбор правильного типа данных кажется очень простым и легким делом, но я часто вижу, что это делается некорректно. Самое основополагающее решение — какой тип использовать для хранения данных — будет оказывать влияние на ваши приложения и данные в течение многих лет. Таким образом, выбор подходящего типа данных является первостепенным. Его также трудно изменить после свершившегося факта — другими словами, однажды реализовав его, вы будете связаны с ним надолго.

В этой главе мы рассмотрим все доступные базовые типы данных Oracle и обсудим, как они реализованы, и когда целесообразно применять тот или иной тип. Мы не будем исследовать типы, определяемые пользователем, поскольку они представляют собой составные объекты, производные от встроенных типов данных Oracle. Мы выясним, что происходит, когда при решении задачи используется ошибочный тип данных или даже неправильный параметр для подходящего типа данных (длина, точность, масштаб и т.д.). К концу этой главы вы получите представление о доступных типах, о том, как они реализованы, когда нужно применять каждый из них и — самое важное — почему ключом к успешному решению задачи является использование подходящего типа данных.

Обзор типов данных Oracle

В Oracle предлагаются 22 разных типа данных SQL. Ниже представлен их перечень с краткими описаниями.

- **CHAR.** Символьная строка фиксированной длины, которая будет дополняться пробелами до своей максимальной длины. Не равное NULL поле типа CHAR(10) всегда будет содержать 10 байтов информации согласно стандартным установкам NLS (National Language Support — поддержка национальных языков). Вскоре мы раскроем реализацию NLS более подробно. Поле типа CHAR может хранить до 2000 байтов информации.
- **NCHAR.** Символьная строка фиксированной длины, содержащая данные в формате Unicode. Здесь Unicode — это стандарт кодирования символов, разработанный консорциумом Unicode (Unicode Consortium) с целью предоставления универсального способа для кодирования символов любого языка, не зависящего от применяемой компьютерной системы или платформы. Тип NCHAR позволяет базе данных хранить информацию в двух разных наборах символов: типы CHAR и NCHAR используют, соответственно, набор символов базы данных

и национальный набор символов. Отличное от NULL поле типа NCHAR(10) будет всегда содержать 10 символов информации (обратите внимание на отличие от типа CHAR в этом отношении). Тип NCHAR может хранить до 2000 байтов информации.

- VARCHAR2. В настоящее время также является синонимом типа VARCHAR. Это символьная строка переменной длины, которая отличается от типа CHAR тем, что не дополняется пробелами до своей максимальной длины. Поле типа VARCHAR(10) может содержать от 0 до 10 байтов информации с применением стандартных установок NLS. Поле VARCHAR2 позволяет хранить до 4000 байтов информации. Начиная с версии Oracle 12c, тип VARCHAR2 может быть сконфигурирован для хранения до 32 767 байтов информации (за дополнительными деталями обращайтесь в раздел “Расширенные типы данных” далее в этой главе).
- NVARCHAR2. Символьная строка переменной длины, содержащая данные в формате Unicode. Поле типа NVARCHAR2(10) может хранить от 0 до 10 символов информации. Тип NVARCHAR2 позволяет содержать до 4000 байтов информации. Начиная с версии Oracle 12c, тип NVARCHAR2 может быть сконфигурирован для хранения до 32 767 байтов информации (за дополнительными деталями обращайтесь в раздел “Расширенные типы данных” далее в этой главе).
- RAW. Двоичный тип данных переменной длины; над данными, хранящимися в этом типе, никаких преобразований между наборами символов не происходит. Он считается строкой двоичных байтов информации, которая будет просто храниться в базе данных. Тип RAW позволяет хранить до 2000 байтов информации. Начиная с версии Oracle 12c, тип RAW может быть сконфигурирован для хранения до 32 767 байтов информации (за дополнительными деталями обращайтесь в раздел “Расширенные типы данных” далее в этой главе).
- NUMBER. Этот тип данных способен хранить числа с точностью до 38 цифр. Числа могут варьироваться от $1,0 \times 10^{-130}$ до $1,0 \times 10^{126}$, не включая его. Каждое число хранится в поле переменной длины от 0 (для NULL) до 22 байтов. Типы NUMBER в Oracle очень точны — намного точнее, чем обычные типы FLOAT и DOUBLE, присутствующие во многих языках программирования.
- BINARY_FLOAT. Это тип, доступный в Oracle 10g Release 1 и последующих версиях. Представляет 32-битовое число с плавающей точкой одинарной точности. Может поддерживать минимум 6 цифр точности и занимает 5 байтов при хранении на диске.
- BINARY_DOUBLE. Это тип, доступный в Oracle 10g Release 1 и последующих версиях. Представляет 64-битовое число с плавающей точкой двойной точности. Может поддерживать минимум 15 цифр точности и занимает 9 байтов при хранении на диске.
- LONG. Этот тип способен хранить вплоть до 2 Гбайт символьных данных (не символов, т.к. каждый символ может занимать несколько байтов в многобайтном наборе символов). Типы LONG имеют множество ограничений (обсуждаются позже), предназначенных для обеспечения обратной совместимости, поэтому его настоятельно не рекомендуется применять в новых приложениях. Когда возможно, преобразовывайте типы LONG в CLOB в существующих приложениях.

- **LONG RAW.** Тип LONG RAW способен хранить до 2 Гбайт двоичной информации. По тем же причинам, указанным для типов LONG, вместо него рекомендуется использовать тип BLOB как в будущих разработках, так и по возможности в существующих приложениях.
- **DATE.** Тип данных фиксированной ширины в 7 байтов для хранения значений даты/времени. Всегда содержит семь атрибутов — век, год внутри века, месяц, день месяца, часы, минуты и секунды.
- **TIMESTAMP.** Тип данных фиксированной ширины в 7 или 11 байтов (в зависимости от точности) для хранения значений даты/времени. Отличается от типа данных DATE тем, что может содержать дробные части секунды, для представления которых отведено до 9 цифр после десятичной точки.
- **TIMESTAMP WITH TIME ZONE.** Тип фиксированной ширины в 13 байтов для хранения значений даты/времени, который также предоставляет поддержку часового пояса (TIME ZONE). Дополнительная информация о часовом поясе хранится вместе с данными даты/времени.
- **TIMESTAMP WITH LOCAL TIME ZONE.** Тип фиксированной ширины в 7 или 11 байтов (в зависимости от точности) для хранения значений даты/времени, похожий на TIMESTAMP, но чувствительный к часовому поясу. При выполнении модификации в базе данных принимается во внимание часовой пояс, указанный с данными, и компонент даты/времени нормализуется в соответствии с локальным часовым поясом базы данных. Таким образом, если вы вставили значение даты/времени с применением часового пояса U.S./Pacific, а часовым поясом базы данных являлся U.S./Eastern, то финальная информация даты/времени должна быть преобразована в часовой пояс U.S./Eastern и сохранена как значение TIMESTAMP. При извлечении хранящееся в базе данных значение TIMESTAMP должно быть преобразовано в значение времени, соответствующее часовому поясу текущего сеанса.
- **INTERVAL YEAR TO MONTH.** Тип данных фиксированной ширины в 5 байтов для хранения продолжительности промежутка времени — в этом случае количества лет и месяцев. Интервалы можно использовать в арифметике дат для добавления или вычитания периода времени из значений типа DATE или TIMESTAMP.
- **INTERVAL DAY TO SECOND.** Тип данных фиксированной ширины в 11 байтов для хранения продолжительности промежутка времени — в этом случае количества дней, часов, минут и секунд и дополнительно до 9 цифр дробной части секунды.
- **BLOB.** Этот тип данных предназначен для хранения до 4 Гбайт данных в Oracle9i и предшествующих версиях или $(4 \text{ Гбайт} - 1) \times$ (размер блока базы данных) байтов данных в Oracle 10g и последующих версиях. Значения типа BLOB содержат “двоичную” информацию, которая не подвергается преобразованиям символьных наборов. Тип BLOB может быть подходящим типом для хранения электронной таблицы, документа текстового процессора, файлов графических изображений и т.п.

- CLOB. Этот тип данных предназначен для хранения до 4 Гбайт данных в Oracle9i и предшествующих версиях или $(4 \text{ Гбайт} - 1) \times (\text{размер блока базы данных})$ байтов данных в Oracle 10g и последующих версиях. Значения типа CLOB содержат информацию, подвергаемую преобразованиям символьных наборов. Тип CLOB может быть подходящим типом для хранения крупного объема простого текста. Обратите внимание, что речь идет о крупном объеме простого текста; этот тип данных не подходит для текста объемом 4000 байтов или меньше — в таком случае следует применять тип VARCHAR2.
- NCLOB. Этот тип данных предназначен для хранения до 4 Гбайт данных в Oracle9i и предшествующих версиях или $(4 \text{ Гбайт} - 1) \times (\text{размер блока базы данных})$ байтов данных в Oracle 10g и последующих версиях. Значения типа NCLOB содержат информацию, закодированную с использованием национального набора символов базы данных, и подобно значениям типа CLOB подвергаются преобразованиям символьных наборов.
- BFILE. Этот тип данных позволяет хранить в столбце базы данных объект каталога Oracle (указатель на каталог операционной системы) и имя файла, а также читать такой файл. В итоге появляется возможность эффективной работы с файлами операционной системы, доступными на сервере базы данных, в режиме только для чтения, как если бы они хранились в самой таблице базы данных.
- ROWID. Это 10-байтный адрес строки в базе данных. Значение ROWID содержит достаточно информации для нахождения строки на диске, а также для идентификации объекта, на который указывает ROWID (таблицы и т.д.).
- UROWID. Это универсальный идентификатор строки (ROWID), применяемый для таблиц, таких как индекс-таблицы и таблицы, доступные через шлюзы к гетерогенным базам данных, которые не имеют фиксированных идентификаторов строк. Значение типа UROWID является представлением значения первичного ключа строки, и потому будет варьироваться по размеру в зависимости от объекта, на который оно указывает.

Вполне очевидно, что в приведенный список не вошли многие типы, такие как INT, INTEGER, SMALLINT, FLOAT, REAL и другие. В действительности эти типы реализованы на основе одного из базовых типов, перечисленных в списке, т.е. являются синонимами собственных типов Oracle. Кроме того, в списке отсутствуют типы данных наподобие XMLType, SYS.ANYTYPE и SDO_GEOMETRY, поскольку в книге они не рассматриваются. Они представляют собой сложные объектные типы, включающие коллекцию атрибутов наряду с методами (функциями), которые оперируют на этих атрибутах. Они состоят из указанных ранее базовых типов и не являются настоящими типами данных в общепринятом смысле, а скорее реализациями, или наборами функциональности, которые можно использовать в своих приложениях.

А теперь давайте рассмотрим перечисленные базовые типы более подробно.

Символьные и двоичные строковые типы

Символьными типами данных в Oracle являются CHAR, VARCHAR2 и их варианты “N”. Типы CHAR и NCHAR способны хранить до 2000 байтов текста. Типы VARCHAR2 и NVARCHAR2 могут хранить до 4000 байтов информации.

На заметку! Начиная с версии Oracle 12c, типы данных VARCHAR2, NVARCHAR2 и RAW могут быть сконфигурированы для хранения вплоть до 32 767 байтов информации. Расширенные типы данных по умолчанию не включены; следовательно, если явная настройка не производилась, то максимальный размер по-прежнему составляет 4000 байтов для типов VARCHAR2 и NVARCHAR2 и 2000 байтов для RAW. Дополнительные сведения приведены в разделе “Расширенные типы данных” далее в этой главе.

Упомянутый текст преобразуется между разнообразными наборами символов в соответствии с нуждами базы данных. *Набор символов* (character set) — это двоичное представление индивидуальных символов в битах и байтах. Доступно много разных наборов символов, каждый из которых способен представлять различные символы. Ниже перечислены примеры.

- Набор символов US7ASCII — стандартное представление ASCII для 128 символов. Для их представления применяются младшие 7 битов байта.
- Набор WE8MSWIN1252 — символьный набор Западной Европы, который может представлять 128 символов ASCII, а также 128 расширенных символов с использованием всех 8 битов байта.

Прежде чем мы погрузимся в детали типов данных CHAR, VARCHAR2 и их вариантов “N”, полезно получить базовое понимание того, что означают различные наборы символов.

Обзор NLS

Как было указано ранее, NLS означает *National Language Support* (поддержка национальных языков). Это очень мощное средство базы данных, но является одним из тех, которые часто понимают недостаточно хорошо. NLS управляет многими аспектами данных. Например, NLS влияет на то, как данные сортируются, присутствуют ли запятые между группами цифр и точка между целой и дробной частями чисел (т.е. 1,000,000.01) либо наоборот — точки между группами цифр и запятая между целой и дробной частями (т.е. 1.000.000,01). Но более важно то, что NLS управляет следующими аспектами:

- кодировка текстовых данных при их постоянном хранении на диске;
- прозрачное преобразование данных из одного набора символов в другой.

Эта прозрачная часть работы обычно смущает людей больше всего — она настолько прозрачна, что вы даже не можете реально увидеть, происходит ли она. Рассмотрим небольшой пример.

Предположим, что вы храните в своей базе 8-битовые данные с применением набора символов WE8MSWIN1252, но есть несколько клиентов, которые подключаются с использованием 7-битового набора символов, такого как US7ASCII. Эти клиенты не ожидают 8-битовых данных и нуждаются в преобразовании данных, поступаю-

ших из базы, в нечто такое, с чем они могут работать. Хотя это звучит удивительно, но если вы не знаете о том, что подобное происходит, то со временем можете обнаружить, что данные теряют символы, т.к. символы, не доступные в US7ASCII, транслируются в другие символы, которые в наборе US7ASCII доступны. Это связано с выполняемой трансляцией символьных наборов. Короче говоря, если вы извлекаете данные из базы в символьном наборе 1, преобразуете их в символьный набор 2, а затем вставляете обратно (выполняя обратный процесс), то имеете высокие шансы существенно модифицировать данные. Как правило, преобразование символьных наборов — это процесс, который изменяет данные, и обычно приходится отображать большие наборы символов (набор 8-битовых символов в текущем примере) на меньшие (7-битовые символы). Это *преобразование с потерей информации* — символы изменяются просто из-за невозможности представления каждого символа. Но такие преобразования должны происходить. Если база хранит данные в однобайтовом наборе символов, но клиент (скажем, Java-приложение, поскольку в языке Java применяется Unicode) ожидает их в многобайтном представлении, то данные должны быть преобразованы, чтобы клиент смог работать с ними.

Увидеть преобразование символьных наборов в действии очень легко. Например, у меня есть база данных, набор символов которой установлен в WE8MSWIN1252 — типичный набор символов для Западной Европы:

```
EODA@ORA12CR1> select *
  2 from nls_database_parameters
  3 where parameter = 'NLS_CHARACTERSET';
```

| PARAMETER | VALUE |
|------------------|--------------|
| NLS_CHARACTERSET | WE8MSWIN1252 |

Теперь я удостоверяюсь, что параметр NLS_LANG установлен в тот же набор символов, который используется в базе данных (пользователи Windows могут изменить/посмотреть эту установку в своих реестрах):

```
EODA@ORA12CR1> host echo $NLS_LANG
AMERICAN_AMERICA.WE8MSWIN1252
```

Затем я могу создать таблицу и поместить в нее какие-нибудь “8-битовые” данные. Эти данные не могут использоваться 7-битовым клиентом, который ожидает только 7-битовых данных ASCII:

```
EODA@ORA12CR1> create table t ( data varchar2(1) );
Table created.
Таблица создана.
EODA@ORA12CR1> insert into t values ( chr(224) );
1 row created.
1 строка создана.
EODA@ORA12CR1> insert into t values ( chr(225) );
1 row created.
1 строка создана.
EODA@ORA12CR1> insert into t values ( chr(226) );
1 row created.
1 строка создана.
```

```

EODA@ORA12CR1> select data, dump(data) dump from t;

D DUMP
-----
à Typ=1 Len=1: 224
á Typ=1 Len=1: 225
â Typ=1 Len=1: 226

EODA@ORA12CR1> commit;

```

На заметку! Если при воспроизведении этого примера вы не увидите приведенный выше вывод, убедитесь, что в программном обеспечении терминального клиента применяется набор символов UTF-8. В противном случае оно может транслировать символы при выводе их на экран! Распространенный эмулятор терминала для UNIX обычно работает с 7-битовым набором символов ASCII. Это касается как пользователей Windows, так и UNIX/Linux. Удостоверьтесь, что ваш терминал способен отображать соответствующие символы.

Если я перейду в другое окно и укажу там клиента с 7-битовым набором ASCII, то получу другие результаты:

```

$ export NLS_LANG=AMERICAN_AMERICA.US7ASCII
$ sqlplus eoda
Enter password:

EODA@ORA12CR1> select data, dump(data) dump from t;

D DUMP
-----
a Typ=1 Len=1: 224
a Typ=1 Len=1: 225
a Typ=1 Len=1: 226

```

Обратите внимание, что в 7-битовом сеансе я трижды получил букву а без диакритических знаков. Тем не менее, функция DUMP показывает, что на самом деле в базе данных хранятся три совершенно разных символа, а не одна и та же буква а. Данные в базе не изменились — поменялись только значения, получаемые клиентом. Если клиент извлечет эти данные в хост-переменные, как показано ниже:

```

EODA@ORA12CR1> variable d varchar2(1)
EODA@ORA12CR1> variable r varchar2(20)
EODA@ORA12CR1> begin
2  select data, rowid into :d, :r from t where rownum = 1;
3  end;
4  /
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.

```

и ничего не сделает с ними, а просто отправит их обратно в базу данных:

```

EODA@ORA12CR1> update t set data = :d where rowid = chartorowid(:r);
1 row updated.
1 строка обновлена.

EODA@ORA12CR1> commit;
Commit complete.
Фиксация завершена.

```

то в исходном 8-битовом сеансе можно будет заметить, что один из исходных символов потерян. Код буквы à, который был раньше, заменяется семью младшими битами, давая в итоге a:

```

EODA@ORA12CR1> select data, dump(data) dump from t;
D DUMP
-----
a Typ=1 Len=1: 97
à Typ=1 Len=1: 225
â Typ=1 Len=1: 226

```

Функция DUMP в Oracle SQL

Функция DUMP в Oracle SQL позволяет отображать код типа данных, длину в байтах и внутреннее представление значения данных (и также дополнительно имя набора символов). Его синтаксис выглядит следующим образом:

```
DUMP( выражение [, формат_возврата] [, начальная_позиция] [, длина] )
```

Стандартным значением параметра формат_возврата является 10 (десятичный) и он может принимать любое из перечисленных далее значений: 8, 10, 16, 17, 1008, 1010, 1016 или 1017. Здесь 8 обозначает восьмеричный формат, 10 — десятичный, 16 — шестнадцатеричный, 17 — одиночные символы, 1008 — восьмеричный формат с именем набора символов, 1010 — десятичный формат с именем набора символов, 1016 — шестнадцатеричный формат с именем набора символов и 1017 — одиночные символы с именем набора символов. В приведенном ниже примере выводится информация, касающаяся символа a:

```

EODA@ORA12CR1> select dump('a'), dump('a',8), dump('a',16) from dual;
DUMP('A')          DUMP('A',8)          DUMP('A',16)
-----
Typ=96 Len=1: 97   Typ=96 Len=1: 141   Typ=96 Len=1: 61

```

Значения 97, 141 и 61 — это соответствующие ASCII-коды для символа a в десятичной, восьмеричной и шестнадцатеричной форме. Возвращенный код типа данных Typ=96 указывает на тип данных CHAR (полный список кодов типов данных Oracle вместе с их описаниями доступен в руководстве по языку SQL для Oracle (*Oracle Database SQL Language Reference*)).

Пример демонстрирует прямое влияние среды с гетерогенным набором символов, когда клиенты и база данных используют разные настройки NLS. Об этом следует помнить, поскольку такой эффект проявляется во многих ситуациях. Например, если администратор базы данных применяет инструмент EXP для извлечения информации, то может наблюдать следующее предупреждение:

```

[tkyte@desktop tkyte] exp userid=eoda tables=t
Export: Release 12.1.0.1.0 - Production on Thu Jan 9 16:11:24 2014
Copyright (c) 1982, 2013, Oracle and/or its affiliates. All rights reserved.
Password:
Connected to: Oracle Database 12c Enterprise Edition Release 12.1.0.1.0 -
64bit Production
With the Partitioning, OLAP, Advanced Analytics and Real Application
Testing options
Export done in US7ASCII character set and UTF8 NCHAR character set
server uses WE8MSWIN1252 character set (possible charset conversion)
About to export specified tables via Conventional Path

```

сервер использует набор символов WE8MSWIN1252 (возможно преобразование набора символов)
 Об обычном экспорте указанных таблиц
 ...

К таким предупреждениям нужно относиться очень серьезно. Если вы экспортируете эту таблицу, намереваясь затем удалить ее и воссоздать с помощью инструмента IMP, то обнаружите, что *все* данные в таблице будут состоять из 7-битовых символов! Остерегайтесь непредвиденных преобразований наборов символов.

На заметку! Проблема с непредвиденными преобразованиями символов затрагивает не каждый инструмент, а если и затрагивает, то по-разному. Например, при использовании рекомендованного процесса экспорта/импорта Data Pump вы обнаружите, что экспорт всегда производится с набором символов базы данных независимо от установок NLS на стороне клиента. Причина в том, что Data Pump запускается на самом сервере базы данных; это не инструмент клиентской стороны. Аналогично, во время импорта Data Pump будет всегда преобразовывать данные в файле, подлежащем импорту, из набора символов исходной базы данных в набор символов целевой базы данных. Другими словами, преобразование набора символов в Data Pump по-прежнему возможно (если исходная и целевая базы данных имеют разные наборы символов), но не в той же манере, как это делалось унаследованными инструментами EXP/IMP.

Но вы также должны понимать, что в общем случае преобразования наборов символов необходимы. Если клиенты ожидают поступления данных в определенном наборе символов, то весьма нежелательно отправлять им информацию с отличающимся набором символов.

На заметку! Я настоятельно рекомендую всем ознакомиться с документом *Oracle Globalization Support Guide* (Руководство по поддержке глобализации в Oracle). Вопросы, связанные с NLS, в нем раскрыты намного глубже, чем это сделано здесь. Любой, кто создает приложения, которые будут применяться по всему миру (или даже просто в нескольких странах), должен овладеть информацией, содержащейся в этом документе.

Теперь, имея беглое представление о наборах символов и оказываемом ими влиянии, давайте обратимся к строковым типам данных, которые предлагаются в Oracle.

Символьные строки

В Oracle имеются четыре базовых строковых типа: CHAR, VARCHAR2, NCHAR и NVARCHAR2. Все строки в Oracle хранятся в одном и том же формате. В блоке базы данных они будут иметь начальное поле длины размером от 1 до 3 байтов, за которым следуют сами данные; когда строка равна NULL, она представляется единственным байтом 0xFF.

На заметку! В Oracle хвостовые столбцы NULL потребляют 0 байтов дискового пространства. Это значит, что если последний столбец в таблице содержит NULL, то Oracle ничего в нем не сохраняет. Если два последних столбца содержат NULL, для них также ничего не хранится. Но если после столбца NULL идет столбец, отличный от NULL, то Oracle будет использовать NULL-флаг, описанный в настоящем разделе, чтобы указать пропущенное значение.

Если длина строки меньше или равна 250 (от 0x01 до 0xFA), то для представления длины Oracle будет применять 1 байт. Все строки, длина которых превышает 250 байтов, снабжены байтом-флагом 0xFE, за которым следуют 2 байта, представляющие длину. Таким образом, столбец типа VARCHAR2(80), содержащий слова Hello World, может выглядеть в блоке так, как показано на рис. 12.1.

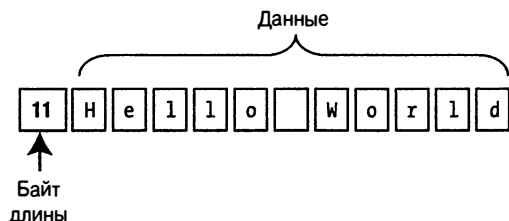


Рис. 12.1. Строка Hello World, хранящаяся в столбце VARCHAR2(80)

С другой стороны, столбец типа CHAR(80), хранящий те же данные, мог бы выглядеть подобно приведенному на рис. 12.2.

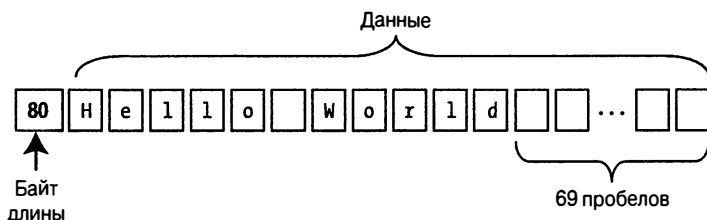


Рис. 12.2. Строка Hello World, хранящаяся в столбце CHAR(80)

Тот факт, что CHAR/NCHAR на самом деле являются всего лишь замаскированными типами VARCHAR2/NVARCHAR2, приводит меня к заключению, что на самом деле существуют только два типа символьных строк, которые нужно принимать во внимание — VARCHAR2 и NVARCHAR2. Ни в одном из своих приложений я никогда не находил повода использовать тип CHAR. Поскольку тип CHAR *всегда* дополняет результирующую строку пробелами до фиксированной длины, нетрудно понять, что он потребляет максимум пространства как в табличном сегменте, так и в любых сегментах индексов. Уже это довольно плохо, но есть и другая причина избегать применения типов CHAR/NCHAR: они приводят к путанице в приложениях, которые должны извлекать эту информацию (многие не могут “найти” свои данные после того, как сохранили их). Причина связана с правилами сравнения символьных строк и строгостью их выполнения. Давайте в целях демонстрации воспользуемся строкой 'Hello World' в простой таблице:

```
EODA@ORA12CR1> create table t
2 ( char_column      char(20),
3   varchar2_column  varchar2(20)
4 )
5 /
Table created.
Таблица создана.
```

```
EODA@ORA12CR1> insert into t values ( 'Hello World', 'Hello World' );
1 row created.
1 строка создана.
```

```
EODA@ORA12CR1> select * from t;
CHAR_COLUMN          VARCHAR2_COLUMN
-----
Hello World          Hello World
```

```
EODA@ORA12CR1> select * from t where char_column = 'Hello World';
CHAR_COLUMN          VARCHAR2_COLUMN
-----
Hello World          Hello World
```

```
EODA@ORA12CR1> select * from t where varchar2_column = 'Hello World';
CHAR_COLUMN          VARCHAR2_COLUMN
-----
Hello World          Hello World
```

Пока что столбцы выглядят идентичными, но на самом деле произошли некоторые неявные преобразования, и литерал 'Hello World' типа CHAR(11) был превращен в CHAR(20) с дополнением пробелами, когда сравнить его со столбцом CHAR. Так и должно быть, поскольку 'Hello World' — это *не то же самое*, что 'Hello World' без хвостовых пробелов. Мы можем подтвердить, что эти две строки по своей сути разные:

```
EODA@ORA12CR1> select * from t where char_column = varchar2_column;
no rows selected
строки не выбраны
```

Строки не равны друг другу. Для того чтобы они стали равными, понадобится либо дополнить пробелами VARCHAR2_COLUMN до длины в 20 символов, либо отбросить хвостовые пробелы в столбце CHAR_COLUMN, как показано ниже:

```
EODA@ORA12CR1> select * from t where trim(char_column) = varchar2_column;
CHAR_COLUMN          VARCHAR2_COLUMN
-----
Hello World          Hello World
```

```
EODA@ORA12CR1> select * from t where char_column = rpad( varchar2_column,
20 );
CHAR_COLUMN          VARCHAR2_COLUMN
-----
Hello World          Hello World
```

На заметку! Существует множество способов дополнения пробелами VARCHAR2_COLUMN, например, с помощью функции CAST().

Эта проблема возникает в приложениях, где применяются строки переменной длины, когда они производят привязку к входным данным, что в результате приводит к невозможности нахождения данных:

```
EODA@ORA12CR1> variable varchar2_bv varchar2(20)
EODA@ORA12CR1> exec :varchar2_bv := 'Hello World';
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.
```

```

EODA@ORA12CR1> select * from t where char_column = :varchar2_bv;
no rows selected
строки не выбраны

EODA@ORA12CR1> select * from t where varchar2_column = :varchar2_bv;
CHAR_COLUMN          VARCHAR2_COLUMN
-----
Hello World          Hello World

```

Итак, здесь поиск строки по столбцу VARCHAR2 работает, а по столбцу CHAR — нет. Переменная привязки VARCHAR2 не расширяется до CHAR(20), как это делает строковый литерал. В данный момент у многих программистов формируется убеждение, что переменные привязки не работают, и нужно использовать литералы. На самом деле это было бы очень неудачным решением. Решение заключается в привязке с применением типа CHAR:

```

EODA@ORA12CR1> variable char_bv char(20)
EODA@ORA12CR1> exec :char_bv := 'Hello World';

PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.

EODA@ORA12CR1>
EODA@ORA12CR1> select * from t where char_column = :char_bv;
CHAR_COLUMN          VARCHAR2_COLUMN
-----
Hello World          Hello World
EODA@ORA12CR1> select * from t where varchar2_column = :char_bv;
no rows selected
строки не выбраны

```

Однако при сочетании типов VARCHAR2 и CHAR вы будете сталкиваться с этой проблемой постоянно. Вдобавок разработчику теперь приходится принимать во внимание ширину поля в своих приложениях. Если разработчик предпочтет реализовать трюк с RPAD для преобразования переменной привязки в тип, который может сравниваться с полем типа CHAR (разумеется, дополнение пробелами переменной привязки предпочтительнее применения функции TRIM к столбцу базы данных, т.к. это легко делает невозможным использование индексов, определенных на этом столбце), то ему все равно придется со временем учитывать изменения ширины столбца. Если размер поля изменится, то это затронет приложение, потому что в нем придется менять ширину поля.

Именно по этим причинам — хранение данных с фиксированной длиной, приводящее к тому, что таблицы и связанные индексы становятся намного больше нормы, вкупе с проблемой с переменными привязки — я избегаю применения типа CHAR в любых обстоятельствах. Я не могу найти довод в его пользу даже в случае односимвольного поля, т.к. здесь действительно нет никакой существенной разницы. Типы VARCHAR2(1) и CHAR(1) идентичны во всех отношениях. Нет никакой важной причины использовать тип CHAR в данной ситуации, поэтому во избежание любой путаницы я просто говорю “нет” полю типа CHAR(1).

Синтаксис символьных строк

Синтаксис для четырех базовых строковых типов является простым (табл. 12.1).

Таблица 12.1. Четыре базовых строковых типа

| Строковый тип | Параметр <i><РАЗМЕР></i> |
|--|--|
| <code>VARCHAR2 (<РАЗМЕР> <BYTE CHAR>)</code> | Число между 1 и 4000 для хранения до 4000 байтов. В следующем разделе мы исследуем отличия и особенности применения модификатора <code>BYTE</code> по сравнению с <code>CHAR</code> в этой конструкции. Начиная с версии Oracle 12c, тип <code>VARCHAR2</code> можно сконфигурировать для хранения до 32 767 байтов информации |
| <code>CHAR (<РАЗМЕР> <BYTE CHAR>)</code> | Число между 1 и 2000 для хранения до 2000 байтов |
| <code>NVARCHAR2 (<РАЗМЕР>)</code> | Число больше 0, верхняя граница которого диктуется используемым национальным набором символов. Начиная с версии Oracle 12c, тип <code>NVARCHAR2</code> можно сконфигурировать для хранения до 32 767 байтов информации |
| <code>NCHAR (<РАЗМЕР>)</code> | Число больше 0, верхняя граница которого диктуется применяемым национальным набором символов |

Байты или символы

Типы `VARCHAR2` и `CHAR` поддерживают два метода указания длины.

- В байтах: `VARCHAR2 (10 BYTE)`. Будет поддерживать до 10 байтов данных, что может соответствовать всего двум символам в многобайтном символьном наборе. Помните, что в многобайтном символьном наборе байты и символы — это не одно и то же!
- В символах: `VARCHAR2 (10 CHAR)`. Будет поддерживать до 10 символов данных, что может соответствовать 40 байтам информации. Кроме того, `VARCHAR2 (4000 CHAR)` теоретически мог бы поддерживать до 4000 символов данных, но поскольку строковый тип данных в Oracle ограничен 4000 байтов, сохранить настолько много символов может не получиться. Ниже в разделе приведен пример.

При использовании многобайтного набора символов вроде UTF8 в определениях `VARCHAR2/CHAR` я советую применять модификатор `CHAR`, т.е. записывать `VARCHAR2 (80 CHAR)`, а не `VARCHAR2 (80)`, т.к. ваше намерение, скорее всего, заключается в определении столбца, который фактически может хранить 80 символов данных. С помощью параметра `NLS_LENGTH_SEMANTICS` на уровне сеанса или системы можно изменить стандартное поведение с `BYTE` на `CHAR`. Я не рекомендую изменять эту настройку на уровне системы; лучше ее использовать как часть команд `ALTER SESSION` в сценариях установки схемы базы данных. Любое приложение, которое требует от базы данных наличие специфического набора установок `NLS`, является недружественным. В большинстве случаев такие приложения не могут быть установлены в базе данных с другими приложениями, которые не требуют таких настроек, а полагаются на стандартные значения.

Важно также помнить о том, что верхний предел количества байтов, хранящихся в `VARCHAR2`, составляет 4000. Однако даже если указать `VARCHAR2 (4000 CHAR)`, то уместить 4000 символов в это поле может не получиться. На самом деле, может оказаться, что в это поле помещаются только 1000 символов, если все они требуют по 4 байта для представления в выбранном наборе символов!

Следующий небольшой пример демонстрирует отличия между BYTE и CHAR, а также работу с верхними границами. Мы создадим таблицу с тремя столбцами, первый из которых имеет длину 1 байт, второй — 1 символ, а третий — 4000 символов. Обратите внимание, что этот тест выполняется в базе данных с многобайтным набором AL23UTF8, который поддерживает последнюю версию стандарта Unicode и кодирует символы с переменной длиной от 1 до 4 байтов для каждого символа:

```
EODA@ORA12CR1> select *
  2   from nls_database_parameters
  3   where parameter = 'NLS_CHARACTERSET';
```

| PARAMETER | VALUE |
|------------------|----------|
| NLS_CHARACTERSET | AL32UTF8 |

```
EODA@ORA12CR1> create table t
  2   ( a varchar2(1),
  3     b varchar2(1 char),
  4     c varchar2(4000 char)
  5   )
  6 /
Table created.
Таблица создана.
```

Если теперь попробовать вставить в эту таблицу одиночный символ длиной 2 байта в UTF, то мы будем наблюдать такое поведение:

```
EODA@ORA12CR1> insert into t (a) values (unistr('\00d6'));
insert into t (a) values (unistr('\00d6'))
      *
```

```
ERROR at line 1:
ORA-12899: value too large for column "EODA"."T"."A" (actual: 2, maximum: 1)
ОШИБКА в строке 1:
ORA-12899: слишком большое значение для столбца "EODA"."T"."A"
(действительное: 2, максимальное: 1)
```

Этот пример иллюстрирует два момента.

- Определение VARCHAR2(1) выражено в байтах, а не символах. Мы имеем одиночный символ Unicode, но он не умещается в один байт.
- Перенеся свое приложение из однобайтного набора символов с фиксированной шириной в многобайтный набор символов, вы можете обнаружить, что текст, который ранее умещался в поля, перестал это делать.

Причина второго момента состоит в том, что 20-символьная строка в однобайтном наборе символов имеет длину 20 байтов и, безусловно, умещается в VARCHAR2(20). Тем не менее, в многобайтном наборе символов 20-символьное поле может иметь длину 80 байтов, поэтому 20 символов Unicode могут не уместиться в 20 байтов. Вам придется либо модифицировать DDL-код, указав определение VARCHAR2(20 CHAR), либо при запуске DDL-кода для создания таблиц применять упомянутый ранее параметр сеанса NLS_LENGTH_SEMANTICS.

Если мы вставим одиночный символ в поле, настроенное на хранение единственного символа, то увидим следующее:

```

EODA@ORA12CR1> insert into t (b) values (unistr('\00d6'));
1 row created.
1 строка создана.

EODA@ORA12CR1> select length(b), lengthb(b), dump(b) dump from t;
LENGTH(B)  LENGTHB(B)  DUMP
-----
1           2  Typ=1 Len=2: 195,150

```

Этот оператор INSERT выполнен успешно, и мы можем видеть, что длина (возвращаемая функцией LENGTH) вставленного значения составляет 1 символ — все строковые функции работают *в отношении символов*. Поэтому длина поля равна 1 символ, но функция LENGTHB (длина в байтах) показывает, что для хранения данного символа используются 2 байта, а функция DUMP в точности отражает то, что это за байты. Таким образом, пример демонстрирует одну очень распространенную проблему, с которой сталкиваются во время применения многобайтных символьных наборов, а именно — поле VARCHAR2 (N) не обязательно содержит N *символов*, но может содержать N *байтов*.

Еще одна проблема, с которой часто имеют дело, заключается в том, что максимальная длина в байтах типа VARCHAR2 составляет 4000, а CHAR — 2000:

```

EODA@ORA12CR1> declare
2     l_data varchar2(4000 char);
3     l_ch   varchar2(1 char) := unistr( '\00d6' );
4 begin
5     l_data := rpad( l_ch, 4000, l_ch );
6     insert into t ( c ) values ( l_data );
7 end;
8 /
declare
*
ERROR at line 1:
ORA-01461: can bind a LONG value only for insert into a LONG column
ORA-06512: at line 6
ОШИБКА в строке 1:
ORA-01461: значение LONG можно привязать только для вставки в столбец LONG
ORA-06512: в строке 6

```

Это говорит о том, что 4000-символьная строка, которая в действительности имеет длину 8000 байтов, не может быть сохранена в поле VARCHAR2 (4000 CHAR) на постоянной основе. Она уместается в переменную PL/SQL, потому что в PL/SQL типу VARCHAR2 разрешено достигать размера 32 Кбайт. Однако когда такая строка сохраняется в таблице, вступает в силу жесткое ограничение в 4000 байтов. Мы можем успешно сохранить 2000 этих символов:

```

EODA@ORA12CR1> declare
2     l_data varchar2(4000 char);
3     l_ch   varchar2(1 char) := unistr( '\00d6' );
4 begin
5     l_data := rpad( l_ch, 2000, l_ch );
6     insert into t ( c ) values ( l_data );
7 end;
8 /

```

PL/SQL procedure successfully completed.

Процедура PL/SQL успешно завершена.

```
EODA@ORA12CR1> select length( c ), lengthb( c )
      2   from t
      3   where c is not null;
```

| LENGTH(C) | LENGTHB(C) |
|-----------|------------|
| ----- | ----- |
| 2000 | 4000 |

Как видите, они занимают в хранилище 4000 байтов.

Варианты “N”

Давайте для полноты изложения посмотрим, когда используются типы NVARCHAR2 и NCHAR. Они применяются в системах, где возникает потребность в управлении и хранении информации с использованием многочисленных символьных наборов. Обычно такое происходит в базе данных, в которой преобладающим набором символов является однобайтный с фиксированной шириной (такой как WE8ISO8859P1), но имеется необходимость в поддержке и хранении определенного объема многобайтных данных. Существует немало систем, которые работают с унаследованными данными, но для ряда новых приложений нуждаются в поддержке многобайтных данных. Вдобавок есть системы, которым требуется эффективность однобайтного набора символов для большинства операций (операции над строками, состоящими из символов фиксированной ширины, более эффективны, чем над строками, где каждый символ может занимать разное количество байтов), но в определенных случаях необходима гибкость многобайтных данных.

Типы данных NVARCHAR2 и NCHAR обеспечивают такие потребности. В целом они ведут себя так же, как их эквиваленты VARCHAR2 и CHAR, но со следующими исключениями.

- Их текст сохраняется и управляется в национальном наборе символов базы данных, а не в стандартном наборе символов.
- Их длины всегда предоставлены в символах, тогда как для CHAR/VARCHAR2 длина может быть указана либо в байтах, либо в символах.

В Oracle9i и последующих версиях национальный набор символов базы данных может принимать одно из двух значений: UTF8 или AL16UTF16 (UTF16 в Oracle9i; AL16UTF16 в Oracle 10g). Это делает типы NCHAR и NVARCHAR подходящими для хранения только многобайтных данных, что является изменением по сравнению с более ранними выпусками базы данных (в Oracle8i и предшествующих версиях можно было выбирать любой набор символов в качестве национального).

Двоичные строки: типы RAW

Наряду с текстом база данных Oracle поддерживает хранение двоичных данных. Двоичные данные не подвергаются преобразованиям символьных наборов, которые обсуждались ранее в отношении типов CHAR и VARCHAR2. Следовательно, двоичные типы данных не пригодны для хранения текста, введенного пользователем, но подходят для хранения зашифрованной информации. Зашифрованные данные — это не “текст”, а двоичное представление исходного текста, документы текстового про-

цессора, содержащие двоичную информацию разметки, и т.д. Любые строки байтов, которые не рассматриваются базой данных как “текст” (или любой другой базовый тип данных, такой как число, дата и т.д.) и не участвуют в преобразовании символьных наборов, должны сохраняться в двоичном типе данных.

В Oracle поддерживаются три типа для хранения двоичных данных.

- Тип RAW, на котором будет сосредоточено внимание в этом разделе, подходит для хранения данных RAW размером до 2000 байтов. Начиная с версии Oracle 12c, тип RAW можно сконфигурировать для хранения вплоть до 32 767 байтов информации.
- Тип BLOB, поддерживающий двоичные данные намного большего размера, о котором пойдет речь в разделе “Типы LOB” далее в главе.
- Тип LONG RAW, который поддерживается в целях обратной совместимости и не должен применяться в новых приложениях.

Синтаксис, используемый для двоичного типа RAW, прост:

```
RAW ( <размер> )
```

Например, следующий код создает таблицу, способную хранить 16 байтов двоичной информации на строку:

```
EODA@ORA12CR1> create table t ( raw_data raw(16) );
Table created.
Таблица создана.
```

В отношении хранения на диске тип RAW во многом похож на тип VARCHAR2. Тип RAW — это двоичная строка переменной длины, т.е., например, только что созданная таблица T может где угодно хранить от 0 до 16 байтов двоичных данных. Данные ничем не дополняются пробелами, как в случае типа CHAR.

При работе с данными RAW вы, скорее всего, обнаружите, что они неявно преобразуются в тип VARCHAR2 — т.е. многие инструменты, такие как SQL*Plus, не будут отображать данные RAW напрямую, а преобразовывать их перед отображением в шестнадцатеричный формат. В приведенном далее примере мы поместим в таблицу T двоичные данные с применением SYS_GUID() — встроенной функции, которая возвращает 16-байтную строку RAW, являющуюся глобально уникальной (GUID означает *globally unique identifier* — глобально уникальный идентификатор):

```
EODA@ORA12CR1> insert into t values ( sys_guid() );
1 row created.
1 строка создана.

EODA@ORA12CR1> select * from t;

RAW_DATA
-----
EEF18AA30B563AF0E043B7D04F0A4A30
```

Здесь сразу же можно отметить два момента. Во-первых, данные RAW выглядят подобно символьной строке. Именно так SQL*Plus извлекает и выводит их; на диске они хранятся по-другому. Инструмент SQL*Plus не должен выводить на экран произвольные данные, т.к. это может вызвать серьезные побочные эффекты на дисплее. Помните, что двоичные данные могут включать управляющие символы вроде

возврата каретки или перевода строки либо, возможно, символ <Ctrl+G>, который приведет к выдаче терминалом звукового сигнала.

Во-вторых, данные RAW выглядят намного длиннее 16 байтов — фактически, в этом примере можно насчитать 32 символа. Это объясняется тем фактом, что каждый двоичный байт требует двух шестнадцатеричных символов для отображения (если головным символом является ноль, то он не выводится). Сохраненные данные RAW действительно имеют длину 16 байтов, и в этом можно удостовериться с помощью функции DUMP. Ниже получен дамп значения двоичной строки с использованием дополнительного параметра для указания основания системы счисления, которая должна применяться при отображении значения каждого байта. Здесь используется основание 16, что позволяет сравнить результат дампа с предыдущей строкой:

```
EODA@ORA12CR1> select dump(raw_data,16) from t;
```

```
DUMP(RAW_DATA,16)
```

```
-----
Typ=23 Len=16: ee,f1,8a,a3,b,56,3a,f0,e0,43,b7,d0,4f,a,4a,30
```

Итак, функция DUMP показывает, что двоичная строка на самом деле имеет длину 16 байтов (Len=16) и отображает двоичные данные байт за байтом. Легко заметить, что дамп соответствует результатам неявного преобразования, которое выполняет SQL*Plus при извлечении данных RAW в строку. Такое неявное преобразование производится также и в обратном направлении:

```
EODA@ORA12CR1> insert into t values ( 'abcdef' );
```

```
1 row created.
```

```
1 строка создана.
```

Этот оператор вставляет не строку abcdef, а значение RAW, состоящее из трех байтов AB, CD, EF, или 171, 205, 239 в десятичном виде. Если вы попытаетесь указать строку, которая не состоит из допустимых шестнадцатеричных цифр, то получите сообщение об ошибке:

```
EODA@ORA12CR1> insert into t values ( 'abcdefgh' );
```

```
insert into t values ( 'abcdefgh' )
```

```
*
```

```
ERROR at line 1:
```

```
ORA-01465: invalid hex number
```

```
ОШИБКА в строке 1:
```

```
ORA-01465: недопустимое шестнадцатеричное число
```

Тип RAW можно индексировать и применять в предикатах — он так же функционален, как любой другой тип данных. Тем не менее, вы должны постараться избегать нежелательных неявных преобразований, памятуя о том, что преобразования будут происходить.

Я предпочитаю и рекомендую использовать во всех случаях *явные* преобразования, которые можно реализовать с применением перечисленных ниже встроенных функций:

- HEXTORAW — для преобразования строк шестнадцатеричных символов в тип RAW;
- RAWTONEH — для преобразования строк типа RAW в шестнадцатеричные строки.

Функция RAWTOHEX неявно вызывается инструментом SQL*Plus, когда он извлекает тип RAW в строку, а функция HEXTORAW неявно вызывается при вставке строки. Избегать неявных преобразований и всегда кодировать явно — хорошая практика. Следовательно, предшествующие примеры можно было бы переписать так:

```
EODA@ORA12CR1> select rawtohex(raw_data) from t;
RAWTOHEX (RAW_DATA)
-----
EEF18AA30B563AF0E043B7D04F0A4A30

EODA@ORA12CR1> insert into t values ( hextoraw('abcdef') );
1 row created.
1 строка создана.
```

Расширенные типы данных

До выхода версии Oracle 12c максимальной длиной, разрешенной для типов данных VARCHAR2 и NVARCHAR2, была 4000 байтов, а для типа данных RAW — 2000 байтов. Начиная с Oracle 12c, эти типы данных могут быть сконфигурированы для хранения вплоть до 32 767 байтов. Ниже приведены шаги по включению расширенных типов данных для неконтейнерной базы данных с единственным экземпляром (типы баз данных рассматривались в главе 2). Эти шаги должны выполняться от имени пользователя SYS:

```
SYS@012CE> shutdown immediate;
SYS@012CE> startup upgrade;
SYS@012CE> alter system set max_string_size=extended;
SYS@012CE> @?/rdbms/admin/utl132k.sql
SYS@012CE> shutdown immediate;
SYS@012CE> startup;
```

На заметку! В руководстве Oracle Database Reference (Справочник по СУБД Oracle) можно найти подробные сведения по реализации расширенных типов данных для всех видов баз данных (с единственным экземпляром, контейнерной, RAC и логической резервной базы данных Data Guard).

После изменения параметра MAX_STRING_SIZE на EXTENDED вы не сможете восстановить его стандартное значение (STANDARD). Такое изменение является однонаправленным. Если вы нуждаетесь в переключении обратно, придется выполнить восстановление в точку, перед которой было сделано это изменение — т.е. вам понадобятся резервные копии RMAN (полученные до внесения изменений) или включенная ретроспективная база данных. Вы также можете выполнить посредством Data Pump экспорт базы данных с расширенными типами данных и затем импортировать результаты в базу данных, в которой расширенные типы данных не включены, но с одним предостережением — любые таблицы с расширенными столбцами импортированы не будут.

После включения расширенных типов данных можно создать таблицу с расширенными столбцами:

```
EODA@012CE> create table t(et varchar2(32727)) tablespace users;
Table created.
Таблица создана.
```

Вот как выглядит описание таблицы:

```
EODA@012CE> desc t
Name                                     Null?      Type
-----
ET                                     VARCHA2 (32727)
```

Расширенным столбцом VARCHAR2 можно манипулировать через SQL, как если бы он был обычным столбцом, например:

```
EODA@012CE> insert into t values(rpad('abc',10000,'abc'));
EODA@012CE> select substr(et,9500,10) from t where UPPER(et) like 'ABC%';
```

Расширенный тип данных внутренне реализован как LOB-объект. Предполагая, что таблица T создана в схеме, не содержащей какие-то другие объекты, в результате запроса USER_OBJECTS вы получите следующий вывод:

```
EODA@012CE> select object_name, object_type from user_objects;
OBJECT_NAME                                OBJECT_TYPE
-----
SYS_LOB0000019479C00001$$                LOB
SYS_IL0000019479C00001$$                INDEX
T                                           TABLE
```

Запросив USER_LOBS, можно просмотреть дополнительные сведения о сегменте LOB:

```
EODA@012CE> select table_name, column_name, segment_name,
  ↳tablespace_name, in_row
  2    from user_lobs where table_name='T';
TABLE_NAME  COLUMN_NAME  SEGMENT_NAME                                TABLESPACE_NAME  IN_
-----
T           ET           SYS_LOB0000019479C00001$$                USERS              YES
```

Вы не имеете прямого контроля над LOB-объектом, который ассоциирован с расширенным столбцом. Это означает, что вы не можете манипулировать лежащим в основе столбцом LOB с помощью пакета DBMS_LOB. Кроме того, внутренний LOB-объект, связанный со столбцом расширенного типа, не будет виден через представление DBA_TAB_COLUMNS или COL\$.

Сегмент LOB и ассоциированный с ним индекс LOB всегда хранятся в табличном пространстве той таблицы, в которой был создан расширенный тип данных. Следуя обычным правилам хранения LOB, первые 4000 байтов сохраняются внутри таблицы как встроенные. Любые данные, выходящие за пределы 4000 байтов, хранятся в сегменте LOB. Если табличное пространство, в котором создан LOB-объект, использует метод ASSM, то LOB создается как SECUREFILE, а иначе как BASICFILE.

На заметку! В разделе "Типы LOB" далее в главе обсуждается хранение строк, а также технические аспекты парадигм хранения SECUREFILE и BASICFILE.

Доступ посредством SQL к любым данным, хранящимся в сегменте LOB расширенного столбца, прозрачным образом обрабатывается Oracle. Из этого вытекают некоторые интересные последствия.

Например, вы можете успешно производить выборку данных, которые хранятся в расширенном столбце, через связь базы данных. Следующий фрагмент кода осуществляет выборку (через связь базы данных) из таблицы по имени T в удаленной базе данных O12CE:

```
EODA@ORA12CR1> select substr(et, 9000,10) from t@O12CE;
SUBSTR(ET,9000,10)
-----
cabcabcab
```

Почему это так важно? Рассмотрим, что происходит, когда в удаленной базе данных O12CE создается таблица со столбцом, который определен с типом данных LOB:

```
EODA@O12CE> create table c(ct clob);
Table created.
Таблица создана.
```

База данных Oracle сгенерирует ошибку, если вы попытаетесь произвести выборку из столбца LOB удаленно посредством связи базы данных:

```
EODA@ORA12CR1> select * from c@O12CE;
ERROR:
ORA-22992: cannot use LOB locators selected from remote tables
ОШИБКА:
ORA-22992: не удастся использовать локаторы LOB, выбранные из удаленных таблиц
```

В отношении расширенных столбцов можно также выполнять операции сравнения над множествами (UNION, UNION ALL, MINUS, INTERSECT), например:

```
EODA@O12CE> select et from t minus select et from t;
```

Несмотря на то что вы попытались сравнить два столбца LOB через операции над множествами, Oracle возвращает сообщение об ошибке:

```
EODA@O12CE> select ct from c minus select ct from c;
select ct from c minus select ct from c
*
ERROR at line 1:
ORA-00932: inconsistent datatypes: expected - got CLOB
ОШИБКА в строке 1:
ORA-00932: несовместимые типы данных: ожидался -, получен CLOB
```

Предшествующие примеры продемонстрировали, что вы располагаете большей гибкостью при работе с расширенным типом данных, чем в случае, если бы пришлось взаимодействовать со столбцом LOB напрямую. Следовательно, если приложение имеет дело с символьными данными, объем которых превышает 4000 байтов, но меньше или равен 32 727 байтам, то имеет смысл подумать о применении расширенных типов данных. К тому же, если вы переходите из среды базы данных, отличной от Oracle (которая поддерживает крупные символьные столбцы), в среду Oracle, то средство расширенных типов данных поможет провести этот переход намного легче, т.к. теперь вы можете определять большие размеры для столбцов VARCHAR2, NVARCHAR2 и RAW внутри Oracle.

Числовые типы

В Oracle 10g и последующих версиях поддерживаются три собственных типа данных, которые подходят для хранения чисел. В Oracle9i Release 2 и предшествующих версиях доступен только один собственный тип данных, предназначенный для хранения числовых данных. В приведенном ниже списке тип NUMBER поддерживается всеми версиями, а другие два типа присутствуют только в Oracle 10g и последующих версиях.

- **NUMBER.** Тип NUMBER в Oracle способен хранить числовые данные с исключительно высокой степенью точности — фактически, 38 цифр. Лежащий в основе формат данных похож на упакованное десятичное представление. Тип NUMBER имеет формат переменной длины — от 0 до 22 байтов. Он подходит для хранения любого числа, начиная с $10e^{-130}$ и заканчивая $10e^{126}$, не включая его. В наши дни это наиболее широко используемый числовой тип.
- **BINARY_FLOAT.** Это собственный, утвержденный IEEE, числовой тип с плавающей точкой одинарной точности. Для его хранения на диске требуется 5 байтов: четыре фиксированных байта для числа с плавающей точкой и один байт для длины. Он способен хранить числа в диапазоне $\pm 1038,53$ с 6 цифрами точности.
- **BINARY_DOUBLE.** Это собственный, утвержденный IEEE, числовой тип с плавающей точкой двойной точности. Для его хранения на диске требуется 9 байтов: восемь фиксированных байтов для числа с плавающей точкой и один байт для длины. Может хранить числа в диапазоне $\pm 10308,25$ с 13 цифрами точности.

В этом кратком обзоре видно, что тип NUMBER в Oracle имеет значительно более высокую точность, чем типы BINARY_FLOAT и BINARY_DOUBLE, но намного меньший диапазон значений, чем BINARY_DOUBLE. То есть с помощью типа NUMBER можно хранить числа очень точно со многими значащими цифрами, но типы BINARY_FLOAT и BINARY_DOUBLE позволяют хранить намного меньшие и большие числа. Для примера создадим таблицу с разнообразными типами данных в ней и посмотрим, что в ней сохраняется для одних и тех же входных данных:

```
EODA@ORA12CR1> create table t
2 ( num_col number,
3   float_col binary_float,
4   dbl_col   binary_double
5 )
6 /
```

Table created.

Таблица создана.

```
EODA@ORA12CR1> insert into t ( num_col, float_col, dbl_col )
2 values ( 1234567890.0987654321,
3         1234567890.0987654321,
4         1234567890.0987654321 );
```

1 row created.

1 строка создана.

```
EODA@ORA12CR1> set numformat 9999999999.9999999999
EODA@ORA12CR1> select * from t;
```

| NUM_COL | FLOAT_COL | DBL_COL |
|------------------------|------------------------|------------------------|
| 1234567890.09876543210 | 1234567940.00000000000 | 1234567890.09876540000 |

Обратите внимание, что `NUM_COL` возвращает точно то же число, которое мы предоставили в качестве входных данных. Входное число содержит менее 38 значащих цифр (на самом деле их 20), поэтому сохранилось точное число. Однако представить точно это число в столбце `FLOAT_COL` типа `BINARY_FLOAT` невозможно. Фактически в нем точно сохранились только 7 цифр. Столбец `DBL_COL` показывает намного лучший результат, точно представляя число в данном случае вплоть до 17 цифр. Тем не менее, в целом это должно быть четким указанием на то, что типы `BINARY_FLOAT` и `BINARY_DOUBLE` не подойдут для финансовых приложений! Поупражнявшись с разными значениями, вы получите разные результаты:

```
EODA@ORA12CR1> delete from t;
1 row deleted.
1 строка удалена.

EODA@ORA12CR1> insert into t ( num_col, float_col, dbl_col )
  2 values ( 9999999999.9999999999,
  3          9999999999.9999999999,
  4          9999999999.9999999999 );
1 row created.
1 строка создана.

EODA@ORA12CR1> select * from t;
```

| NUM_COL | FLOAT_COL | DBL_COL |
|-----------------------|-------------------------|-------------------------|
| 9999999999.9999999990 | 100000000000.0000000000 | 100000000000.0000000000 |

И снова столбец NUM_COL точно представил число, но столбцы FLOAT_COL и DBL_COL — нет. Это вовсе не означает, что тип NUMBER способен обеспечить бесконечную точность — просто он обладает намного большей точностью. Результаты, аналогичные типам BINARY_FLOAT и BINARY_DOUBLE, довольно легко получить и с типом NUMBER:

[illegible][illegible]

Как видите, когда мы складывает вместе очень большое (123×10^{20}) и очень малое (123×10^{-20}) числа, то теряем точность, потому что такая арифметика требует более 38 цифр. Само большое число может быть представлено точно, равно как и малое число, но результат сложения малого и большого чисел — нет. В том, что это не является проблемой, связанной с отображением/форматированием, можно удостовериться следующим образом:

```
EODA@ORA12CR1> select num_col from t where num_col = 123*1e20;
```

```

                                NUM_COL
-----
123000000000000000000000000000.0000000000000000000000000000000000
```

Величина NUM_COL равна 123×10^{20} , а не значению, которое мы пытались вставить.

Синтаксис и использование типа NUMBER

Синтаксис типа NUMBER достаточно прост:

```
NUMBER( p, s )
```

Здесь P и S являются необязательными и применяются для указания описанных ниже характеристик.

- P — точность (precision), или общее количество цифр. По умолчанию точность составляет 38, а ее допустимые значения входят в диапазон от 1 до 38. Для представления точности 38 можно также использовать символ *.
- S — масштаб (scale), или количество цифр справа от десятичной точки. Допустимыми значениями являются от -84 до 127, а стандартное значение зависит от того, указана ли точность. Если точность не задана, то масштаб по умолчанию принимает максимальное значение. Если точность указана, то стандартное значение масштаба равно 0 (цифры справа от десятичной точки отсутствуют). Таким образом, например, столбец, определенный с типом NUMBER, хранит числа с плавающей точкой (с цифрами после десятичной точки), а с типом NUMBER(38) — только целые числа (без цифр после десятичной точки), т.к. во втором случае масштабом по умолчанию является 0.

Вы должны воспринимать точность и масштаб как средства *редактирования* данных — в некотором смысле инструменты для обеспечения целостности данных. Точность и масштаб вообще *не влияют* на то, как данные хранятся на диске, а только на то, какие значения допускаются и каким образом числа округляются. Например, если значение превышает разрешенную точность, то Oracle сообщает об ошибке:

```
EODA@ORA12CR1> create table t ( num_col number(5,0) );
```

```
Table created.
```

```
Таблица создана.
```

```
EODA@ORA12CR1> insert into t (num_col) values ( 12345 );
```

```
1 row created.
```

```
1 строка создана.
```

```
EODA@ORA12CR1> insert into t (num_col) values ( 123456 );
```

```
insert into t (num_col) values ( 123456 )
```

```
*
```

```
ERROR at line 1:
```

```
ORA-01438: value larger than specified precision allowed for this column
```

ОШИБКА в строке 1:

ORA-01438: значение больше допускаемого точностью, которая указана для этого столбца

Итак, вы можете применять точность для обеспечения некоторых ограничений целостности данных. В этом случае столбцу NUM_COL не разрешено иметь более пяти цифр.

С другой стороны, масштаб используется для управления округлением числа, например:

```
EODA@ORA12CR1> create table t ( msg varchar2(10), num_col number(5,2) );
```

Table created.

Таблица создана.

```
EODA@ORA12CR1> insert into t (msg,num_col) values ( '123.45', 123.45 );
```

1 row created.

1 строка создана.

```
EODA@ORA12CR1> insert into t (msg,num_col) values ( '123.456', 123.456 );
```

1 row created.

1 строка создана.

```
EODA@ORA12CR1> select * from t;
```

| MSG | NUM_COL |
|---------|---------|
| 123.45 | 123.45 |
| 123.456 | 123.46 |

Обратите внимание, что число 123.456, имеющее более пяти цифр, на этот раз обработано успешно. Причина в том, что применяемый в этом примере масштаб округлил 123.456 до двух цифр, давая в результате 123.46, *после* чего 123.46 было проверено на предмет точности, проверка прошла, и число было вставлено в таблицу. Однако если мы попытаемся выполнить следующую вставку, это закончится неудачей, т.к. число 1234.00 содержит всего более пяти цифр:

```
EODA@ORA12CR1> insert into t (msg,num_col) values ( '1234', 1234 );
```

```
insert into t (msg,num_col) values ( '1234', 1234 )
```

*

ERROR at line 1:

ORA-01438: value larger than specified precision allowed for this column

ОШИБКА в строке 1:

ORA-01438: значение больше допускаемого точностью, которая указана для этого столбца

Когда указан масштаб 2, допускается максимум три цифры слева от десятичной точки и две — справа от нее. Поэтому предложенное число не подходит. Столбец NUMBER(5,2) способен содержать все значения между 999.99 и -999.99.

Может показаться странным, что масштабу разрешено варьироваться от -84 до 127. Для какой цели мог бы понадобиться отрицательный масштаб? Он позволяет округлять значения слева от десятичной точки. Подобно тому, как NUMBER(5,2) округляет значения до ближайшего .01, NUMBER(5,-2) будет округлять до ближайшего 100. Например:

```
EODA@ORA12CR1> create table t ( msg varchar2(10), num_col number(5,-2) );
```

Table created.

Таблица создана.

```
EODA@ORA12CR1> insert into t (msg,num_col) values ( '123.45', 123.45 );
1 row created.
```

1 строка создана.

```
EODA@ORA12CR1> insert into t (msg,num_col) values ( '123.456', 123.456 );
1 row created.
```

1 строка создана.

```
EODA@ORA12CR1> select * from t;
```

| MSG | NUM_COL |
|---------|---------|
| 123.45 | 100 |
| 123.456 | 100 |

Числа были округлены до ближайшей сотни. Мы по-прежнему имеем пять цифр точности, но теперь слева от десятичной точки разрешены семь цифр (включая два хвостовых нуля):

```
EODA@ORA12CR1> insert into t (msg,num_col) values ( '1234567', 1234567 );
1 row created.
```

1 строка создана.

```
EODA@ORA12CR1> select * from t;
```

| MSG | NUM_COL |
|---------|---------|
| 123.45 | 100 |
| 123.456 | 100 |
| 1234567 | 1234600 |

```
EODA@ORA12CR1> insert into t (msg,num_col) values ( '12345678', 12345678 );
insert into t (msg,num_col) values ( '12345678', 12345678 )
```

*

```
ERROR at line 1:
```

```
ORA-01438: value larger than specified precision allowed for this column
```

ОШИБКА в строке 1:

ORA-01438: значение больше допускаемого точностью, которая указана для этого столбца

Таким образом, точность диктует количество цифр, которые допускаются в числе после округления, используя масштаб для определения того, как производить округление. Точность — это ограничение целостности, а масштаб — средство редактирования.

Интересно и полезно отметить, что тип NUMBER фактически является типом данных переменной длины на диске, и он будет потреблять от 0 до 22 байтов для хранения. Очень часто программисты воспринимают числовые типы данных как типы фиксированной длины, поскольку обычно им приходится иметь дело с 2- или 4-байтными целыми числами и 4- или 8-байтными числами с плавающей точкой. Тип NUMBER в Oracle подобен символьной строке переменной длины. Мы можем посмотреть, что происходит с числами, которые содержат разное количество значащих цифр. Для этого создадим таблицу с двумя столбцами NUMBER и заполним первый столбец множеством чисел, имеющих 2, 4, 6, ..., 28 значащих цифр. Затем просто добавим 1 к каждому из них:

Последний факт объясняет, почему полезно знать, что числа хранятся в полях переменной ширины. При попытке определить размер таблицы (т.е. выяснить, сколько места в хранилище понадобится для 1 000 000 строк таблицы), вам следует внимательно учитывать поля NUMBER. Сколько ваши числа займут — 2 байта или 20 байтов? Каков будет средний размер? Все это существенно затрудняет точное определение размера таблицы, если нет репрезентативного тестового набора данных. Можно подсчитать размеры для худшего и лучшего случаев, но реальный размер окажется где-то между ними.

Синтаксис и использование типов BINARY_FLOAT/BINARY_DOUBLE

В версии Oracle 10g были введены два числовых типа для хранения данных; в предшествующих версиях они отсутствовали. Это соответствующие стандарту IEEE числа с плавающей точкой, с которыми постоянно имели дело многие программисты. С описанием этих типов и их реализации можно ознакомиться по ссылке http://ru.wikipedia.org/wiki/Число_с_плавающей_запятой. В определении числа с плавающей точкой интересно отметить следующий момент.

Число с плавающей точкой — это цифровое представление числа из определенного подмножества рациональных чисел, которое часто применяется в компьютере для приближения к произвольному вещественному числу. В частности, оно представляет целое или число с фиксированной точкой (значащая часть, или, неформально, мантисса), умноженное на основание системы счисления (в компьютерах обычно 2) в определенной целой степени (порядок). Когда основанием системы счисления является 2, это двоичный аналог экспоненциального представления (с основанием 10).

Они используются для приближения к числам и не настолько точны, как описанный ранее встроенный тип NUMBER. Числа с плавающей точкой обычно применяются в научных приложениях, но полезны также в приложениях многих других видов благодаря тому факту, что позволяют производить арифметические вычисления с помощью оборудования (центрального процессора), а не подпрограмм Oracle. По этой причине арифметические вычисления выполняются намного быстрее, если научное приложение оперирует вещественными числами, но вы вряд ли захотите использовать числа с плавающей точкой для хранения финансовой информации. Например, предположим, что необходимо просуммировать числа с плавающей точкой 0.3 и 0.1. Может показаться, что результат должен быть равен 0.4. Однако в арифметике с плавающей точкой это не так — результат будет чуть больше, чем 0.4:

```
EODA@ORA12CR1> select to_char( 0.3f + 0.1f, '0.99999999999999' ) from dual;
TO_CHAR(0.3F+0.1F
-----
0.400000000600000
```

Это не ошибка, а способ работы чисел с плавающей точкой IEEE. Такая манера функционирования пригодна для решения определенного класса задач, но только не тех, в которых участвуют денежные величины!

Синтаксис объявления столбцов этого типа в таблице совершенно прямолинеен:

BINARY_FLOAT
BINARY_DOUBLE

Вот и все. Для указанных типов никаких опций не предусмотрено.

Несобственные числовые типы

В дополнение к типам NUMBER, BINARY_FLOAT и BINARY_DOUBLE база данных Oracle синтаксически поддерживает описанные ниже числовые типы данных.

- NUMERIC (p, s). Отображается в точности на NUMBER (p, s). Если p не указано, по умолчанию принимается 38.
- DECIMAL (p, s) или DEC (p, s). Отображается в точности на NUMBER (p, s). Если p не указано, по умолчанию принимается 38.
- INTEGER или INT. Отображается в точности на тип NUMBER (38).
- SMALLINT. Отображается в точности на тип NUMBER (38).
- FLOAT (p). Отображается на тип NUMBER.
- DOUBLE PRECISION. Отображается на тип NUMBER.
- REAL. Отображается на тип NUMBER.

На заметку! Когда я применяю формулировку “синтаксически поддерживает”, то имею в виду, что оператор CREATE может использовать эти типы данных, но “за кулисами” они на самом деле представляют собой тип NUMBER. В Oracle 10g Release 1 и последующих версиях есть всего три собственных числовых формата, а в Oracle9i Release 2 и предшествующих версиях — только один такой формат. Любой другой числовой тип данных всегда отображается на собственный тип NUMBER, поддерживаемый Oracle.

Соображения по поводу производительности

Вообще говоря, тип NUMBER в Oracle является наилучшим выбором для большинства приложений. Тем не менее, с этим типом связаны последствия в отношении производительности. Тип NUMBER — это *программный тип данных*, т.е. он реализован в самом программном обеспечении Oracle. Мы не можем применять встроенные аппаратные операции для сложения двух значений NUMBER, т.к. тип NUMBER эмулируется программным обеспечением. Однако типы с плавающей точкой не располагают такой реализацией. Для выполнения операции сложения двух чисел с плавающей точкой Oracle будет использовать оборудование.

Сказанное несложно проверить. Создадим таблицу, содержащую около 70 000 строк, и поместим в ее столбцы типов NUMBER и BINARY_FLOAT/BINARY_DOUBLE одни и те же данные:

```
EODA@ORA12CR1> create table t
2 ( num_type    number,
3   float_type  binary_float,
4   double_type binary_double
5 )
6 /
```

Table created.

Таблица создана.


```

EODA@ORA12CR1> insert /*+ APPEND */ into t
2  select rownum, rownum, rownum
3  from all_objects
4  /
72089 rows created.
72089 строк создано.

EODA@ORA12CR1> commit;
Commit complete.
Фиксация завершена.

```

Затем мы запустим в отношении каждого столбца один и тот же запрос, в котором применяется сложная математическая функция, такая как LN (натуральный логарифм). В результате мы будем наблюдать совершенно разную утилизацию центрального процессора:

```

select sum(ln(num_type)) from t
call      count      cpu      elapsed
-----
total      4          4.45      4.66

select sum(ln(float_type)) from t
call      count      cpu      elapsed
-----
total      4          0.07      0.08

select sum(ln(double_type)) from t
call      count      cpu      elapsed
-----
total      4          0.06      0.06

```

Тип NUMBER в этом примере использует процессор почти в 63 раза интенсивнее, чем типы с плавающей точкой. Но вы должны помнить, что мы не получаем в точности один и тот же ответ из всех трех запросов!

```

EODA@ORA12CR1> set numformat 999999.999999999999999999
EODA@ORA12CR1> select sum(ln(num_type)) from t;
              SUM(LN(NUM_TYPE))
-----
734280.3209126472927309

EODA@ORA12CR1> select sum(ln(double_type)) from t;
              SUM(LN(DOUBLE_TYPE))
-----
734280.3209126447300000

```

Числа с плавающей точкой были приближениями целевого числа с точностью между 6 и 13 цифрами. Ответ, полученный от столбца типа NUMBER, намного точнее, чем от столбцов, имеющих типы с плавающей точкой. Тем не менее, при проведении глубинного анализа данных (Data Mining) или сложного числового анализа научных данных такая потеря точности обычно приемлема, а выигрыш производительности может оказаться весьма существенным.

На заметку! Если вы интересуетесь подробностями арифметики с плавающей точкой и связанной с ней потерей точности, почитайте документ по ссылке http://docs.sun.com/source/806-3568/ncg_goldberg.html.

Необходимо отметить, что в этом случае можно есть мед, да еще и ложкой. С помощью встроенной функции CAST можно выполнить преобразование типа NUMBER в тип с плавающей точкой на лету — до того, как с ним будут выполнены сложные вычисления. Это позволит значительно сократить загрузку процессора, приблизив ее к той, которой требуют собственные типы с плавающей точкой:

```
select sum(ln(cast( num_type as binary_double ) )) from t
call      count      cpu      elapsed
-----
total      4          0.08      0.08
```

Это означает, что мы можем хранить данные в очень точном виде, а при возникновении потребности в низкоуровневой скорости, когда типы с плавающей точкой превосходят тип NUMBER, для достижения этой цели использовать функцию CAST.

Типы LONG

Типы LONG в Oracle предлагаются в двух видах.

- Текстовый тип LONG, способный хранить 2 Гбайт текста. Текст, сохраненный в поле типа LONG, подвергается преобразованиям символьных наборов, что очень похоже на типы VARCHAR2 и CHAR.
- Тип LONG RAW, способный хранить до 2 Гбайт низкоуровневых двоичных данных (данных, к которым не применяются преобразования символьных наборов).

История типов LONG восходит к версии Oracle 6, где они ограничивались 64 Кбайт данных. В версии Oracle 7 они были расширены для поддержки вплоть до 2 Гбайт информации, но в версии Oracle 8 их заменили типами LOB, которые мы вскоре обсудим.

Вместо того чтобы демонстрировать использование типа LONG я объясню, почему применять LONG (или LONG RAW) в приложениях лишено смысла. Прежде всего, в документации Oracle очень четко выражено отношение к типам LONG. В руководстве *Oracle Database SQL Language Reference* заявлено следующее.

Не создавайте таблицу со столбцами LONG. Взамен используйте столбцы LOB (CLOB, NCLOB, BLOB). Столбцы LONG поддерживаются только в целях обратной совместимости.

Ограничения типов LONG и LONG RAW

С типами LONG и LONG RAW связаны ограничения, приведенные в табл. 12.2. Хотя, возможно, я забегаю вперед, но здесь указано, подчиняются ли типы LOB, являющиеся заменой типов LONG и LONG RAW, тем же самым ограничениям.

Таблица 12.2. Типы LONG в сравнении с LOB

| Тип LONG/LONG RAW | Тип CLOB/BLOB |
|---|--|
| В таблице может быть только один столбец типа LONG или LONG RAW | В таблице допускается иметь до 1000 столбцов типа CLOB или BLOB |
| Определяемые пользователем типы не могут иметь атрибутов типа LONG или LONG RAW | Определяемые пользователем типы могут в полной мере задействовать типы CLOB и BLOB |
| На столбцы типа LONG нельзя ссылаться в конструкции WHERE | На столбцы LOB можно ссылаться в конструкции WHERE, а для манипуляций ими предусмотрено семейство функций в пакете DBMS_LOB |
| Типы LONG не поддерживают распределенные транзакции | Типы LOB поддерживают распределенные транзакции |
| Типы LONG не могут реплицироваться с применением базовой или расширенной репликации | Типы LOB полностью поддерживают репликацию |
| Столбцы типа LONG не могут находиться в конструкциях GROUP BY, ORDER BY и CONNECT BY или в запросах, которые используют DISTINCT, UNIQUE, INTERSECT, MINUS либо UNION | Столбцы типа LOB могут участвовать во всех перечисленных конструкциях при условии, что к ним применяется функция, которая преобразует их в скалярные типы SQL (содержащие атомарные значения), такие как VARCHAR2, NUMBER или DATE |
| Функции/процедуры PL/SQL не могут принимать аргументы типа LONG | Функции/процедуры PL/SQL полноценно работают с типами LOB |
| Встроенные функции SQL (например, SUBSTR) не могут применяться к столбцам типа LONG | Встроенные функции SQL могут использоваться со столбцами типов LOB |
| Тип LONG нельзя применять в операторе CREATE TABLE AS SELECT | Типы LOB допускаются в операторе CREATE TABLE AS SELECT |
| Оператор ALTER TABLE MOVE нельзя использовать для таблиц, содержащих в себе типы LONG | Таблицы со столбцами типов LOB можно перемещать |

Как видите, в табл. 12.2 представлен довольно длинный список; существует много того, что при наличии столбца LONG в таблице делать попросту невозможно. При разработке новых приложений даже не рассматривайте возможность применения типа LONG. Вместо этого используйте подходящий тип LOB. Что касается существующих приложений: если вы сталкиваетесь с любым ограничением из числа указанных в табл. 12.2, то стоит серьезно подумать о преобразовании типа LONG в соответствующий тип LOB. Были предприняты необходимые меры для обеспечения обратной совместимости, так что приложение, написанное для типов LONG, будет прозрачным образом работать с типами LOB.

На заметку! Само собой разумеется, перед переводом производственной системы с типов LONG на типы LOB вы должны провести полное тестирование функциональности своих приложений.

Копирование с участием унаследованных типов LONG

Часто возникает вопрос: а как насчет словаря данных Oracle? Он изобилует столбцами LONG, и это затрудняет работу со столбцами словаря. Например, невозмож-

но применять SQL для выполнения поиска в словарном представлении ALL_VIEWS, чтобы найти все представления, которые содержат текст HELLO:

```
EODA@ORA12CR1> select *
  2 from all_views
  3 where text like '%HELLO%';
where text like '%HELLO%'
      *
```

ERROR at line 3:
ORA-00932: inconsistent datatypes: expected CHAR got LONG
ОШИБКА в строке 3:
ORA-00932: несовместимые типы данных: ожидался CHAR, получен LONG

Эта проблема не ограничивается представлением ALL_VIEWS; она затрагивает многие другие представления:

```
EODA@ORA12CR1> select table_name, column_name
  2 from dba_tab_columns
  3 where data_type in ( 'LONG', 'LONG RAW' )
  4 and owner = 'SYS'
  5 and table_name like 'DBA%'
  6 order by table_name;
```

| TABLE_NAME | COLUMN_NAME |
|------------------------------|-------------------|
| DBA_ADVISOR_SQLPLANS | OTHER |
| DBA_ARGUMENTS | DEFAULT_VALUE |
| DBA_CLUSTER_HASH_EXPRESSIONS | HASH_EXPRESSION |
| DBA_CONSTRAINTS | SEARCH_CONDITION |
| DBA_IND_EXPRESSIONS | COLUMN_EXPRESSION |
| DBA_IND_PARTITIONS | HIGH_VALUE |
| DBA_IND_SUBPARTITIONS | HIGH_VALUE |
| DBA_MVIEWS | QUERY |
| DBA_MVIEW_AGGREGATES | MEASURE |
| DBA_MVIEW_ANALYSIS | QUERY |
| DBA_NESTED_TABLE_COLS | DATA_DEFAULT |
| DBA_OUTLINES | SQL_TEXT |
| DBA_REGISTERED_MVIEWS | QUERY_TXT |
| DBA_REGISTERED_SNAPSHOTS | QUERY_TXT |
| DBA_SNAPSHOTS | QUERY |
| DBA_SQLSET_PLANS | OTHER |
| DBA_SQLTUNE_PLANS | OTHER |
| DBA_SUBPARTITION_TEMPLATES | HIGH_BOUND |
| DBA_SUMMARIES | QUERY |
| DBA_SUMMARY_AGGREGATES | MEASURE |
| DBA_TAB_COLS | DATA_DEFAULT |
| DBA_TAB_COLS_V\$ | DATA_DEFAULT |
| DBA_TAB_COLUMNS | DATA_DEFAULT |
| DBA_TAB_PARTITIONS | HIGH_VALUE |
| DBA_TAB_SUBPARTITIONS | HIGH_VALUE |
| DBA_TRIGGERS | TRIGGER_BODY |
| DBA_VIEWS | TEXT |
| DBA_VIEWS_AE | TEXT |
| DBA_ZONEMAPS | QUERY |
| DBA_ZONEMAP_MEASURES | MEASURE |

30 rows selected.

30 строк выбрано.

Итак, каково же решение? Если вы хотите использовать эти столбцы в SQL, их понадобится преобразовать в дружелюбный к SQL тип. Для этого можно применить функцию, определенную пользователем. В следующем примере показано, как достигнуть этой цели с помощью функции SUBSTR_OF, которая позволит эффективно преобразовывать любые 4000 байтов типа LONG в VARCHAR2 для использования в SQL. Запрос выглядит следующим образом:

```
EODA@ORA12CR1> select *
2   from (
3   select owner, view_name,
4           long_help.substr_of( 'select text
5                               from dba_views
6                               where owner = :owner
7                               and view_name = :view_name',
8                               1, 4000,
9                               'owner', owner,
10                              'view_name', view_name ) substr_of_view_text
11   from dba_views
12   where owner = user
13   )
14   where upper(substr_of_view_text) like '%INNER%'
15 /
```

Здесь выполнено преобразование первых 4000 байтов столбца TEXT из типа LONG в тип VARCHAR2, и теперь результат можно применять в предикате. С помощью такого же приема можно было бы также реализовать собственные функции INSTR, LIKE и им подобные для типа LONG. В этой книге демонстрируется получение подстроки из типа LONG.

Пакет, который мы реализуем, имеет показанную ниже спецификацию:

```
EODA@ORA12CR1> create or replace package long_help
2   authid current_user
3   as
4       function substr_of
5       ( p_query in varchar2,
6         p_from in number,
7         p_for in number,
8         p_name1 in varchar2 default NULL,
9         p_bind1 in varchar2 default NULL,
10        p_name2 in varchar2 default NULL,
11        p_bind2 in varchar2 default NULL,
12        p_name3 in varchar2 default NULL,
13        p_bind3 in varchar2 default NULL,
14        p_name4 in varchar2 default NULL,
15        p_bind4 in varchar2 default NULL )
16   return varchar2;
17 end;
18 /
Package created.
Пакет создан.
```

Обратите внимание, что в строке 2 указано AUTHID CURRENT_USER. Это позволяет выполнять пакет от имени вызывающего, со всеми предоставленными ему ро-

лями и привилегиями, и важно по двум причинам. Во-первых, мы предпочитаем не нарушать безопасность базы данных — пакет будет возвращать только подстроки столбцов, которые разрешено видеть вызывающему. В частности, пакет неуязвим к атакам внедрением SQL — он запускается от имени вызывающего, а не владельца пакета. Во-вторых, мы хотим установить этот пакет в базе данных один раз и сделать его функциональность доступной для использования всеми; это становится возможным благодаря применению прав вызывающего. В случае использования стандартной модели безопасности PL/SQL пакет запускался бы с правами владельца и видел бы только данные, разрешенные владельцу, которые вполне могли не включать набор данных, доступных вызывающему.

Концепция, положенная в основу функции SUBSTR_OF, предусматривает применение запроса, который извлекает максимум одну строку и один столбец: интересующее нас значение LONG. Функция SUBSTR_OF при необходимости разберет запрос, привяжет к нему любые входные параметры и программно извлечет результаты, возвратив нужную часть значения LONG.

Тело пакета, т.е. его реализация, начинается с определения двух глобальных переменных. Переменная G_CURSOR содержит постоянный курсор, открытый на протяжении сеанса. Это позволит избежать необходимости в многократном открытии и закрытии курсора, а также в разборе SQL больше, чем необходимо. Вторая глобальная переменная, G_QUERY, служит для запоминания текста последнего SQL-запроса, разобранный в пакете. До тех пор, пока запрос остается неизменным, он разбирается только один раз. Таким образом, даже в случае запрашивания 5000 строк, если переданный функции запрос не изменился, то понадобится только однократный разбор:

```
EODA@ORA12CR1> create or replace package body long_help
2 as
3
4     g_cursor number := dbms_sql.open_cursor;
5     g_query  varchar2(32765);
6
```

Далее в пакете находится закрытая процедура BIND_VARIABLE, которая будет использоваться для привязки входных параметров, переданных вызывающим. Мы реализуем отдельную закрытую процедуру только ради облегчения. Привязка должна производиться, когда имя параметра является NOT NULL. Вместо того чтобы четыре раза выполнять проверку для каждого входного параметра, мы делаем это лишь один раз внутри процедуры:

```
7 procedure bind_variable( p_name in varchar2, p_value in varchar2 )
8 is
9 begin
10     if ( p_name is not null )
11     then
12         dbms_sql.bind_variable( g_cursor, p_name, p_value );
13     end if;
14 end;
15
```

Затем в теле пакета следует действительная реализация SUBSTR_OF. Функция начинается со своего объявления, указанного в спецификации пакета, и объявления нескольких локальных переменных. Переменная L_BUFFER будет применяться для

возврата значения, а `L_BUFFER_LEN` — для хранения длины, которая возвращается функцией, предоставляемой Oracle:

```

16
17 function substr_of
18 ( p_query in varchar2,
19   p_from   in number,
20   p_for    in number,
21   p_name1  in varchar2 default NULL,
22   p_bind1  in varchar2 default NULL,
23   p_name2  in varchar2 default NULL,
24   p_bind2  in varchar2 default NULL,
25   p_name3  in varchar2 default NULL,
26   p_bind3  in varchar2 default NULL,
27   p_name4  in varchar2 default NULL,
28   p_bind4  in varchar2 default NULL )
29 return varchar2
30 as
31     l_buffer      varchar2(4000);
32     l_buffer_len  number;
33 begin

```

Первое, что делает код — проверяет корректность `P_FROM` и `P_FOR`. Значение `P_FROM` должно быть числом, которое больше или равно 1, а значение `P_FOR` должно находиться между 1 и 4000 — в точности как для встроенной функции `SUBSTR`:

```

34     if ( nvl(p_from,0) <= 0 )
35     then
36         raise_application_error
37         (-20002, 'From must be >= 1 (positive numbers)' );
38     end if;
39     if ( nvl(p_for,0) not between 1 and 4000 )
40     then
41         raise_application_error
42         (-20003, 'For must be between 1 and 4000' );
43     end if;
44

```

Далее производится проверка, получен ли новый запрос, который нуждается в разборе. Если последний разобранный запрос совпадает с текущим, то этот шаг можно пропустить. Очень важно отметить, что в строке 47 производится проверка, что в переданном значении `P_QUERY` действительно содержится `SELECT` — мы будем использовать этот пакет *только* для выполнения SQL-операторов `SELECT`. Вот как реализована эта проверка:

```

45     if ( p_query <> g_query or g_query is NULL )
46     then
47         if ( upper(trim(nvl(p_query,'x'))) not like 'SELECT%' )
48         then
49             raise_application_error
50             (-20001, 'This must be a select only' );
51         end if;
52         dbms_sql.parse( g_cursor, p_query, dbms_sql.native );
53         g_query := p_query;
54     end if;

```

Итак, все готово к привязке входных параметров к этому запросу. Любые переданные имена, не равные NULL, будут привязаны к запросу, так что при выполнении он находит корректные строки:

```

55     bind_variable( p_name1, p_bind1 );
56     bind_variable( p_name2, p_bind2 );
57     bind_variable( p_name3, p_bind3 );
58     bind_variable( p_name4, p_bind4 );
59

```

Наконец, мы можем запустить запрос и извлечь строку. С помощью DBMS_SQL.COLUMN_VALUE_LONG мы извлекаем необходимую подстроку из LONG и возвращаем ее:

```

60     dbms_sql.define_column_long(g_cursor, 1);
61     if (dbms_sql.execute_and_fetch(g_cursor)>0)
62     then
63         dbms_sql.column_value_long
64         (g_cursor, 1, p_for, p_from-1,
65          l_buffer, l_buffer_len );
66     end if;
67     return l_buffer;
68 end substr_of;
69
70 end;
71 /

```

Package body created.

Тело пакета создано.

Теперь у вас появилась возможность применения этого пакета с *любым* унаследованным столбцом LONG в базе данных, что позволит выполнять многие операции внутри конструкции WHERE, которые ранее были недоступными. Например, вы можете найти все секции в схеме, содержащие в HIGH_VALUE год 2014 (помните, что если таблицы с HIGH_VALUE, равным 2014, отсутствуют, то не следует ожидать возвращения чего-либо):

```

EODA@ORA12CR1> select *
2   from (
3   select table_owner, table_name, partition_name,
4          long_help.substr_of
5          ( 'select high_value
6            from all_tab_partitions
7            where table_owner = :o
8              and table_name = :n
9              and partition_name = :p',
10          1, 4000,
11          'o', table_owner,
12          'n', table_name,
13          'p', partition_name ) high_value
14   from all_tab_partitions
15   where table_owner = user
16   )
17  where high_value like '%2014%'
18 /

```


| TABLE_OWNER | TABLE_NAME | PARTITION_NAME | HIGH_VALUE |
|-------------|------------|----------------|------------|
| ----- | ----- | ----- | ----- |
| EODA | F_CONFIGS | CONFIG_P_7 | 20140101 |

Используя тот же самый прием, т.е. обработку внутри функции результатов запроса, который возвращает единственную строку с одним столбцом LONG, при необходимости можно реализовать собственные конструкции INSERT, LIKE и т.д.

Такая реализация хорошо работает с типом LONG, но не с типами LONG RAW. Типы LONG RAW не являются доступными повсеместно (в пакете DBMS_SQL нет функции вроде COLUMN_VALUE_LONG_RAW). К счастью, это не слишком серьезное ограничение, поскольку столбцы типов LONG RAW в словаре не задействованы, и необходимость извлечения “подстроки” из LONG RAW, чтобы можно было выполнять в них поиск, возникает редко. Однако если в этом все же есть потребность, то вы не будете применять PL/SQL для столбцов LONG RAW, превышающих 32 Кбайт, т.к. в PL/SQL отсутствуют методы работы с типами LONG RAW, которые имеют объем более 32 Кбайт. В этом случае должен использоваться Java, C, C++, Visual Basic или какой-то другой язык.

Другой подход предусматривает временное преобразование LONG или LONG RAW в CLOB или BLOB с применением встроенной функции TO_LOB и глобальной временной таблицы. Процедура PL/SQL могла бы выглядеть следующим образом:

```
Insert into global_temp_table ( blob_column )
select to_lob(long_raw_column) from t where...
```

Такой прием был бы удобным в приложении, которое иногда нуждается в работе с одиночным значением LONG RAW. Тем не менее, поступать так постоянно нежелательно из-за большого объема сопутствующих работ. Если вы обнаруживаете потребность в частом обращении к этому приему, то определенно должны преобразовать тип LONG RAW в BLOB и в дальнейшем работать с ним.

Типы DATE, TIMESTAMP и INTERVAL

Собственные типы данных DATE, TIMESTAMP и INTERVAL в Oracle тесно взаимосвязаны. Типы DATE и TIMESTAMP хранят фиксированные значения даты/времени с разной степенью точности. Тип INTERVAL используется для хранения периода времени, такого как “8 часов” или “30 дней”. Результатом вычитания двух отметок времени может быть интервал, а сложение интервала в 8 часов и TIMESTAMP дает в результате новое значение TIMESTAMP, которое отражает время на 8 часов позднее.

Тип данных DATE присутствовал во многих выпусках Oracle — когда я еще только начинал работать с этой базой данных, т.е., по крайней мере, в версии 5 или даже раньше. Типы TIMESTAMP и INTERVAL являются относительно новыми, если сравнивать их с остальными, т.к. они были введены в выпуске Oracle9i Release 1. По этой простой причине вы обнаружите, что тип DATE оказался наиболее распространенным типом данных, выбираемым для хранения информации о дате/времени. Однако во многих новых приложениях тип TIMESTAMP применяется по двум причинам: он поддерживает дробные части секунды (а тип DATE нет) и располагает поддержкой часовых поясов (что также отсутствует в DATE).

Мы рассмотрим каждый из этих типов после обсуждения форматов DATE/TIMESTAMP и случаев их использования.

Форматы

Я не собираюсь здесь раскрывать все форматы DATE, TIMESTAMP и INTERVAL. Они подробно описаны в руководстве *Oracle Database SQL Language Reference*, которое бесплатно доступно всем желающим. В нашем распоряжении имеется огромный объем форматов, поэтому жизненно важно понимать то, что они собой представляют. Я настоятельно рекомендую исследовать их.

Я хотел бы обсудить, что делают форматы, из-за многочисленных недоразумений, связанных с этой темой. Форматы применяются для решения двух задач:

- форматирование данных по пути из базы в желаемом стиле;
- сообщение базе данных о том, как преобразовывать входную строку в значение типа DATE, TIMESTAMP или INTERVAL.

На этом все. Распространенное заблуждение, с которым я сталкиваюсь на протяжении многих лет, заключается в том, что используемый формат каким-то образом влияет на то, что сохраняется на диске, и как в действительности сохраняются данные. *Формат никак не воздействует на то, каким образом сохраняются данные. Формат применяется только для преобразования единственного двоичного формата, используемого для хранения DATE, в строку или для преобразования строки в этот единственный двоичный формат хранения DATE.* То же самое справедливо для типов TIMESTAMP и INTERVAL.

Мой совет относительно форматов прост: применяйте их. Используйте их, когда отправляете в базу данных строку, которая представляет значение DATE, TIMESTAMP или INTERVAL. *Не* полагайтесь на стандартные форматы для дат — принятые по умолчанию правила когда-нибудь в будущем могут измениться.

На заметку! В главе 1 можно найти действительно серьезное основание, связанное с безопасностью, чтобы никогда не применять функции TO_CHAR и TO_DATE без явного указания формата. Там был описан пример атаки внедрением SQL, которая стала доступной конечному пользователю просто потому, что разработчик забыл использовать явный формат. Кроме того, выполнение операций с данными без применения явного формата даты может и будет приводить к получению некорректных ответов. Чтобы уловить это: скажите, какую дату представляет строка '01-02-03'? Что бы вы ни ответили, я замечу, что вы ошибаетесь. Никогда не полагайтесь на правила по умолчанию!

Если вы рассчитывали на стандартный формат для даты, а он изменился, это может негативно повлиять на приложение. Если преобразовать дату не удастся, может возникнуть ошибка, передаваемая конечному пользователю, образоваться серьезная брешь в защите или, что еще хуже, в таблицу могут молча вставиться *некорректные данные*. Рассмотрим следующий оператор INSERT, полагающийся на стандартную маску даты:

```
Insert into t ( date_column ) values ( '01/02/03' );
```

Предположим, что приложение рассчитывало на стандартную маску даты DD/MM/YY (день/месяц/год). Показанный оператор даст 1 февраля 2003 года (в случае запуска запроса после 2000 года; мы вскоре вернемся к этому обстоятельству).

Теперь представим, что кто-то решил, что правильным форматом даты должен быть ММ/ДД/YY (месяц/день/год). Тогда предыдущая дата превратится в 2 января 2003 года. А если кто-то решит установить в качестве стандартного формата YY/ММ/ДД (год/месяц/день), то получится 3 февраля 2001 года. Короче говоря, в отсутствие формата даты, сопровождающего эту строку, существует много способов ее интерпретации. Оператор INSERT должен выглядеть, например, так:

```
Insert into t ( date_column ) values ( to_date( '01/02/03', 'DD/MM/YY' ) );
```

А если вы поинтересуетесь моим мнением, то лучше будет поступить следующим образом:

```
Insert into t ( date_column ) values ( to_date( '01/02/2003', 'DD/MM/YYYY' ) );
```

То есть необходимо указывать год в виде четырех символов. Несколько лет назад наша индустрия на собственном горьком опыте убедилась, насколько много времени и сил пришлось потратить, чтобы исправить программное обеспечение, в котором разработчики пытались “сэкономить” 2 байта. Похоже, что с течением времени этот урок благополучно забыли. Я считаю непростительным в наши дни *не* использовать четырехсимвольное представление года! То, что 2000 год канул в небытие, вовсе не означает, что теперь опять можно применять двухзначный год. Подумайте, к примеру, о датах рождения. Если вы введете дату рождения так, как показано ниже, то какой будет эта дата — 1 февраля 2010 года или 1 февраля 1910 года?

```
Insert into t ( DOB ) values ( to_date( '01/02/10', 'DD/MM/YY' ) );
```

Любая из двух дат является допустимой; уверенно выбрать какую-то одну не удастся.

То же самое касается данных, покидающих базу. Если вы запускаете запрос SELECT DATE_COLUMN FROM T и извлекаете столбец в строку внутри своего приложения, то должны применить к ней явный формат даты. Каким бы ни был формат, ожидаемый приложением, его следует указывать в нем явно. В противном случае, когда кто-то в будущем изменит стандартный формат даты, работа вашего приложения может быть нарушена.

Теперь давайте рассмотрим более подробно сами типы данных.

Тип DATE

DATE — это тип данных фиксированной ширины в 7 байтов, предназначенный для хранения даты/времени. Он всегда содержит семь атрибутов — век, год внутри века, месяц, день месяца, часы, минуты и секунды. Для представления этой информации Oracle использует внутренний формат, поэтому на самом деле 25 июня 2005 года, 12:01:00 хранится вовсе не как 20, 05, 06, 25, 12, 01, 00. С помощью встроенной функции DUMP можно посмотреть, что в действительности сохраняет Oracle:

```
EOGA@ORA12CR1> create table t ( x date );
Table created.
Таблица создана.

EOGA@ORA12CR1> insert into t (x) values
  2 ( to_date( '25-jun-2005 12:01:00',
  3 'dd-mon-yyyy hh24:mi:ss' ) );
1 row created.
1 строка создана.
```

```

EODA@ORA12CR1> select x, dump(x,10) d from t;
X          D
-----
25-JUN-05  Typ=12 Len=7: 120,105,6,25,13,2,1

```

Байты века и года (120,105 в выводе DUMP) хранятся в нотации с избытком 100. Для получения корректного века и года из них понадобится вычесть 100. Причина применения такой нотации — поддержка дат до нашей эры (BC) и нашей эры (AD). Если в результате вычитания 100 из байта века получено отрицательное число, то это дата до нашей эры. Например:

```

EODA@ORA12CR1> insert into t (x) values
  2 ( to_date( '01-jan-4712bc',
  3           'dd-mon-yyyybc hh24:mi:ss' ) );
1 row created.
1 строка создана.

EODA@ORA12CR1> select x, dump(x,10) d from t;
X          D
-----
25-JUN-05  Typ=12 Len=7: 120,105,6,25,13,2,1
01-JAN-12  Typ=12 Len=7: 53,88,1,1,1,1,1

```

Таким образом, когда мы вставляем дату 01-JAN-4712BC, байт века получает значение 53, а $53 - 100 = -47$, т.е. век, который вставлен. Поскольку он отрицательный, мы знаем, что речь идет о дате до нашей эры. Этот формат хранения также позволяет естественным образом сортировать даты в двоичном виде. Поскольку 4712 год до нашей эры *меньше чем* 4710 год до нашей эры, мы хотели бы иметь двоичное представление, которое поддерживает это. Получив дампы для указанных двух дат, мы можем видеть, что 01-JAN-4710BC *больше чем* тот же день в 4712 году до нашей эры, так что эти значения эффективно сортируются и сравниваются:

```

EODA@ORA12CR1> insert into t (x) values
  2 ( to_date( '01-jan-4710bc',
  3           'dd-mon-yyyybc hh24:mi:ss' ) );
1 row created.
1 строка создана.

EODA@ORA12CR1> select x, dump(x,10) d from t;
X          D
-----
25-JUN-05  Typ=12 Len=7: 120,105,6,25,13,2,1
01-JAN-12  Typ=12 Len=7: 53,88,1,1,1,1,1
01-JAN-10  Typ=12 Len=7: 53,90,1,1,1,1,1

```

Байты месяца и дня — следующие два поля — сохраняются естественным образом без каких-либо модификаций. Поэтому 25 июня представлено байтом месяца со значением 6 и байтом дня со значением 25. Поля часов, минут и секунд сохраняются с нотации с избытком 1, т.е. для получения реального значения нужно вычесть единицу из значения каждого компонента. Следовательно, полночь в поле даты представляется как 1,1,1.

Этот 7-байтный формат сортируется естественным образом, как вы уже видели — он является 7-байтным полем, которое может быть очень эффективно отсортировано в двоичной манере от меньшего к большему (и наоборот). Вдобавок его структура допускает любые усечения без преобразования в какой-то другой формат.

Например, усечение только что сохраненной даты (25-JUN-2005 12:01:00) до дня (отбрасывание часов, минут и секунд) делается очень просто. Для этого достаточно установить три завершающих байта в 1,1,1 и компонент времени исключается. Рассмотрим новую таблицу T со следующими вставками:

```
EOA@ORA12CR1> create table t ( what varchar2(10), x date );
Table created.
```

Таблица создана.

```
EOA@ORA12CR1> insert into t (what, x) values
2  ( 'orig',
3   to_date( '25-jun-2005 12:01:00',
4           'dd-mon-yyyy hh24:mi:ss' ) );
```

1 row created.

1 строка создана.

```
EOA@ORA12CR1> insert into t (what, x)
2  select 'minute', trunc(x,'mi') from t
3  union all
4  select 'day', trunc(x,'dd') from t
5  union all
6  select 'month', trunc(x,'mm') from t
7  union all
8  select 'year', trunc(x,'y') from t
9  /
```

4 rows created.

4 строки создано.

```
EOA@ORA12CR1> select what, x, dump(x,10) d from t;
```

| WHAT | X | D |
|--------|-----------|-----------------------------------|
| orig | 25-JUN-05 | Typ=12 Len=7: 120,105,6,25,13,2,1 |
| minute | 25-JUN-05 | Typ=12 Len=7: 120,105,6,25,13,2,1 |
| day | 25-JUN-05 | Typ=12 Len=7: 120,105,6,25,1,1,1 |
| month | 01-JUN-05 | Typ=12 Len=7: 120,105,6,1,1,1,1 |
| year | 01-JAN-05 | Typ=12 Len=7: 120,105,1,1,1,1,1 |

Чтобы усечь дату до года, базе данных потребуется лишь поместить единицы в последние пять байтов — очень быстрая операция. Теперь мы имеем сортируемое и сравниваемое значение поля DATE, которое усечено до уровня года, причем получаем его настолько эффективно, насколько это вообще возможно.

Добавление и вычитание времени из значения DATE

Мне часто задают вопрос о добавлении и вычитании значений времени из типа DATE. Например, каким образом добавить к значению DATE один день, восемь часов, один год, один месяц и т.д.? Существуют три приема, которые обычно используются для этого.

- Просто добавить NUMBER к DATE. Добавление 1 к DATE — это способ добавить 1 день. Таким образом, добавление 1/24 к DATE прибавляет 1 час и т.д.
- Для добавления единиц времени можно применять тип INTERVAL, как вскоре будет описано. Типы INTERVAL поддерживают два уровня детализации — годы и месяцы или дни/часы/минуты/секунды. То есть у вас может быть интервал из лет и месяцев либо интервал из дней, часов, минут и секунд.

- Добавить месяцы с использованием встроенной функции `ADD_MONTHS`. Поскольку добавление месяца обычно не сводится просто к прибавлению от 28 до 31 дня, для упрощения этой задачи была предусмотрена специальная функция.

В табл. 12.3 демонстрируются приемы, которые будут применяться для добавления (или вычитания) N единиц времени к дате.

Таблица 12.3. Добавление времени к дате

| Единица времени | Операция | Описание |
|-----------------|---|---|
| N секунд | <code>DATE + n/24/60/60</code> <code>DATE + n/86400</code> <code>DATE +</code> <code>NUMTODSINTERVAL(n, 'second')</code> | В сутках находится 86 400 секунд. Поскольку добавление 1 означает прибавление одного дня, добавление 1/86 400 добавляет к дате одну секунду. Я предпочитаю записывать <code>n/24/60/60</code> вместо <code>n/86400</code> . Эти варианты эквивалентны. Но еще более читабельный метод предполагает использование <code>NUMTODSINTERVAL</code> (преобразование числа в интервал в днях/секундах) для добавления N секунд |
| N минут | <code>DATE + n/24/60</code> <code>DATE + n/1440</code> <code>DATE +</code> <code>NUMTODSINTERVAL(n, 'minute')</code> | В сутках находится 1440 минут. Таким образом, добавление 1/1440 прибавляет одну минуту к <code>DATE</code> . Более читабельный метод заключается в применении функции <code>NUMTODSINTERVAL</code> |
| N часов | <code>DATE + n/24</code> <code>DATE +</code> <code>NUMTODSINTERVAL(n, 'hour')</code> | В сутках находится 24 часа. Поэтому добавление 1/24 прибавляет к <code>DATE</code> один час. Более читабельный метод предусматривает использование функции <code>NUMTODSINTERVAL</code> |
| N дней | <code>DATE + n</code> | Чтобы добавить (или вычесть) N дней, нужно просто добавить (или вычесть) N к значению <code>DATE</code> |
| N недель | <code>DATE + 7*n</code> | Неделя составляет 7 дней, поэтому для добавления или вычитания N недель нужно просто прибавить или вычесть $7*N$ |
| N месяцев | <code>ADD_MONTHS (DATE, n)</code> <code>DATE +</code> <code>NUMTOYMINTERVAL(n, 'month')</code> | Для прибавления N месяцев к <code>DATE</code> можно применять встроенную функцию <code>ADD_MONTHS</code> или же добавлять интервал из N месяцев. Ознакомьтесь с приведенными ниже важными предостережениями относительно использования интервалов в месяцах со значениями <code>DATE</code> |
| N лет | <code>ADD_MONTHS (DATE, 12*n)</code> <code>DATE +</code> <code>NUMTOYMINTERVAL(n, 'year')</code> | Чтобы прибавить или вычесть N лет, можно применять встроенную функцию <code>ADD_MONTHS</code> с передачей ей $12*N$. Аналогичного результата можно достичь с интервалом в годах. Ознакомьтесь с приведенными ниже важными предостережениями относительно использования интервалов в годах со значениями <code>DATE</code> |

Ниже перечислены рекомендации по применению типа DATE в Oracle.

- Для добавления часов, минут и секунд используйте встроенную функцию NUMTODSINTERVAL.
- Для добавления дней просто добавляйте их количество.
- Для добавления месяцев и лет применяйте встроенную функцию ADD_MONTHS.

Я не рекомендую пользоваться функцией NUMTOYMINTERVAL (для добавления месяцев и лет). Причина связана с тем, как эта функция ведет себя в конце месяцев.

Функция ADD_MONTHS трактует конец дней месяца особым образом. В сущности, она округляет даты — если мы добавляем один месяц к месяцу, содержащему 31 день, а следующий месяц имеет меньше 31 дня, то ADD_MONTHS возвратит последний день следующего месяца. Кроме того, добавление одного месяца к последнему дню месяца дает в результате последний день следующего месяца. Это можно увидеть при добавлении одного месяца к месяцу с 30 или меньшим числом дней:

```
EOA@ORA12CR1> alter session set nls_date_format = 'dd-mon-yyyy hh24:mi:ss';
Session altered.
```

Сеанс изменен.

```
EOA@ORA12CR1> select dt, add_months(dt,1)
  2   from (select to_date('29-feb-2000','dd-mon-yyyy') dt from dual )
  3 /
DT                                ADD_MONTHS(DT,1)
-----
29-feb-2000 00:00:00  31-mar-2000 00:00:00
```

```
EOA@ORA12CR1> select dt, add_months(dt,1)
  2   from (select to_date('28-feb-2001','dd-mon-yyyy') dt from dual )
  3 /
DT                                ADD_MONTHS(DT,1)
-----
28-feb-2001 00:00:00  31-mar-2001 00:00:00
```

```
EOA@ORA12CR1> select dt, add_months(dt,1)
  2   from (select to_date('30-jan-2001','dd-mon-yyyy') dt from dual )
  3 /
DT                                ADD_MONTHS(DT,1)
-----
30-jan-2001 00:00:00  28-feb-2001 00:00:00
```

```
EOA@ORA12CR1> select dt, add_months(dt,1)
  2   from (select to_date('30-jan-2000','dd-mon-yyyy') dt from dual )
  3 /
DT                                ADD_MONTHS(DT,1)
-----
30-jan-2000 00:00:00  29-feb-2000 00:00:00
```

Заметили, что добавление одного месяца к 29 февраля 2000 года дает 31 марта 2000 года? 29 февраля — последний день месяца, поэтому функция ADD_MONTHS возвращает последний день следующего месяца. К тому же обратите внимание, что результатом добавления одного месяца к 30 января 2000 года и 30 января 2001 года будет последний день февраля в 2000 и 2001 годах соответственно.

Если мы сравним это с тем, как работает добавление интервалов, то получим совсем другие результаты:

```
EODA@ORA12CR1> select dt, dt+numtoyminterval(1,'month')
  2   from (select to_date('29-feb-2000','dd-mon-yyyy') dt from dual )
  3 /
DT                                DT+NUMTOYMINTERVAL(1
-----
29-feb-2000 00:00:00  29-mar-2000 00:00:00

EODA@ORA12CR1> select dt, dt+numtoyminterval(1,'month')
  2   from (select to_date('28-feb-2001','dd-mon-yyyy') dt from dual )
  3 /
DT                                DT+NUMTOYMINTERVAL(1
-----
28-feb-2001 00:00:00  28-mar-2001 00:00:00
```

Обратите внимание, что результирующая дата не является последним днем следующего месяца, а *тем же самым* днем следующего месяца. Сомнительно, чтобы такое поведение оказалось приемлемым, но посмотрим, что происходит, когда результирующий месяц не содержит такое количество дней:

```
EODA@ORA12CR1> select dt, dt+numtoyminterval(1,'month')
  2   from (select to_date('30-jan-2001','dd-mon-yyyy') dt from dual )
  3 /
select dt, dt+numtoyminterval(1,'month')
      *
ERROR at line 1:
ORA-01839: date not valid for month specified
ОШИБКА в строке 1:
ORA-01839: недопустимая дата для указанного месяца

EODA@ORA12CR1> select dt, dt+numtoyminterval(1,'month')
  2   from (select to_date('30-jan-2000','dd-mon-yyyy') dt from dual )
  3 /
select dt, dt+numtoyminterval(1,'month')
      *
ERROR at line 1:
ORA-01839: date not valid for month specified
ОШИБКА в строке 1:
ORA-01839: недопустимая дата для указанного месяца
```

Согласно моему опыту, это делает применение интервалов в месяцах при вычислениях с датами вообще невозможным. Похожая проблема возникает с интервалами в годах: добавление одного года к 29 февраля 2000 года приводит к возникновению ошибки во время выполнения, т.к. в 2001 году нет 29 февраля.

Получение разности между двумя значениями DATE

Другой часто задаваемый вопрос касается получения разности между двумя датами. Ответ предельно прост: нужно вычесть одну дату из другой. Результатом будет число, представляющее количество дней между двумя датами. Кроме того, доступна встроенная функция MONTHS_BETWEEN, которая возвращает число, представляющее количество месяцев (включая дробную часть) между двумя датами. И, наконец, типы данных INTERVAL дают еще один метод получения времени, истекшего между дву-

мя датами. В следующем SQL-запросе демонстрируется результат вычитания дат (с выводом количества дней между ними), использование функции MONTHS_BETWEEN и применение двух функций, работающих с типами INTERVAL:

```
EODA@ORA12CR1> select dt2-dt1 ,
2      months_between(dt2,dt1) months_btwn,
3      numtodsinterval(dt2-dt1,'day') days,
4      numtoyminterval(trunc(months_between(dt2,dt1)), 'month') months
5 from (select to_date('29-feb-2000 01:02:03','dd-mon-yyyy hh24:mi:ss') dt1,
6      to_date('15-mar-2001 11:22:33','dd-mon-yyyy hh24:mi:ss') dt2
7      from dual )
8 /
```

| DT2-DT1 | MONTHS_BTWN | DAYS | MONTHS |
|------------|-------------|-------------------------------|---------------|
| 380.430903 | 12.5622872 | +000000380 10:20:30.000000000 | +000000001-00 |

Все эти значения корректны, но не особенно пригодны для использования. В большинстве приложений желательно отображать годы, месяцы, дни, часы, минуты и секунды между датами. Достичь такой цели можно с применением комбинации показанных выше функций. Мы выберем два интервала: один для лет и месяцев, а другой — только для дней, часов и т.д. Мы будем использовать встроенную функцию MONTH_BETWEEN для определения целого числа месяцев между двумя датами, после чего с помощью встроенной функции NUMTOYMINTERVAL преобразуем это число в годы и месяцы. Вдобавок мы будем применять MONTHS_BETWEEN для вычитания целого числа месяцев между двумя датами из большей даты среди них, чтобы в итоге получить количество дней и часов между датами:

```
EODA@ORA12CR1> select numtoyminterval
2      (trunc(months_between(dt2,dt1)), 'month')
3      years_months,
4      numtodsinterval
5      (dt2-add_months( dt1, trunc(months_between(dt2,dt1)) ),
6      'day' )
7      days_hours
8 from (select to_date('29-feb-2000 01:02:03','dd-mon-yyyy hh24:mi:ss') dt1,
9      to_date('15-mar-2001 11:22:33','dd-mon-yyyy hh24:mi:ss') dt2
10     from dual )
11 /
```

| YEARS_MONTHS | DAYS_HOURS |
|---------------|-------------------------------|
| +000000001-00 | +000000015 10:20:30.000000000 |

Теперь совершенно ясно, что между двумя датами прошел 1 год, 15 дней, 10 часов, 20 минут и 30 секунд.

Вариации типа TIMESTAMP

Тип TIMESTAMP очень похож на DATE, но предлагает дополнительную поддержку дробной части секунд и часовых поясов. Мы рассмотрим тип TIMESTAMP в последующих трех разделах: в одном речь пойдет только о поддержке дробных частей секунд без поддержки часовых поясов, а в других — о двух методах сохранения TIMESTAMP с поддержкой часовых поясов.

Тип TIMESTAMP

Синтаксис базового типа данных TIMESTAMP прост:

TIMESTAMP (n)

Здесь N является необязательным параметром и служит для указания масштаба компонента секунд в TIMESTAMP; допустимое значение N лежит в диапазоне от 0 до 9. Если задано 0, это говорит о том, что значение TIMESTAMP будет функционально эквивалентным DATE и фактически хранить те же значения в той же самой манере:

```
EODA@ORA12CR1> create table t
```

```
2 ( dt date,
3 ts timestamp(0)
4 )
5 /
```

Table created.

Таблица создана.

```
EODA@ORA12CR1> insert into t values ( sysdate, systimestamp );
```

1 row created.

1 строка создана.

```
EODA@ORA12CR1> select dump(dt,10) dump, dump(ts,10) dump from t;
```

| DUMP | DUMP |
|-----------------------------------|------------------------------------|
| ----- | |
| Тип=12 Len=7: 120,110,4,12,20,4,8 | Тип=180 Len=7: 120,110,4,12,20,4,8 |

Типы данных отличаются (на это указывает поле Тип=), но стиль хранения данных тот же самый. В случае указания некоторой дробной части секунд тип данных TIMESTAMP будет иметь другую длину по сравнению с типом DATE, например:

```
EODA@ORA12CR1> create table t
```

```
2 ( dt date,
3 ts timestamp(9)
4 )
5 /
```

Table created.

Таблица создана.

```
EODA@ORA12CR1> insert into t values ( sysdate, systimestamp );
```

1 row created.

1 строка создана.

```
EODA@ORA12CR1> select dump(dt,10) dump, dump(ts,10) dump
```

```
2 from t;
```

| DUMP | DUMP |
|----------------------------------|---|
| ----- | |
| Тип=12 Len=7: 120,114,1,2,8,20,1 | Тип=180 Len=11: 120,114,1,2,8,20,1,53,55,172,40 |

Теперь TIMESTAMP потребляет 11 байтов хранилища, при этом 4 дополнительных байта в конце содержат дробные части секунды, в чем можно убедиться, взглянув на сохраненное значение времени:

```
EODA@ORA12CR1> alter session set nls_date_format = 'dd-mon-yyyy hh24:mi:ss';
```

Session altered.

Сеанс изменен.

```

EODA@ORA12CR1> select * from t;
DT                                TS
-----
02-jan-2014 07:19:00 02-JAN-14 07.19.00.892841000 AM
EODA@ORA12CR1> select dump(ts,16) dump from t;
DUMP
-----
Typ=180 Len=11: 78,72,1,2,8,14,1,35,37,ac,28
EODA@ORA12CR1> select to_number('3537ac28', 'xxxxxxxx' ) from dual;
TO_NUMBER('3537AC28','XXXXXXXX')
-----
892841000

```

Здесь можно видеть дробную часть секунд, сохраненную в последних 4 байтах. На этот раз мы использовали функцию DUMP, чтобы просмотреть данные в шестнадцатеричном виде, но эти 4 байта можно легко преобразовать в десятичное представление.

Сложение и вычитание времени с помощью **TIMESTAMP**

Те же самые приемы, которые применялись к DATE для выполнения арифметических действий с датами, работают с типом **TIMESTAMP**, но во многих случаях использования этих приемов значение **TIMESTAMP** преобразуется в тип **DATE**. Например:

```

EODA@ORA12CR1> alter session set nls_date_format = 'dd-mon-yyyy hh24:mi:ss';
Session altered.
Сеанс изменен.
EODA@ORA12CR1> select systimestamp ts, systimestamp+1 dt
2 from dual;
TS                                DT
-----
02-JAN-14 07.30.37.627678 AM -07:00 03-jan-2014 07:30:37

```

Обратите внимание, что добавление 1 фактически сдвигает **SYSTEMEASTAMP** на один день, но дробная часть секунд при этом исчезает, как будет и с информацией о часовом поясе. Поэтому применение здесь значений **INTERVAL** становится более важным:

```

EODA@ORA12CR1> select systimestamp ts, systimestamp
+numtodsinterval(1,'day') dt
2 from dual;
TS                                DT
-----
02-JAN-14 07.31.45.451317 AM -07:00 03-JAN-14 07.31.45.451317000 AM -07:00

```

Использование функции, возвращающей тип **INTERVAL**, предохраняет точность **TIMESTAMP**. Во время применения **TIMESTAMP** следует проявлять осторожность, чтобы избежать неявных преобразований. Но помните о предостережении относительно добавления интервалов в месяцах или годах к **TIMESTAMP**. Если результирующий день не является допустимой датой, то операция завершится неудачно (добавление одного месяца посредством **INTERVAL** к последнему дню января не достигнет успеха).

Получение разности между двумя *TIMESTAMP*

Здесь типы *DATE* и *TIMESTAMP* существенно отличаются. В то время как результат вычитания *DATE* из *DATE* имеет тип *NUMBER*, результатом той же операции над *TIMESTAMP* будет *INTERVAL*:

```
EODA@ORA12CR1> select dt2-dt1
2   from (select to_timestamp('29-feb-2000 01:02:03.122000',
3                        'dd-mon-yyyy hh24:mi:ss.ff') dt1,
4                        to_timestamp('15-mar-2001 11:22:33.000000',
5                        'dd-mon-yyyy hh24:mi:ss.ff') dt2
6   from dual )
7 /

DT2-DT1
-----
+0000000380 10:20:29.878000000
```

Разность между двумя значениями *TIMESTAMP* представляет собой значение *INTERVAL*, которое показывает количество дней и часов/минут/секунд между ними. Если нужно получить разность в годах, месяцах и т.п., придется снова использовать запрос, подобный тому, что применялся с датами:

```
EODA@ORA12CR1> select numtoyminterval
2   (trunc(months_between(dt2,dt1)), 'month')
3   years_months,
4   dt2-add_months(dt1,trunc(months_between(dt2,dt1)))
5   days_hours
6   from (select to_timestamp('29-feb-2000 01:02:03.122000',
7                        'dd-mon-yyyy hh24:mi:ss.ff') dt1,
8                        to_timestamp('15-mar-2001 11:22:33.000000',
9                        'dd-mon-yyyy hh24:mi:ss.ff') dt2
10  from dual )
11 /

YEARS_MONTHS  DAYS_HOURS
-----
+0000000001-00 +0000000015 10:20:30.000000000
```

Обратите внимание, что в этом случае из-за использования *ADD_MONTHS* значение *DT1* было неявно преобразовано в тип *DATE* и мы потеряли дробную часть секунд. Чтобы сохранить ее, необходимо предусмотреть дополнительный код. Мы могли бы применить *NUMTOYMINTERVAL*, чтобы добавить месяцы и предохранить *TIMESTAMP*, но тогда возникнет ошибка времени выполнения:

```
EODA@ORA12CR1> select numtoyminterval
2   (trunc(months_between(dt2,dt1)), 'month')
3   years_months,
4   dt2-(dt1 + numtoyminterval( trunc(months_between(dt2,dt1)), 'month' ))
5   days_hours
6   from (select to_timestamp('29-feb-2000 01:02:03.122000',
7                        'dd-mon-yyyy hh24:mi:ss.ff') dt1,
8                        to_timestamp('15-mar-2001 11:22:33.000000',
9                        'dd-mon-yyyy hh24:mi:ss.ff') dt2
10  from dual )
11 /
```

```
dt2-(dt1 + numtoyminterval( trunc(months_between(dt2,dt1)), 'month' ))
*
```

ERROR at line 4:

ORA-01839: date not valid for month specified

ОШИБКА в строке 4:

ORA-01839: недопустимая дата для указанного месяца

Лично я нахожу это неприемлемым. Фактически на момент отображения информации с годами и месяцами точность `TIMESTAMP` уже нарушена. Длительность года не постоянна (он может иметь 365 или 366 дней), и то же самое касается месяца. Если вы отображаете информацию в годах и месяцах, то потеря микросекунд несущественна; в таком случае вывода сведений вплоть до секунд более чем достаточно.

Тип `TIMESTAMP WITH TIME ZONE`

Тип `TIMESTAMP WITH TIME ZONE` наследует все качества типа `TIMESTAMP` и добавляет поддержку часового пояса. Он требует для хранения 13 байтов, с двумя дополнительными байтами, выделенными под сведения о часовом поясе.

Тип `TIMESTAMP WITH TIME ZONE` структурно отличается от `TIMESTAMP` только наличием этих двух байтов:

```
EODA@ORA12CR1> create table t
2 (
3   ts      timestamp,
4   ts_tz   timestamp with time zone
5 )
6 /
```

Table created.

Таблица создана.

```
EODA@ORA12CR1> insert into t ( ts, ts_tz )
2 values ( systimestamp, systimestamp );
1 row created.
1 строка создана.
```

```
EODA@ORA12CR1> select * from t;
```

| TS | TS_TZ |
|------------------------------|-------------------------------------|
| 02-JAN-14 03.02.51.890565 PM | 02-JAN-14 03.02.51.890565 PM -07:00 |

```
EODA@ORA12CR1> select dump(ts) dump, dump(ts_tz) dump from t;
```

DUMP

DUMP

Typ=180 Len=11: 120,114,1,2,16,3,52,53,20,241,136

Typ=181 Len=13: 120,114,1,2,23,3,52,53,20,241,136,13,60

При извлечении стандартный формат `TIMESTAMP` включал информацию о часовом поясе (когда выполнялся этот запрос, часовым поясом был U.S. Mountain Standard Time (стандартное горное время в США)).

Значения `TIMESTAMP WITH TIME ZONE` хранят данные в часовом поясе, который был указан при сохранении данных. Часовой пояс становится частью самих данных. Обратите внимание, что поле `TIMESTAMP WITH TIME ZONE` хранит ...23,3,52... для часов, минут и секунд (в нотации с избытком 1, т.е. на самом деле это 22:02:51),

тогда как поле `TIMESTAMP` хранит просто `...16,3,52...`, что соответствует 15:02:51 — точному времени вставки. Значение типа `TIMESTAMP WITH TIME ZONE` добавляет к нему семь часов, чтобы сохранить время в часовом поясе GMT (также известном как UTC). Два хвостовых байта используются при извлечении, чтобы соответствующим образом скорректировать значение `TIMESTAMP`.

В мои намерения здесь не входит раскрытие всех нюансов часовых поясов; данной теме посвящено много других материалов. Поэтому я только подчеркиваю, что в этом типе данных имеется поддержка часовых поясов. В наши дни поддержка подобного рода в приложениях становится более важной, чем когда-либо ранее. В далеком прошлом приложения даже близко не были настолько глобальными, как теперь. До широкого распространения Интернета приложения были намного более распределенными и децентрализованными, и часовые пояса неявно основывались на том, где находился сервер. Сегодня, когда крупные централизованные системы применяются людьми по всему миру, потребность в отслеживании и использовании часовых поясов является очень важной.

До того, как поддержка часовых поясов стала встроенной в тип данных, нужно было иметь функцию приложения, предназначенную для сохранения `DATE` в одном столбце и часового пояса в другом, а также функции для преобразования значений `DATE` из одного часового пояса в другой. Теперь это работа базы данных, и она может хранить данные во множестве часовых поясов:

```
EODA@ORA12CR1> create table t
  2 ( ts1 timestamp with time zone,
  3   ts2 timestamp with time zone
  4 )
  5 /
Table created.
Таблица создана.

EODA@ORA12CR1> insert into t (ts1, ts2)
  2 values ( timestamp'2014-02-27 16:02:32.212 US/Eastern',
  3          timestamp'2014-02-27 16:02:32.212 US/Pacific' );
1 row created.
1 строка создана.
```

Вдобавок база данных может выполнять с ними корректные арифметические действия для типа `TIMESTAMP`:

```
EODA@ORA12CR1> select ts1-ts2 from t;

TS1-TS2
-----
-0000000000 03:00:00.000000
```

Поскольку между этими двумя часовыми поясами существует разница в 3 часа, несмотря на то, что они показывают одно и то же время 16:02:32.212, интервал сообщает о наличии трехчасовой разницы. При выполнении арифметических действий для `TIMESTAMP` над типами `TIMESTAMP WITH TIME ZONE` оба типа сначала автоматически преобразуются Oracle в значение времени UTC, а затем происходит операция.

Тип `TIMESTAMP WITH LOCAL TIME ZONE`

Этот тип работает во многом подобно столбцу `TIMESTAMP`. Это 7- или 11-байтное поле (в зависимости от точности `TIMESTAMP`), но нормализованное для хранения с локальным часовым поясом базы данных. Чтобы увидеть это, мы применим команду `DUMP` еще раз. Для начала создадим таблицу с тремя столбцами типов `DATE`, `TIMESTAMP WITH TIME ZONE` и `TIMESTAMP WITH LOCAL TIME ZONE`, после чего вставим одно и то же значение во все три столбца:

```
EODA@ORA12CR1> create table t
  2 ( dt date,
  3   ts1 timestamp with time zone,
  4   ts2 timestamp with local time zone
  5 )
  6 /
Table created.
Таблица создана.

EODA@ORA12CR1> insert into t (dt, ts1, ts2)
  2 values ( timestamp'2014-02-27 16:02:32.212 US/Pacific',
  3          timestamp'2014-02-27 16:02:32.212 US/Pacific',
  4          timestamp'2014-02-27 16:02:32.212 US/Pacific' );
1 row created.
1 строка создана.

EODA@ORA12CR1> select dbtimezone from dual;

DBTIMEZONE
-----
-07:00
```

Теперь запросим дампы этих трех значений следующим образом:

```
EODA@ORA12CR1> select dump(dt), dump(ts1), dump(ts2) from t;

DUMP(DT)
-----
DUMP(TS1)
-----
DUMP(TS2)
-----
Typ=12 Len=7: 120,114,2,27,17,3,33
Typ=181 Len=13: 120,114,2,28,1,3,33,12,162,221,0,137,156
Typ=231 Len=11: 120,114,2,27,18,3,33,12,162,221,0
```

В этом случае мы видим, что были сохранены три совершенно разных представления даты/времени.

- **DT.** Этот столбец сохранил дату 27 февраля 2014 года и время 16:02:32. Часовой пояс и дробная часть секунд были утеряны, потому что использовался тип `DATE`. Никаких преобразований часовых поясов вообще не происходит. Мы сохранили точную дату/время вставки, но потеряли часовой пояс.
- **TS1.** Этот столбец сохранил информацию часового пояса и благодаря этому был нормализован ко времени UTC. Вставленное значение `TIMESTAMP` относилось к часовому поясу US/Pacific, которое отличалось от UTC на 8 часов. Таким образом, сохранилась дата 28 февраля 2014 года и время 00:02:32.

Это сдвинуло введенное значение времени вперед на 8 часов, чтобы сделать его временем UTC, и сохранило в последних двух байтах информацию о часовом поясе US/Pacific, так что позже эти данные могут быть должным образом интерпретированы.

- TS2. Предполагается, что этот столбец находится в часовом поясе базы данных, которым является US/Mountain. Теперь 16:02:32 US/Pacific соответствует 17:02:32 US/Mountain, и это было сохранено в байтах ...18, 3, 33... (нотация с избытком 1; не забудьте вычесть 1).

Поскольку столбец TS1 хранит исходный часовой пояс в последних двух байтах, в результате извлечения мы увидим следующий вывод:

```
EODA@ORA12CR1> select ts1, ts2 from t;
```

```
TS1
```

```
TS2
```

```
27-FEB-14 04.02.32.212000 PM US/PACIFIC
```

```
27-FEB-14 05.02.32.212000 PM
```

База данных должна быть в состоянии отобразить эту информацию, но столбец TS2 типа `TIMESTAMP WITH LOCAL TIME ZONE` показывает время в часовом поясе базы данных, который является подразумеваемым часовым поясом для этого столбца (и фактически для всех столбцов типа `TIMESTAMP WITH LOCAL TIME ZONE` в этой базе данных). Моя база данных находилась в часовом поясе US/Mountain, поэтому 16:02:32 US/Pacific теперь отображается как 17:02:32 в стандартном горном времени.

На заметку! Вы можете получить слегка отличающиеся результаты, если дата сохранялась, когда было актуальным стандартное время, и затем извлекалась, когда действовало летнее время (Daylight Savings). Тогда вывод в предыдущем примере отобразил бы разницу в два часа, а не один час, как могло бы показаться логичным. Я лишь хочу окончательно подчеркнуть тот факт, что математика часовых поясов намного сложнее, чем выглядит!

Если вы не нуждаетесь в запоминании исходного часового пояса, а необходим только тип данных, обеспечивающий согласованную мировую поддержку типов даты/времени, то тип `TIMESTAMP WITH LOCAL TIME ZONE` предлагает достаточную поддержку для большинства приложений. Кроме того, `TIMESTAMP(0) WITH LOCAL TIME ZONE` предоставляет эквивалент типа `DATE` с поддержкой часовых поясов; он потребляет 7 байтов хранилища и обладает возможностью хранения дат в нормализованной форме UTC.

Относительно типа `TIMESTAMP WITH LOCAL TIMEZONE` следует сделать одно предостережение: как только вы создадите таблицу с таким столбцом, то обнаружите, что часовой пояс базы данных заморожен — вы не в состоянии изменить его:

```
EODA@ORA12CR1> alter database set time_zone = 'PST';
```

```
alter database set time_zone = 'PST'
```

```
*
```

```
ERROR at line 1:
```

```
ORA-30079: cannot alter database timezone when database has  
TIMESTAMP WITH LOCAL TIME ZONE columns
```


ОШИБКА в строке 1:

ORA-30079: невозможно изменить часовой пояс базы данных, когда
в ней содержатся столбцы *TIMESTAMP WITH LOCAL TIME ZONE*

```
EODA@ORA12CR1> !oerr ora 30079
30079, 00000, "cannot alter database timezone when database has
      TIMESTAMP WITH LOCAL TIME ZONE columns"
// *Cause: An attempt was made to alter database timezone with
//      TIMESTAMP WITH LOCAL TIME ZONE column in the database.
// *Action: Either do not alter database timezone or first drop all the
//      TIMESTAMP WITH LOCAL TIME ZONE columns.
// *Причина: Предпринята попытка изменить часовой пояс базы данных
//      со столбцом TIMESTAMP WITH LOCAL TIME ZONE.
// *Действие: Либо не изменяйте часовой пояс базы данных, либо сначала
//      удалите все столбцы TIMESTAMP WITH LOCAL TIME ZONE.
```

Причина должна быть очевидной: если бы можно было изменить часовой пояс базы данных, то пришлось бы переписать каждую таблицу со столбцами типа *TIMESTAMP WITH LOCAL TIME ZONE*, т.к. их текущие значения станут неправильными с учетом нового часового пояса!

Тип INTERVAL

Вы уже видели тип *INTERVAL* — он применялся в предыдущем разделе. Этот тип обеспечивает способ представления продолжительности или интервала времени. В настоящем разделе мы обсудим два типа интервалов: *INTERVAL YEAR TO MONTH*, который способен хранить продолжительность, выраженную в годах и месяцах, и *INTERVAL DATE TO SECOND*, который предназначен для хранения продолжительности, заданной в днях, часах, минутах и секундах (включая дробные части секунды).

Прежде чем углубляться в специфику этих двух типов *INTERVAL*, давайте взглянем на встроенную функцию *EXTRACT*, которая может быть очень полезной при работе с такими типами. Встроенная функция *EXTRACT* имеет дело со значениями *TIMESTAMP* и *INTERVAL* и возвращает разнообразные фрагменты входящей в них информации, такие как часовой пояс из *TIMESTAMP* или часы/дни/минуты из *INTERVAL*. Давайте воспользуемся предыдущим примером, где мы получили интервал в 380 дней, 10 часов, 20 минут и 29,878 секунды:

```
EODA@ORA12CR1> select dt2-dt1
2   from (select to_timestamp('29-feb-2000 01:02:03.122000',
3                        'dd-mon-yyyy hh24:mi:ss.ff') dt1,
4                        to_timestamp('15-mar-2001 11:22:33.000000',
5                        'dd-mon-yyyy hh24:mi:ss.ff') dt2
6   from dual )
7 /

DT2-DT1
-----
+0000000380 10:20:29.878000000
```

Мы можем применить функцию *EXTRACT* для извлечения каждого фрагмента информации:

```

EODA@ORA12CR1> select extract( day    from dt2-dt1 ) day,
2      extract( hour    from dt2-dt1 ) hour,
3      extract( minute  from dt2-dt1 ) minute,
4      extract( second  from dt2-dt1 ) second
5  from (select to_timestamp('29-feb-2000 01:02:03.122000',
6                'dd-mon-yyyy hh24:mi:ss.ff') dt1,
7                to_timestamp('15-mar-2001 11:22:33.000000',
8                'dd-mon-yyyy hh24:mi:ss.ff') dt2
9        from dual )
10 /

```

| DAY | HOUR | MINUTE | SECOND |
|-----|------|--------|--------|
| 380 | 10 | 20 | 29.878 |

Вдобавок вы уже видели применение функций `NUMTOYMINTERVAL` и `NUMTODSINTERVAL` для создания интервалов `YEAR TO MONTH` и `DAY TO SECOND`. Я считаю, что эти функции предлагают простейший способ создания экземпляров типов `INTERVAL` — гораздо проще, чем функции преобразования строк. Вместо конкатенации множества чисел, представляющих дни, часы, минуты и секунды, для образования интервалов я предпочитаю для выполнения той же самой работы использовать четыре вызова `NUMTODSINTERVAL`.

Тип `INTERVAL` может применяться для хранения не только продолжительностей, но и показаний времени. Например, если требуется сохранить конкретную дату и время, то есть типы `DATE` или `TIMESTAMP`. Но что если нужно сохранить только время 8:00? Для этого может оказаться удобным тип `INTERVAL` (в частности, `INTERVAL DAY TO SECOND`).

Тип `INTERVAL YEAR TO MONTH`

Синтаксис использования типа `INTERVAL YEAR TO MONTH` очень прост:

```
INTERVAL YEAR(n) TO MONTH
```

Здесь `N` — необязательное количество цифр для поддержки числа лет, которое варьируется от 0 до 9, со стандартным значением 2 (для сохранения количества лет от 0 до 99). Этот тип позволяет хранить любое количество лет (вплоть до 9 цифр) и месяцев. Функцией, которую я предпочитаю применять для создания экземпляров `INTERVAL` этого типа, является `NUMTOYMINTERVAL`. Например, вот как можно создать интервал из 5 лет и 2 месяцев:

```

EODA@ORA12CR1> select numtoyminterval(5,'year')+numtoyminterval(2,'month')
from dual;
NUMTOYMINTERVAL(5,'YEAR')+NUMTOYMINTERVAL(2,'MONTH')
-----
+0000000005-02

```

Либо, используя единственный вызов и учитывая тот факт, что год содержит 12 месяцев, можно поступить следующим образом:

```

EODA@ORA12CR1> select numtoyminterval(5*12+2,'month') from dual;
NUMTOYMINTERVAL(5*12+2,'MONTH')
-----
+0000000005-02

```

Оба подхода работают хорошо. Еще одна функция, `TO_YMINTERVAL`, может применяться для преобразования строки в тип `INTERVAL YEAR TO MONTH`:

```
EODA@ORA12CR1> select to_yminterval( '5-2' ) from dual;
TO_YMINTERVAL('5-2')
-----
+0000000005-02
```

Но поскольку подавляющее большинство времени в своем приложении я храню год и месяцы в двух полях `NUMBER`, то нахожу использование функции `NUMTOYMINTERVAL` более удобным, чем построение форматированной строки из чисел. И, наконец, тип `INTERVAL` можно просто указать непосредственно в SQL-операторе, вообще пропуская вызовы функций:

```
EODA@ORA12CR1> select interval '5-2' year to month from dual;
INTERVAL'5-2'YEARTOMONTH
-----
+05-02
```

Тип `INTERVAL DAY TO SECOND`

Синтаксис применения типа `INTERVAL DAY TO SECOND` прост:

```
INTERVAL DAY(n) TO SECOND(m)
```

Здесь `N` — необязательное количество цифр для поддержки компонента дней, которое варьируется от 0 до 9 и имеет стандартное значение 2. Кроме того, `M` — это количество цифр для хранения дробной части секунд; оно варьируется от 0 до 9 и по умолчанию равно 6. Функцией, которую я предпочитаю использовать для создания экземпляров типа `INTERVAL DAY TO SECOND`, является `NUMTODSINTERVAL`:

```
EODA@ORA12CR1> select numtodsinterval( 10, 'day' )+
2 numtodsinterval( 2, 'hour' )+
3 numtodsinterval( 3, 'minute' )+
4 numtodsinterval( 2.3312, 'second' )
5 from dual;
NUMTODSINTERVAL(10,'DAY')+NUMTODSINTERVAL(2,'HOUR')+NUMTODSINTERVAL(3,'MINU
-----
+0000000010 02:03:02.331200000
```

или просто

```
EODA@ORA12CR1> select numtodsinterval( 10*86400+2*3600+3*60+2.3312, 'second' )
from dual;
NUMTODSINTERVAL(10*86400+2*3600+3*60+2.3312, 'SECOND')
-----
+0000000010 02:03:02.331200000
```

с учетом того факта, что в сутках содержится 86 400 секунд, в часе — 3600 секунд и т.д. В качестве альтернативы, как и ранее, можно применить функцию `TO_DSINTERVAL` для преобразования строки в тип `INTERVAL DAY TO SECOND`:

```
EODA@ORA12CR1> select to_dsinterval( '10 02:03:02.3312' ) from dual;
TO_DSINTERVAL('1002:03:02.3312')
-----
+0000000010 02:03:02.331200000
```

Можно также использовать литерал типа INTERVAL DAY TO SECOND в самом операторе SQL:

```

EODA@ORA12CR1> select interval '10 02:03:02.3312' day to second from dual;
INTERVAL '1002:03:02.3312' DAYTOSECOND
-----
+10 02:03:02.331200

```

Типы LOB

Согласно моему опыту, LOB, или *большие объекты* (large object), являются источником значительной путаницы. Это неправильно понимаемый тип данных — как в отношении его реализации, так и в отношении его эффективного применения. В настоящем разделе представлен обзор особенностей физического хранения значений LOB и соображения, которые вы должны принимать во внимание при использовании этих типов. Типы LOB имеют множество необязательных настроек, и правильное их сочетание критически важно для ваших приложений.

В Oracle поддерживаются четыре типа LOB.

- CLOB. Символьный LOB. Этот тип данных служит для хранения крупных объемов текстовой информации, такой как XML-разметка или простой текст. Он подвержен преобразованиям символьных наборов, т.е. символы в поле данного типа при извлечении будут преобразовываться из символьного набора базы данных в клиентский символьный набор, а при модификации — из символьного набора клиента в символьный набор базы данных.
- NCLOB. Еще один тип символьного LOB. Символьным набором данных, хранящихся в этом столбце, является национальный символьный набор, а не стандартный набор базы данных.
- BLOB. Двоичный LOB. Этот тип применяется для хранения крупных объемов двоичной информации, такой как документы текстового процессора, графические изображения и любые другие данные, которые только можно себе вообразить. Эти данные не подвергаются преобразованиям символьных наборов. Все биты и байты, записываемые приложением в поле BLOB, возвращаются из него без изменений.
- BFILE. Двоичный файловый LOB. Это в большей степени указатель, чем сущность, хранящаяся в базе данных. Единственное, что сохраняется в базе с BFILE — указатель на файл в операционной системе. Файл поддерживается за пределами базы данных и вообще не является ее частью. Тип BFILE предоставляет доступ к содержимому файла только по чтению.

При обсуждении типов LOB приведенный список будет разбит на две части: типы LOB, сохраняемые в базе данных, или внутренние типы LOB, к которым относятся CLOB, BLOB и NCLOB, и типы LOB, хранящиеся вне базы данных — тип BFILE. Мы не станем рассматривать типы CLOB, BLOB и NCLOB по отдельности, т.к. с точки зрения способа и вариантов хранения они одинаковы. Их отличает лишь то, что типы CLOB и NCLOB поддерживают текстовую информацию, а BLOB — нет. Но независимо от базового типа для них указываются те же самые опции CHUNK, RETENTION и т.д. Поскольку тип BFILE существенно отличается, он обсуждается отдельно.

Внутренние типы LOB

Начиная с версии Oracle 11g, стала доступной новая архитектура для LOB, именуемая SECUREFILE. Ранее существовавшая архитектура для LOB была известна под названием BASICFILE. По умолчанию в версии Oracle 11g столбец LOB создается как LOB вида BASICFILE. В версии Oracle 12c, если столбец LOB создается в табличном пространстве ASSM, то по умолчанию он будет LOB вида SECUREFILE.

Забегая вперед, я рекомендую использовать вариант SECUREFILE, а не BASICFILE, по следующим причинам.

- В документации Oracle утверждается, что в будущем выпуске вариант BASICFILE будет объявлен устаревшим.
- Для управления SECUREFILE предусмотрено меньшее число параметров, в частности, к SECUREFILE не применяются атрибуты CHUNK, PCTVERSION, FREEPOOLS, FREELISTS и FREELIST GROUPS.
- Объекты LOB вида SECUREFILE позволяют использовать расширенное шифрование, сжатие и дедупликацию. Если вы собираетесь применять эти расширенные средства LOB, то должны получить лицензию для опции расширенной безопасности (Advanced Security Option) и/или опции расширенного сжатия (Advanced Compression Option). Если использование расширенных средств LOB не планируется, можете работать с объектами LOB вида SECUREFILE без необходимости в дополнительной лицензии.

В последующих разделах мы подробно рассмотрим нюансы применения вариантов SECUREFILE и BASICFILE.

Создание LOB вида SECUREFILE

Синтаксис LOB вида SECUREFILE на первый взгляд очень прост — обманчиво прост. Вы можете создавать таблицы со столбцами типа CLOB, BLOB или NCLOB и все.

```
EODA@ORA12CR1> create table t
2 ( id int primary key,
3   txt clob
4 )
5 segment creation immediate
6 /
```

Table created.

Таблица создана.

Проверить, что столбец был создан как LOB вида SECUREFILE, можно следующим образом:

```
EODA@ORA12CR1> select column_name, securefile from user_lobs where table_name='T';
COLUMN_NAME  SECUREFILE
-----
TXT          YES
```

Если вы используете версию Oracle 11g, то стандартным типом LOB является BASICFILE, поэтому для создания LOB вида SECUREFILE понадобится применять конструкцию STORE AS SECUREFILE:

```

EODA@ORA11GR2> create table t
2 ( id int primary key,
3   txt clob
4 )
5 segment creation immediate
6 lob(txt) store as securefile
7 /

```

Table created.

Таблица создана.

На вид столбцы LOB столь же просты в использовании, как столбцы типа NUMBER, DATE или VARCHAR2. Но так ли это? Приведенные выше небольшие примеры продемонстрировали только верхушку айсберга — абсолютный минимум, который можно указывать для LOB. С помощью DBMS_METADATA можем получить полную картину:

```

EODA@ORA12CR1> select dbms_metadata.get_ddl( 'TABLE', 'T' ) from dual;
DBMS_METADATA.GET_DDL('TABLE','T')

```

```

-----
CREATE TABLE "EODA"."T"
(   "ID" NUMBER(*,0),
    "TXT" CLOB,
    PRIMARY KEY ("ID")
  USING INDEX PCTFREE 10 INITRANS 2 MAXTRANS 255
  STORAGE(INITIAL 65536 NEXT 1048576 MINEXTENTS 1 MAXEXTENTS 2147483645
  PCTINCREASE 0 FREELISTS 1 FREELIST GROUPS 1
  BUFFER_POOL DEFAULT FLASH_CACHE DEFAULT CELL_FLASH_CACHE DEFAULT)
  TABLESPACE "USERS" ENABLE
  ) SEGMENT CREATION IMMEDIATE
  PCTFREE 10 PCTUSED 40 INITRANS 1 MAXTRANS 255
  NOCOMPRESS LOGGING
  STORAGE(INITIAL 65536 NEXT 1048576 MINEXTENTS 1 MAXEXTENTS 2147483645
  PCTINCREASE 0 FREELISTS 1 FREELIST GROUPS 1
  BUFFER_POOL DEFAULT FLASH_CACHE DEFAULT CELL_FLASH_CACHE DEFAULT)
  TABLESPACE "USERS"
  LOB ("TXT") STORE AS SECUREFILE (
    TABLESPACE "USERS" ENABLE STORAGE IN ROW CHUNK 8192
    NOCACHE LOGGING NOCOMPRESS KEEP_DUPLICATES
    STORAGE(INITIAL 106496 NEXT 1048576 MINEXTENTS 1 MAXEXTENTS 2147483645
    PCTINCREASE 0
    BUFFER_POOL DEFAULT FLASH_CACHE DEFAULT CELL_FLASH_CACHE DEFAULT) )

```

Как видите, параметров довольно много. Перед тем, как перейти к детальному анализу этих параметров, в следующем разделе я сгенерирую тот же самый тип вывода для LOB вида BASICFILE. Это послужит основой при обсуждении разнообразных атрибутов LOB.

Создание LOB вида BASICFILE

В версиях, предшествующих Oracle 12c, приведенный ниже код создаст LOB вида BASICFILE:

```
EODA@ORA11GR2> create table t
2 ( id int primary key,
3   txt clob
4 )
5 segment creation immediate
6 /
```

Table created.

Таблица создана.

В версии Oracle 12c для создания LOB вида BASICFILE необходимо применять конструкцию STORE AS BASICFILE:

```
EODA@ORA12CR1> create table t
2 ( id int primary key,
3   txt clob
4 )
5 segment creation immediate
6 lob(txt) store as basicfile
7 /
```

Table created.

Таблица создана.

Используя пакет DBMS_METADATA, мы можем просмотреть сведения о LOB вида BASICFILE:

```
EODA@ORA12CR1> select dbms_metadata.get_ddl( 'TABLE', 'T' ) from dual;
DBMS_METADATA.GET_DDL('TABLE','T')
```

```
-----
CREATE TABLE "EODA"."T"
(
  "ID" NUMBER(*,0),
  "TXT" CLOB,
  PRIMARY KEY ("ID")
USING INDEX PCTFREE 10 INITRANS 2 MAXTRANS 255
STORAGE(INITIAL 65536 NEXT 1048576 MINEXTENTS 1 MAXEXTENTS 2147483645
PCTINCREASE 0 FREELISTS 1 FREELIST GROUPS 1
BUFFER_POOL DEFAULT FLASH_CACHE DEFAULT CELL_FLASH_CACHE DEFAULT)
TABLESPACE "USERS" ENABLE
) SEGMENT CREATION IMMEDIATE
PCTFREE 10 PCTUSED 40 INITRANS 1 MAXTRANS 255
NOCOMPRESS LOGGING
STORAGE(INITIAL 65536 NEXT 1048576 MINEXTENTS 1 MAXEXTENTS 2147483645
PCTINCREASE 0 FREELISTS 1 FREELIST GROUPS 1
BUFFER_POOL DEFAULT FLASH_CACHE DEFAULT CELL_FLASH_CACHE DEFAULT)
TABLESPACE "USERS"
LOB ("TXT") STORE AS BASICFILE (
TABLESPACE "USERS" ENABLE STORAGE IN ROW CHUNK 8192 RETENTION
NOCACHE LOGGING
STORAGE(INITIAL 65536 NEXT 1048576 MINEXTENTS 1 MAXEXTENTS 2147483645
PCTINCREASE 0 FREELISTS 1 FREELIST GROUPS 1
BUFFER_POOL DEFAULT FLASH_CACHE DEFAULT CELL_FLASH_CACHE DEFAULT))
```

Большинство параметров для LOB вида BASICFILE идентично параметрам LOB вида SECUREFILE. Главное отличие заключается в том, что конструкция хранения LOB вида SECUREFILE содержит меньше параметров (скажем, отсутствуют параметры FREELISTS и FREELIST GROUPS).

Компоненты LOB

Как показано в выводе DBMS_METADATA в предшествующих разделах, столбец LOB имеет несколько интересных атрибутов:

- табличное пространство (USERS в этом примере);
- ENABLE STORAGE IN ROW в качестве стандартного атрибута;
- CHUNK 8192;
- RETENTION;
- NOCACHE;
- полная конструкция STORAGE.

Перечисленные атрибуты говорят о том, что “за кулисами” LOB происходит много интересного, и на самом деле так и есть. Столбец LOB всегда дает в результате то, что я называю *много сегментным объектом*, т.е. таблица будет применять множество физических сегментов. Если мы создадим такую таблицу в пустой схеме, то обнаружим следующее:

```
EODA@ORA12CR1> select segment_name, segment_type from user_segments;
```

| SEGMENT_NAME | SEGMENT_TY |
|-----------------------------|------------|
| T | TABLE |
| SYS_LOB0000020053C00002\$\$ | LOBSEGMENT |
| SYS_IL0000020053C00002\$\$ | LOBINDEX |
| SYS_C005432 | INDEX |

В поддержку ограничения первичного ключа был создан индекс, и это нормально, но как насчет остальных двух сегментов — LOBINDEX (индекс LOB) и LOBSEGMENT (сегмент LOB)? Они были созданы для поддержки нашего столбца LOB. В сегменте LOBSEGMENT будут храниться действительные данные (так, они могут храниться также и в таблице T, но мы раскроем это более подробно, когда доберемся до конструкции ENABLE STORAGE IN ROW). Сегмент LOBINDEX используется для навигации по столбцу LOB, чтобы находить его фрагменты. Когда мы создаем столбец LOB, в общем случае в строке сохраняется *указатель*, или *локатор LOB* (LOB locator). Локатор LOB — это именно то, что извлекает наше приложение. Когда мы запрашиваем байты с 1000 по 2000 из LOB, локатор LOB применяется к LOBINDEX, чтобы найти, где хранятся эти байты, после чего производится доступ к LOBSEGMENT. Сегмент LOBINDEX служит для эффективного нахождения фрагментов LOB. Таким образом, LOB можно воспринимать как разновидность отношения “главный-подчиненный”. Столбец LOB хранится в виде порций или фрагментов, и каждый фрагмент доступен нам. Если бы пришлось реализовывать LOB с использованием только таблиц, например, то можно было бы поступить так:

```
Create table parent
( id int primary key,
  другие_данные...
);
Create table lob
( id references parent on delete cascade,
  chunk_number int,
  data <тип данных>(n),
  primary key (id, chunk_number)
);
```


Концептуально столбец LOB хранится очень похожим образом: после создания этих двух таблиц в таблице LOB мы будем иметь первичный ключ на столбцах ID, CHUNK_NUMBER (аналогично сегменту LOBINDEX, создаваемому Oracle), а в таблице LOB — фрагменты данных (аналогично сегменту LOBSEGMENT). Столбец LOB реализует эту структуру “главный-подчиненный” прозрачным образом. Идею поможет прояснить рис. 12.3.

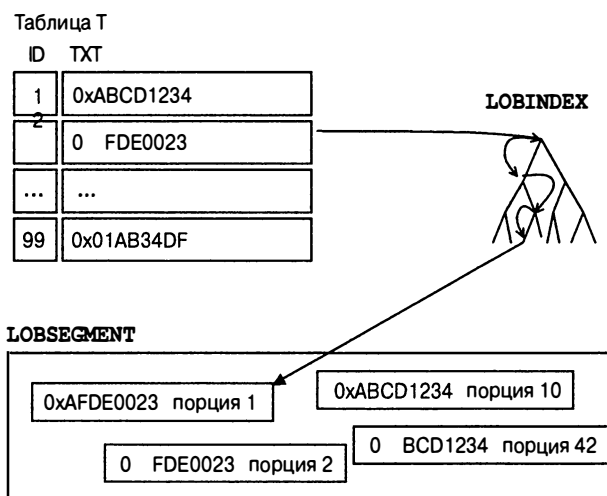


Рис. 12.3. Отношения между таблицей, LOBINDEX и LOBSEGMENT

Локаатор LOB в таблице действительно просто указывает на LOBINDEX, а LOBINDEX, в свою очередь, указывает на все фрагменты самих данных LOB. Чтобы получить байты с N до M из LOB, понадобится разыменовать указатель в таблице (локаатор LOB), пройти по структуре LOBINDEX для нахождения необходимых фрагментов, после чего по очереди обратиться к ним. Это делает произвольный доступ к любой части LOB одинаково быстрым — получить начало, середину или конец LOB можно с равной скоростью, т.к. не приходится каждый раз начинать с начала и проходить по всему LOB.

Теперь, когда вы понимаете, как концептуально хранится LOB, вспомним все перечисленные выше необязательные параметры и посмотрим, для чего они служат и что конкретно подразумевают.

Табличное пространство LOB

Оператор CREATE TABLE, возвращенный посредством DBMS_METADATA для вариантов SECUREFILE и BASICFILE, включает следующую конструкцию:

```
LOB ("TXT") STORE AS ... ( TABLESPACE "USERS" ...
```

Здесь TABLESPACE представляет табличное пространство, в котором будут сохранены сегменты LOBSEGMENT и LOBINDEX, и оно может отличаться от табличного пространства, где находится сама таблица. То есть табличное пространство, которое содержит данные LOB, может быть отдельным и отличным от табличного пространства, хранящего действительные данные таблицы.

Основные причины, по которым может быть принято решение применять для хранения данных LOB другое табличное пространство вместо той же самой таблицы, главным образом связаны с администрированием и производительностью. С точки зрения администрирования тип данных LOB представляет немалый объем информации. Если таблица содержит миллионы строк, и каждая строка имеет связанный с ней столбец LOB порядочного размера, то данные LOB могут оказаться поистине гигантскими. В этом случае имеет смысл отделить таблицу от данных LOB — просто чтобы облегчить резервное копирование и восстановление, а также управление пространством. Например, для данных LOB может быть установлен другой размер экстенда, который отличается от используемого для обычных табличных данных.

Еще одна причина касается производительности ввода-вывода. По умолчанию данные LOB не кешируются в буферном кеше (подробнее об этом позже). Следовательно, по умолчанию каждое обращение к LOB, будь то чтение или запись, означает физическую операцию ввода-вывода — прямое чтение с диска или прямую запись на диск.

На заметку! Столбцы LOB могут быть внутристрочными или храниться в таблице. В таком случае данные LOB должны кешироваться, но это применимо только к столбцам LOB размером 4000 байтов и менее. Мы обсудим это дополнительно в разделе “Конструкция IN ROW” далее в главе.

Поскольку каждый доступ — это физический ввод-вывод, имеет смысл вынести на отдельные диски те объекты, о которых известно, что они потребуют больше физических операций ввода-вывода в реальном времени (при обращении пользователя к ним), чем другие. Следует отметить, что сегменты LOBINDEX и LOBSEGMENT *всегда будут находиться в одном и том же табличном пространстве*. Разнести их по разным табличным пространствам не удастся. Очень ранние версии Oracle разрешали это делать, но начиная с версии Oracle8i Release 3, это стало невозможным. В действительности, как вы вскоре увидите, все характеристики хранения сегмента LOBINDEX наследуются от LOBSEGMENT.

Конструкция IN ROW

Оператор CREATE TABLE, возвращенный с помощью DBMS_METADATA для вариантов SECUREFILE и BASICFILE, включает такую конструкцию:

```
LOB ("TXT") STORE AS ... (... ENABLE STORAGE IN ROW ...
```

Эта конструкция управляет тем, хранятся ли данные LOB всегда отдельно от таблицы в сегменте LOBSEGMENT, или же иногда они могут храниться прямо в таблице без их помещения в LOBSEGMENT. Если мы установим ENABLE STORAGE IN ROW как противоположность DISABLE STORAGE IN ROW, то небольшие столбцы LOB до 4000 байтов будут храниться в самой таблице, что очень похоже на VARCHAR2. Только когда размер столбцов LOB превышает 4000 байтов, они будут перемещены из строки в сегмент LOBSEGMENT.

Разрешение хранения в строке действует по умолчанию и обычно так и необходимо поступать, если известно, что во многих случаях столбцы LOB умещаются в самой таблице. Например, у вас может быть приложение с полем описания определенного рода. Описание может содержать от 0 до 32 Кбайт данных (возможно,

даже больше, но в основном до 32 Кбайт). Известно, что многие описания будут очень короткими, состоящими из пары сотен символов. Вместо того чтобы страдать от накладных расходов, связанных с хранением описаний за пределами строк, и получать к ним доступ через индекс при каждом извлечении, описания можно хранить внутри строк, т.е. в самой таблице. Более того, если столбец LOB использует стандартную установку NOCACHE (содержимое LOBSEGMENT не кешируется в буферном кеше), то данные LOB, хранящиеся в сегменте таблицы, позволят избежать физического ввода-вывода, необходимого для их извлечения.

На заметку! Начиная с версии Oracle 12c, можно создавать столбец типа VARCHAR2, NVARCHAR2 или RAW, который будет хранить вплоть до 32 767 байтов информации. За детальными сведениями обращайтесь в раздел “Расширенные типы данных” ранее в этой главе.

Эффект можно увидеть на довольно простом примере. Создадим таблицу со столбцом LOB, который может хранить данные в строке, и столбцом LOB, который не может:

```
EODA@ORA12CR1> create table t
2  ( id int    primary key,
3    in_row clob,
4    out_row clob
5  )
6  lob (in_row) store as ( enable storage in row )
7  lob (out_row) store as ( disable storage in row )
8  /
```

Table created.

Таблица создана.

Мы вставим в эту таблицу строковые данные с длиной, не превышающей 4000 байтов:

```
EODA@ORA12CR1> insert into t
2  select rownum,
3         owner || ' ' || object_name || ' ' || object_type || ' ' || status,
4         owner || ' ' || object_name || ' ' || object_type || ' ' || status
5  from all_objects
6  /
```

72085 rows created.

72085 строк создано.

```
EODA@ORA12CR1> commit;
```

Commit complete.

Фиксация завершена.

Если теперь попробовать прочитать каждую строку, то с применением пакета DBMS_MONITOR при включенном средстве SQL_TRACE можно оценить производительность извлечения данных для каждого столбца:

```
EODA@ORA12CR1> declare
2      l_cnt    number;
3      l_data   varchar2(32765);
4  begin
5      select count(*)
6      into l_cnt
7      from t;
```

```

8
9      dbms_monitor.session_trace_enable;
10     for i in 1 .. l_cnt
11     loop
12         select in_row into l_data from t where id = i;
13         select out_row into l_data from t where id = i;
14     end loop;
15 end;
16 /

```

PL/SQL procedure successfully completed.

Процедура PL/SQL успешно завершена.

Когда мы просмотрим отчет TKPROF для этой небольшой эмуляции, то результаты окажутся достаточно очевидными:

SELECT IN_ROW FROM T WHERE ID = :B1

| call | count | cpu | elapsed | disk | query | current | rows |
|---------|-------|------|---------|------|-------|---------|-------|
| Parse | 1 | 0.00 | 0.00 | 0 | 0 | 0 | 0 |
| Execute | 18240 | 0.23 | 0.25 | 0 | 0 | 0 | 0 |
| Fetch | 18240 | 0.22 | 0.27 | 0 | 54720 | 0 | 18240 |
| total | 36481 | 0.46 | 0.53 | 0 | 54720 | 0 | 18240 |

SELECT OUT_ROW FROM T WHERE ID = :B1

| call | count | cpu | elapsed | disk | query | current | rows |
|---------|-------|------|---------|-------|-------|---------|-------|
| Parse | 1 | 0.00 | 0.00 | 0 | 0 | 0 | 0 |
| Execute | 18240 | 0.23 | 0.24 | 0 | 0 | 0 | 0 |
| Fetch | 18240 | 1.95 | 1.67 | 18240 | 72960 | 0 | 18240 |
| total | 36481 | 2.18 | 1.91 | 18240 | 72960 | 0 | 18240 |

Elapsed times include waiting on following events:

| Event waited on | Times Waited | Max. Wait | Total Waited |
|------------------|-----------------|-----------|--------------|
| direct path read | 18240 | 0.00 | 0.14 |

Извлечение столбца IN_ROW происходит значительно быстрее и потребляет намного меньше ресурсов. Здесь видно, что оно потребовало 54 720 логических операций ввода-вывода (извлечений в режиме запроса), тогда как для извлечения столбца OUT_ROW понадобилось существенно больше таких операций. На первый взгляд не вполне понятно, откуда взялись эти дополнительные операции ввода-вывода, но если вы вспомните, каким образом хранятся столбцы LOB, то все станет совершенно ясным. Эти операции ввода-вывода в сегменте LOBINDEX, выполняемые для поиска фрагментов LOB. Все дополнительные логические операции ввода-вывода касаются сегмента LOBINDEX.

Кроме того, легко заметить, что извлечение 18 240 строк при хранении столбцов LOB за пределами строк выливается в 18 240 физических операций ввода-вывода и приводит к 18 240 ожиданиям чтения в прямом режиме. Это были операции чтения некешированных данных LOB. В рассматриваемом случае их можно сократить, если включить кеширование данных LOB, но тогда необходимо удостовериться, что для этого хватит объема буферного кеша. Вдобавок, если бы в таблице присутствовали

действительно большие столбцы LOB, то вряд ли захочется обеспечивать кеширование для таких данных.

Этот режим хранения в строке или вне строки наряду с чтением затрагивает также модификации. Обновим первые 100 строк таблицы короткими значениями столбцов и вставим 100 новых строк с короткими значениями, используя тот же прием мониторинга производительности:

```
EODA@ORA12CR1> create sequence s start with 100000;
```

```
Sequence created.
```

Последовательность создана.

```
EODA@ORA12CR1> declare
```

```
2         l_cnt      number;
```

```
3         l_data     varchar2(32765);
```

```
4 begin
```

```
5     dbms_monitor.session_trace_enable;
```

```
6     for i in 1 .. 100
```

```
7     loop
```

```
8         update t set in_row =
```

```
to_char(sysdate,'dd-mon-yyyy hh24:mi:ss') where id = i;
```

```
9         update t set out_row =
```

```
to_char(sysdate,'dd-mon-yyyy hh24:mi:ss') where id = i;
```

```
10        insert into t (id, in_row) values ( s.nextval, 'Hello World' );
```

```
11        insert into t (id,out_row) values ( s.nextval, 'Hello World' );
```

```
12    end loop;
```

```
13 end;
```

```
14 /
```

```
PL/SQL procedure successfully completed.
```

Процедура PL/SQL успешно завершена.

В результирующем отчете TKPROF обнаружится следующее:

```
UPDATE T SET IN_ROW = TO_CHAR(SYSDATE,'dd-mon-yyyy hh24:mi:ss') WHERE ID = :B1
```

| call | count | cpu | elapsed | disk | query | current | rows |
|---------|-------|------|---------|------|-------|---------|------|
| Parse | 1 | 0.00 | 0.00 | 0 | 0 | 0 | 0 |
| Execute | 100 | 0.00 | 0.01 | 0 | 200 | 214 | 100 |
| Fetch | 0 | 0.00 | 0.00 | 0 | 0 | 0 | 0 |
| total | 101 | 0.00 | 0.01 | 0 | 200 | 214 | 100 |

```
Misses in library cache during parse: 1
```

```
Misses in library cache during execute: 1
```

```
Optimizer mode: ALL_ROWS
```

```
Parsing user id: 66 (recursive depth: 1)
```

```
Number of plan statistics captured: 1
```

```
Rows (1st) Rows (avg) Rows (max) Row Source Operation
```

| Rows (1st) | Rows (avg) | Rows (max) | Row Source Operation |
|------------|------------|------------|---|
| 0 | 0 | 0 | UPDATE T (cr=2 pr=0 pw=0 time=463 us) |
| 1 | 1 | 1 | INDEX UNIQUE SCAN SYS_C005434 (cr=2 pr=0 pw=0 time=16... |

```
*****
UPDATE T SET OUT_ROW = TO_CHAR(SYSDATE,'dd-mon-yyyy hh24:mi:ss') WHERE ID = :B1
```

| call | count | cpu | elapsed | disk | query | current | rows |
|---------|-------|------|---------|------|-------|---------|------|
| Parse | 1 | 0.00 | 0.00 | 0 | 0 | 0 | 0 |
| Execute | 100 | 0.03 | 0.99 | 0 | 200 | 302 | 100 |
| Fetch | 0 | 0.00 | 0.00 | 0 | 0 | 0 | 0 |
| total | 101 | 0.03 | 0.99 | 0 | 200 | 302 | 100 |

Misses in library cache during parse: 1

Misses in library cache during execute: 1

Optimizer mode: ALL_ROWS

Parsing user id: 66 (recursive depth: 1)

Number of plan statistics captured: 1

Rows (1st) Rows (avg) Rows (max) Row Source Operation

| | | | |
|---|---|---|--|
| 0 | 0 | 0 | UPDATE T (cr=2 pr=0 pw=1 time=8759 us) |
| 1 | 1 | 1 | INDEX UNIQUE SCAN SYS_C005434 (cr=2 pr=0 pw=0 time=6... |

Elapsed times include waiting on following events:

| Event waited on | Times Waited | Max. Wait | Total Waited |
|--------------------------|-----------------|-----------|--------------|
| Disk file operations I/O | 1 | 0.00 | 0.00 |
| direct path write | 163 | 0.01 | 0.96 |

Как видите, обновление столбца LOB, хранящегося вне строки, потребляет заметно больше ресурсов. Некоторое время расходуется на выполнение записей в прямом режиме (физических операций ввода-вывода) и производится намного больше извлечений текущего режима. Это объясняется тем фактом, что сегменты LOBINDEX и LOBSEGMENT нуждаются в обслуживании, как и сама таблица. Действие INSERT демонстрирует то же самое расхождение:

INSERT INTO T (ID, IN_ROW) VALUES (S.NEXTVAL, 'Hello World')

| call | count | cpu | elapsed | disk | query | current | rows |
|---------|-------|------|---------|------|-------|---------|------|
| Parse | 1 | 0.00 | 0.00 | 0 | 0 | 0 | 0 |
| Execute | 100 | 0.00 | 0.00 | 0 | 4 | 317 | 100 |
| Fetch | 0 | 0.00 | 0.00 | 0 | 0 | 0 | 0 |
| total | 101 | 0.00 | 0.00 | 0 | 4 | 317 | 100 |

INSERT INTO T (ID,OUT_ROW) VALUES (S.NEXTVAL, 'Hello World')

| call | count | cpu | elapsed | disk | query | current | rows |
|---------|-------|------|---------|------|-------|---------|------|
| Parse | 1 | 0.00 | 0.00 | 0 | 0 | 0 | 0 |
| Execute | 100 | 0.02 | 0.61 | 0 | 4 | 440 | 100 |
| Fetch | 0 | 0.00 | 0.00 | 0 | 0 | 0 | 0 |
| total | 101 | 0.02 | 0.61 | 0 | 4 | 440 | 100 |

...

Elapsed times include waiting on following events:

| Event waited on | Times Waited | Max. Wait | Total Waited |
|-------------------|-----------------|-----------|--------------|
| direct path write | 100 | 0.01 | 0.60 |

Обратите внимание на увеличившееся количество операций ввода-вывода — как чтения, так и записи. Все это доказывает, что если применяется столбец типа CLOB, и многие значения предположительно должны уместиться в строку таблицы (т.е. будут меньше 4000 байтов), то использование стандартной установки ENABLE STORAGE IN ROW будет неплохой идеей.

Конструкция CHUNK

Столбцы LOB сохраняются порциями; индекс, указывающий на данные LOB, указывает на отдельные порции данных. *Порции* (chunk) — это логически непрерывные наборы блоков, которые являются минимальными единицами выделения для LOB, в то время как обычно наименьшей единицей выделения является один блок. Размер CHUNK должен быть целым числом, кратным размеру блока в Oracle — допускаются только такие значения.

На заметку! Конструкция CHUNK применима только к LOB вида BASICFILE. Она присутствует в синтаксисе для SECUREFILE лишь с целью обратной совместимости.

Вы должны позаботиться о выборе размера CHUNK с учетом двух точек зрения. Прежде всего, каждый экземпляр LOB (каждое значение LOB, сохраненное за пределами строки) будет потреблять минимум одну порцию CHUNK. Одна порция CHUNK используется одиночным значением LOB. Если таблица содержит 100 строк, и каждая строка включает столбец LOB с 7 Кбайт данных в нем, можно с уверенностью утверждать, что выделятся 100 порций. Если установить размер CHUNK в 32 Кбайт, то получится 100 выделенных порций по 32 Кбайт. Если установить размер CHUNK в 8 Кбайт, то получится (вероятно) 100 порций по 8 Кбайт. Суть в том, что порция задействуется только одной записью LOB (две записи LOB не будут применять одну порцию CHUNK). В случае выбора размера CHUNK, который не соответствует ожидаемым размерам столбцов LOB, в конечном итоге придется расходувать впустую чрезмерный объем пространства. Например, если есть таблица со столбцом LOB, содержащим в среднем 7 Кбайт данных, а размер CHUNK установлен в 32 Кбайт, то будет неоправданно тратить впустую примерно 25 Кбайт пространства на каждый экземпляр LOB. С другой стороны, в случае использования размера CHUNK в 8 Кбайт излишний расход пространства будет сведен к минимуму.

Необходимо также проявлять осторожность, когда планируется минимизировать количество порций CHUNK на экземпляре LOB. Как вы уже видели, для указания на отдельные порции применяется сегмент LOBINDEX, и чем больше имеется порций, тем больше этот индекс. Если есть столбец LOB с размером 4 Мбайт, а размер CHUNK составляет 8 Кбайт, то для сохранения этой информации потребуется 512 порций CHUNK. Это значит, что для указания на эти порции понадобится, по крайней мере, достаточное количество записей в LOBINDEX. Может показаться, что количество этих записей не настолько велико, но вспомните, что речь идет о каждом экземпляре LOB; при наличии тысяч экземпляров LOB по 4 Мбайт получится много тысяч записей. В итоге также будет оказано влияние на производительность извлечения, поскольку чтение и управление множеством небольших порций займет больше времени, чем чтение и управление меньшим числом более крупных порций. Конечная цель заключается в выборе такого размера CHUNK, который минимизирует излишние расходы пространства, но также позволяет эффективно хранить данные.

Конструкция *RETENTION*

Конструкция *RETENTION* отличается в зависимости от того, какая парадигма хранения используется — *SECUREFILE* или *BASICFILE*. Если вы снова взглянете на вывод *DBMS_METADATA*, приведенный в начале раздела “Внутренние типы *LOB*”, то заметите, что в операторе *CREATE TABLE* для таблицы со столбцом *LOB* вида *SECUREFILE* конструкция *RETENTION* отсутствует, но она указана в случае создания таблицы со столбцом *LOB* вида *BASICFILE*. Причина в том, что для *SECUREFILE* конструкция *RETENTION* включается автоматически.

Конструкция *RETENTION* применяется для управления согласованностью чтения столбца *LOB*. Отличия в обработке *RETENTION* для случаев *SECUREFILE* и *BASICFILE* будут объяснены в последующих разделах.

Согласованность чтения для столбцов *LOB*

В предшествующих главах мы обсуждали согласованность чтения, многоверсионность и роль во всем этом информации отмены. Когда речь идет о столбцах *LOB*, способ обеспечения согласованности чтения меняется. Сегмент *LOBSEGMENT* не использует журнал отмены для записи своих изменений; версии его информации помещаются непосредственно в *LOBSEGMENT*. Сегмент *LOBINDEX* генерирует информацию отмены, как и любой другой сегмент, но *LOBSEGMENT* — нет. Вместо этого при модификации столбца *LOB* база данных Oracle выделяет новую порцию *CHUNK* и оставляет старую порцию *CHUNK* на месте. Если вы производите откат своей транзакции, то изменения в индексе *LOB* отменяются, поэтому индекс снова будет указывать на старую порцию *CHUNK*. Таким образом, поддержание информации отмены осуществляется прямо в сегменте *LOBSEGMENT*. При модификации данных старые данные остаются на месте, а новые данные создаются.

Это также вступает в игру во время чтения данных *LOB*. Сегменты *LOB* согласованы по чтению, как и все прочие сегменты. Если вы извлекаете локатор *LOB* в 9:00, то полученные из него данные *LOB* будут находиться в состоянии на 9:00. Подобно открытию курсора (результатирующего набора) в 9:00, входящие в него строки будут такими, какими они были в указанный момент времени. Даже если кто-то другой модифицирует данные *LOB* и затем фиксирует (или не фиксирует) изменения, ваш локатор *LOB* останется в состоянии на 9:00, как было бы в случае результатирующего набора. Здесь Oracle применяет сегмент *LOBSEGMENT* вместе с согласованным по чтению представлением *LOBINDEX*, чтобы отменить изменения в *LOB* и предоставить вам данные *LOB* в том виде, в каком они находились при извлечении локатора *LOB*. Информация отмены для *LOBSEGMENT* не используется, потому что для самого сегмента *LOBSEGMENT* она не генерируется.

Продемонстрировать согласованность чтения столбцов *LOB* довольно легко. Рассмотрим следующую небольшую таблицу со столбцом *LOB*, хранящимся за пределами строки (в *LOBSEGMENT*):

```
EODA@ORA12CR1> create table t
2 ( id int    primary key,
3   txt      clob
4 )
5 lob( txt) store as ( disable storage in row )
6 /
Table created.
```


Таблица создана.

```
EODA@ORA12CR1> insert into t values ( 1, 'hello world' );
1 row created.
1 строка создана.
EODA@ORA12CR1> commit;
Commit complete.
Фиксация завершена.
```

Если мы извлечем локатор LOB и откроем курсор на этой таблице:

```
EODA@ORA12CR1> declare
2         l_clob clob;
3
4         cursor c is select id from t;
5         l_id number;
6 begin
7         select txt into l_clob from t;
8         open c;
```

а затем модифицируем строку и произведем фиксацию:

```
9
10        update t set id = 2, txt = 'Goodbye';
11        commit;
12
```

то при работе с локатором LOB и открытым курсором увидим, что данные представлены в состоянии, в котором они находились в момент их извлечения или открытия:

```
13        dbms_output.put_line( dbms_lob.substr( l_clob, 100, 1 ) );
14        fetch c into l_id;
15        dbms_output.put_line( 'id = ' || l_id );
16        close c;
17 end;
18 /
hello world
id = 1
```

PL/SQL procedure successfully completed.

Процедура PL/SQL успешно завершена.

Но в базе данных информация все же была обновлена:

```
EODA@ORA12CR1> select * from t;

      ID TXT
-----
      2 Goodbye
```

Согласованные по чтению образы для курсора C поступают из сегментов отмены, тогда как согласованные по чтению образы для LOB берутся из самого сегмента LOB. Таким образом, это создает нам причину для беспокойства: если сегменты отмены не применяются для хранения информации отката, связанной со столбцами LOB, и сегменты LOB поддерживают согласованность чтения, как можно предотвратить возникновение пресловутой ошибки ORA-01555: snapshot too old

(ORA-1555: устаревший снимок)? И, что более важно — как проконтролировать объем пространства, используемого этими старыми версиями? Именно здесь в игру вступает конструкция `RETENTION` и ее альтернатива `PCTVERSION`.

Конструкция `RETENTION` для случая `BASICFILE`

Конструкция `RETENTION` сообщает базе данных о необходимости удержания модифицированных данных в сегменте `LOB` согласно параметру `UNDO_RETENTION` базы данных. Если вы установите `UNDO_RETENTION` в 2 дня, то Oracle попытается не применять повторно пространство `LOB`, освободившееся в результате модификации. То есть, если вы удалите все строки, указывающие на столбцы `LOB`, то Oracle предпримет попытку удерживать данные в сегменте `LOB` (удаленные данные) в течение 2 дней, чтобы соблюсти политику `UNDO_RETENTION`, подобно тому, как на протяжении 2 дней в табличном пространстве `UNDO` удерживается информация отмены для структурированных данных (реляционных строк и столбцов). Важно понимать, что освобожденное пространство в сегменте `LOB` не будет немедленно повторно использоваться последующими операциями `INSERT` или `UPDATE`. Это часто вызывает вопросы вроде того, почему сегмент `LOB` постоянно растет и растет? Массовая очистка, за которой следует загрузка информации заново, приведет к тому, что сегмент `LOB` будет просто разрастаться, поскольку время хранения пока не истекло.

На заметку! Для применения конструкции `RETENTION` столбец `LOB` вида `BASICFILE` должен находиться в табличном пространстве `ASSM`. Если столбец `LOB` вида `BASICFILE` расположен в табличном пространстве `MSSM`, то параметр `RETENTION` игнорируется. Обсуждение методов управления `ASSM` и `MSSM` можно найти в главе 10.

В качестве альтернативы конструкция хранения `LOB` вида `BASICFILE` может использовать параметр `PCTVERSION`, управляющий процентом выделенного (занимаемого столбцами `LOB` в определенный момент и блоков до `HWM`-уровня в сегменте `LOBSEGMENT`) пространства `LOB`, которое должно применяться для поддержки версий данных `LOB`. Стандартное значение 10% будет адекватным для многих ситуаций, т.к. в большинстве случаев данные `LOB` только вставляются и извлекаются (обновление столбцов `LOB` обычно не производится; они вставляются один раз, после чего многократно извлекаются). Следовательно, для поддержки версий данных `LOB` требуется не особо много пространства.

Тем не менее, при наличии приложения, которое интенсивно модифицирует столбцы `LOB`, стандартное значение 10% может оказаться слишком малым, если столбцы `LOB` часто читаются в то же самое время, когда другой сеанс их модифицирует. Если в процессе обработки данных `LOB` возникает ошибка `ORA-22924`, то решение не должно заключаться в увеличении размера табличного пространства отмены, увеличении длительности хранения информации отмены или добавлении сегментов отката при ручном управлении пространством отмены. Вместо этого необходимо с помощью оператора

```
ALTER TABLE имя_таблицы MODIFY LOB (имя_столбца_LOB) ( PCTVERSION n );
```

расширить объем пространства, используемого в сегменте `LOBSEGMENT` для хранения версий данных.

Конструкция RETENTION для случая SECUREFILE

В случае SECUREFILE конструкция RETENTION применяется для управления согласованностью чтения (подобно случаю BASICFILE). В выводе DBMS_METADATA оператора CREATE TABLE для таблицы со столбцом LOB вида SECUREFILE конструкция RETENTION отсутствует. Причина в том, что по умолчанию параметр RETENTION установлен в AUTO, и это инструктирует Oracle о необходимости хранения информации отмены на протяжении периода, достаточного для обеспечения согласованности чтения.

Скорректировать стандартное поведение RETENTION можно с помощью следующих параметров.

- Используйте параметр MAX для указания на то, что информация отмены должна удерживаться до тех пор, пока размер сегмента LOB не достигнет величины, заданной в параметре MAXSIZE конструкции хранения (поэтому параметр MAX применяется обязательно в сочетании с MAXSIZE).
- Устанавливайте MIN N, если включена ретроспективная база данных, чтобы ограничить продолжительность хранения информации отмены для сегмента LOB до N секунд.
- Устанавливайте NONE, если информация отмены не требуется для согласованных чтений или ретроспективных запросов.

Если вы не установите параметр RETENTION для столбца LOB вида SECUREFILE или укажете просто RETENTION, то он будет установлен в DEFAULT (что эквивалентно AUTO).

Конструкция CACHE

Оператор CREATE TABLE, возвращенный посредством DBMS_METADATA для вариантов SECUREFILE и BASICFILE, включает следующую конструкцию:

```
LOB ("TXT") STORE AS ... (... NOCACHE ...)
```

Альтернативой NOCACHE является CACHE или CACHE READS. Эта конструкция управляет тем, сохраняются ли данные LOBSEGMENT в буферном кеше. Стандартное значение NOCACHE подразумевает, что каждый доступ будет прямым чтением с диска и каждая запись/модификация также будет прямым чтением с диска. Значение CACHE READS позволяет буферизовать данные LOB, прочитанные с диска, но запись данных LOB будет производиться непосредственно на диск. Значение CACHE разрешает кеширование данных LOB как при чтении, так и при записи.

Во многих случаях стандартное значение может оказаться не тем, чем нужно. Кеширование столбцов LOB небольших и средних размеров (используемых, например, для хранения описаний из пары килобайтов) имеет смысл. Если они не кешируются, то когда пользователь обновляет такое поле описания, ему придется также ожидать завершения операций ввода-вывода, записывающих данные (размером CHUNK) на диск. Если выполняется крупная загрузка многочисленных столбцов LOB, понадобится ожидать завершения ввода-вывода при загрузке каждой строки. Для таких столбцов LOB есть основание включить кеширование. Кеширование легко включать и отключать, чтобы оценить результаты, которые оно обеспечивает:

```
ALTER TABLE tablename MODIFY LOB (lobname) ( CACHE );
ALTER TABLE tablename MODIFY LOB (lobname) ( NOCACHE );
```

Для крупной начальной загрузки имеет смысл включить кеширование столбцов LOB и позволить процессу DBWR записывать данные LOB на диск в фоновом режиме, пока клиентское приложение продолжает загрузку. Для столбцов LOB небольших и средних размеров, к которым часто производится доступ или модификация, кеширование имеет смысл, чтобы не заставлять конечных пользователей ожидать завершения физического ввода-вывода в реальном времени. Однако для столбца LOB размером в 50 Мбайт, скорее всего, в кешировании нет смысла.

Совет. Имейте в виду, что вы можете здесь с успехом применять удерживающий или рециклирующий пул (см. главу 4). Вместо кеширования данных LOBSEGMENT в стандартном кеше вместе со всеми обычными данными, можно использовать удерживающий или рециклирующий пул, чтобы отделить их. Это позволит достичь цели кеширования данных LOB, не влияя на кеширование существующих данных в системе.

Конструкция STORAGE

Наконец, оператор CREATE TABLE, возвращенный посредством DBMS_METADATA для варианта SECUREFILE, содержал следующую конструкцию:

```
LOB ("TXT") STORE AS SECUREFILE (...
STORAGE (INITIAL 106496 NEXT 1048576 MINEXTENTS 1 MAXEXTENTS 2147483645
PCTINCREASE 0
BUFFER_POOL DEFAULT FLASH_CACHE DEFAULT CELL_FLASH_CACHE DEFAULT))
```

Для варианта BASICFILE конструкция выглядела так:

```
LOB ("TXT") STORE AS BASICFILE (...
STORAGE (INITIAL 65536 NEXT 1048576 MINEXTENTS 1 MAXEXTENTS 2147483645
PCTINCREASE 0 FREELISTS 1 FREELIST GROUPS 1
BUFFER_POOL DEFAULT FLASH_CACHE DEFAULT CELL_FLASH_CACHE DEFAULT))
```

Оба варианта, SECUREFILE и BASICFILE, имеют полную конструкцию STORAGE, которую можно задействовать для управления характеристиками физического хранения. Следует отметить, что эта конструкция применяется к сегментам LOBSEGMENT и LOBINDEX в равной мере — установка для одного из них используется другим.

Управление хранилищем в случае SECUREFILE реализуется менее сложно, чем для BASICFILE. Вспомните, что столбец LOB вида SECUREFILE должен быть создан внутри табличного пространства ASSM, поэтому следующие атрибуты больше не применяются: FREELISTS, FREELIST GROUPS и FREEPOOLS.

Для LOB вида BASICFILE подходящими установками будут FREELISTS, FREELIST GROUPS и BUFFER_POOL (если не используется ASSM, как обсуждалось в главе 10). Те же самые правила применяются к сегменту LOBINDEX, потому что LOBINDEX управляется таким же образом, как любой другой индексный сегмент. Если у вас происходят интенсивные параллельные модификации столбцов LOB, можно рекомендовать использование нескольких списков свободных блоков в индексном сегменте.

Как упоминалось в предыдущем разделе, применение удерживающего или рециклирующего пула для сегментов LOB может оказаться полезным приемом, который позволит кешировать данные LOB, не повреждая существующий стандартный буферный кеш. Вместо того чтобы заводить для столбцов LOB отдельные от нор-

мальных таблиц буферы блоков, можно выделить специальную область памяти в SGA, предназначенную только для таких объектов. Для этого может использоваться конструкция `BUFFER_POOL`.

Тип BFILE

Последним из типов LOB, который мы рассмотрим, является BFILE. Тип BFILE — это просто указатель на файл в операционной системе. Он применяется для предоставления доступа *только по чтению* к файлам операционной системы.

На заметку! Встроенный пакет `UTL_FILE` обеспечивает также доступ по записи и чтению к файлам операционной системы. Тем не менее, он не использует тип BFILE.

Когда применяется тип BFILE, то также используется объект DIRECTORY, предлагаемый Oracle. Объект DIRECTORY просто отображает каталог операционной системы на строку либо имя в базе данных (он предназначен для переносимости; вы ссылаетесь на строки в BFILE, не имея дело со специфичными для операционной системы именами файлов). В качестве небольшого примера давайте создадим таблицу со столбцом BFILE и объект DIRECTORY, а затем вставим в таблицу строку, ссылающуюся на файл внутри файловой системы:

```
EODA@ORA12CR1> create table t
  2 ( id      int primary key,
  3   os_file bfile
  4 )
  5 /
Table created.
Таблица создана.

EODA@ORA12CR1> create or replace directory my_dir as '/tmp/';
Directory created.
Каталог создан.

EODA@ORA12CR1> insert into t values ( 1, bfilename( 'MY_DIR', 'test.dmp' ) );
1 row created.
1 строка создана.
```

Для этого примера в каталоге `/tmp` будет создан файл `test.dmp` с помощью команды `dd` операционной системы UNIX/Linux:

```
dd if=/dev/zero of=/tmp/test.dmp bs=1056768 count=1
```

Теперь столбец BFILE можно трактовать, как если бы он был столбцом LOB — поскольку он таковым является. Например:

```
EODA@ORA12CR1> select dbms_lob.getlength(os_file) from t;
DBMS_LOB.GETLENGTH(OS_FILE)
-----
1056768
```

Как видите, файл, на который указывает столбец BFILE, имеет размер 1 Мбайт. Обратите внимание, что в операторе `INSERT` строка `MY_DIR` применялась намеренно. В случае использования смешанных регистров или только нижнего регистра мы получили бы такой результат:

```

EODA@ORA12CR1> update t set os_file = bfilename( 'my_dir', 'test.dmp' );
1 row updated.
1 строка обновлена.

EODA@ORA12CR1> select dbms_lob.getlength(os_file) from t;
select dbms_lob.getlength(os_file) from t
*
ERROR at line 1:
ORA-22285: non-existent directory or file for GETLENGTH operation
ORA-06512: at "SYS.DBMS_LOB", line 850
ОШИБКА в строке 1:
ORA-22285: несуществующий каталог или файл в операции GETLENGTH
ORA-06512: в SYS.DBMS_LOB, строка 850

```

Пример показывает, что объекты DIRECTORY в Oracle являются идентификаторами, а идентификаторы по умолчанию сохраняются в верхнем регистре. Встроенная функция BFILENAME принимает строку, и регистр символов в этой строке должен совпадать с регистром символов в названии объекта DIRECTORY, как оно сохранено в словаре данных. Таким образом, нужно либо применять верхний регистр в строке, передаваемой функции BFILENAME, либо указывать идентификаторы в кавычках при создании объекта DIRECTORY:

```

EODA@ORA12CR1> create or replace directory "my_dir" as '/tmp/';
Directory created.
Каталог создан.

EODA@ORA12CR1> select dbms_lob.getlength(os_file) from t;
DBMS_LOB.GETLENGTH(OS_FILE)
-----
1056768

```

Я не рекомендую использовать идентификаторы в кавычках; взамен при вызове функции BFILENAME указывайте имя в верхнем регистре. Идентификаторы в кавычках обычно не применяются и чреваты путаницей в будущем.

В зависимости от длины имени объекта DIRECTORY и имени файла объект BFILE (объект указателя в базе данных, а не действительный двоичный файл на диске) потребляет разный объем дискового пространства. В предыдущем примере результирующий объект BFILE имел длину около 35 байтов. В общем случае вы обнаружите, что с объектом BFILE связано примерно 20 байтов накладных расходов *плюс* длина имени объекта DIRECTORY *плюс* длина самого имени файла.

На заметку! Данные BFILE не являются согласованными по чтению как данные других типов LOB. Поскольку объект BFILE управляется вне базы данных, то при разыменовании BFILE вы получите все, что содержалось на тот момент в файле. Поэтому повторяющиеся чтения из одного и того же объекта BFILE могут давать разные результаты — в отличие от локатора LOB, используемого со столбцом типа CLOB, BLOB или NCLOB.

Типы ROWID и UROWID

В этой главе нам осталось рассмотреть еще типы ROWID и UROWID. Тип ROWID — это адрес строки в таблице (вспомните из главы 10, что для уникальной идентификации строки в базе данных требуется ROWID и имя таблицы). В ROWID закодировано

достаточно информации, чтобы найти строку на диске, а также идентифицировать объект, на который указывает ROWID (таблица и т.д.). Близкий родственник ROWID — UROWID — представляет собой универсальный идентификатор ROWID и применяется для таких таблиц, как индекс-таблицы, и таблиц, доступных через шлюзы в гетерогенных базах данных, которые не имеют фиксированных ROWID. Тип UROWID — это представление значения первичного ключа строки, и потому его размер варьируется в зависимости от объекта, на который указывает.

Каждая строка в любой таблице имеет либо ROWID, либо UROWID, ассоциированный с ней. При извлечении из таблицы они рассматриваются как *псевдостолбцы* — в том смысле, что в действительности они не хранятся вместе со строкой, а являются ее производными атрибутами. ROWID генерируется на основе физического местоположения строки и с ней не сохраняется. UROWID генерируется на основе первичного ключа строки, поэтому в определенном смысле он хранится со строкой, но не на самом деле, поскольку UROWID не существует в виде отдельного столбца, а скорее является функцией от существующих столбцов.

Раньше считалось, что для строк с ROWID (наиболее распространенный тип строк в Oracle; за исключением строк в индекс-таблицах все строки имеют ROWID) значения ROWID должны быть неизменяемыми. Когда строка вставлялась, она должна была ассоциироваться с ROWID — адресом — и этот адрес должен был быть связанным с ней до тех пор, пока строка физически не удалялась из базы данных. Со временем это становится все менее верно, т.к. теперь есть операции, которые могут вызвать изменение ROWID строки; их примеры перечислены ниже.

- Обновление ключа секционирования строки в секционированной таблице, в результате чего строка должна переместиться из одной секции в другую.
- Использование команды FLASHBACK TABLE для восстановления таблицы базы данных в состояние, которое она имела в какой-то момент времени в прошлом.
- Выполнение операций MOVE, а также многих операций с секциями, таких как разделение и слияние секций.
- Применение команды ALTER TABLE SHRINK SPACE для усечения сегмента.

Теперь с учетом того, что ROWID могут со временем изменяться (они больше не являются неизменяемыми), физически хранить их в виде столбцов таблиц базы данных не рекомендуется. То есть использование ROWID в качестве типа данных для столбца считается плохой практикой, которой следует избегать. Взамен должен применяться первичный ключ строки (который должен быть неизменяемым), и для поддержания целостности данных может быть определено ограничение ссылочной целостности. Сделать это с помощью типа ROWID не удастся — невозможно создать внешний ключ от дочерней таблицы к родительской по ROWID, равно как и обеспечить в таком случае целостность данных между таблицами. Должно использоваться ограничение первичного ключа.

Так когда же применяется тип ROWID? Он по-прежнему пригоден в приложениях, которые позволяют конечному пользователю взаимодействовать с данными — тип ROWID, будучи физическим адресом строки, предоставляет самый быстрый способ доступа к отдельной строке в любой таблице. Приложение, которое читает данные таблицы и отображает их конечному пользователю, может использовать ROWID при

попытке обновления этой строки. Такое приложение должно применять ROWID в комбинации с другими полями или контрольными суммами (за дополнительной информацией о блокировке приложений обращайтесь в главу 7). В подобной манере можно обновить интересующую строку с наименьшим объемом работы (например, не требуется просмотр индекса для повторного нахождения строки) и гарантировать, что строка осталась в том же виде, в котором она ранее была прочитана, за счет проверки неизменности значений в ее столбцах. Таким образом, тип ROWID удобен в приложениях, эксплуатирующих оптимистическую блокировку.

Резюме

В настоящей главе мы рассмотрели множество базовых типов данных, предлагаемых Oracle; мы увидели, как они физически хранятся и какие опции доступны для каждого типа. Мы начали с символьных строк, наиболее базового типа, и ознакомились с соображениями относительно многобайтных символов и низкоуровневых двоичных данных. Затем мы обсудили расширенные типы данных (доступные в Oracle 12c и последующих версиях) и узнали, что это средство позволяет определять столбцы типов VARCHAR2, NVARCHAR2 и RAW с размером 32 727 байтов. Далее мы изучили числовые типы, включая очень точный тип NUMBER, и типы с плавающей точкой, предоставляемые Oracle 10g и последующими версиями.

Также мы обсудили соображения по поводу унаследованных типов LONG и LONG RAW, сосредоточив внимание на том, как обойти их существование, поскольку их функциональность далека от той, что предлагают типы LOB. После этого мы рассмотрели типы данных, способные хранить дату и время. Мы раскрыли основы арифметики дат — довольно озадачивающую тему, пока не увидать, что она собой представляет на практике. Наконец, в разделе, посвященном типам DATE и TIMESTAMP, мы взглянули на тип INTERVAL и наилучшие приемы работы с ним.

Самой подробной частью главы с точки зрения физического хранения был раздел, в котором рассматривались типы LOB. Типы LOB часто неправильно понимаются разработчиками и администраторами баз данных, поэтому основное внимание уделялось их физической реализации, а также соображениям относительно производительности и отличиям между способами хранения SECUREFILE и BASICFILE.

Последними мы исследовали типы ROWID и UROWID. По причинам, которые теперь должны быть очевидными, вы не должны использовать эти типы данных для столбцов таблиц, т.к. столбец ROWID перестал быть неизменяемым, и на нем нельзя определять ограничения целостности для отношений “родительский—дочерний”. Если необходимо указывать на другую строку, следует сохранять первичные ключи.

Секционирование

Секционирование (partitioning), впервые представленное в версии Oracle 8.0 — это процесс физического разбиения таблицы или индекса на множество мелких, более управляемых частей. С точки зрения приложения, обращающегося к базе данных, логически существует одна таблица или один индекс, но физически таблица и индекс могут состоять из многих десятков физических секций. Каждая *секция* (partition) является независимым объектом, которым можно манипулировать либо отдельно, либо в качестве части более крупного объекта.

На заметку! Секционирование — дополнительная опция редакции Enterprise базы данных Oracle, предусматривающая отдельную плату. В редакции Standard она не доступна.

В этой главе мы выясним, зачем может понадобиться секционирование. Причины простираются от повышения доступности данных до сокращения нагрузки на администратора базы данных и увеличения производительности в определенных ситуациях. Когда вы достигнете хорошего понимания причин, по которым используется секционирование, будет показано, как секционировать таблицы и соответствующие им индексы. Цель этого материала заключается не в том, чтобы научить вас деталям администрирования секций, а скорее предоставить практическое руководство по реализации приложений, работающих с секциями.

Мы также обсудим тот важный факт, что секционирование таблиц и индексов вовсе не гарантирует ускорения работы базы данных. Мне приходилось встречаться со многими разработчиками и администраторами баз данных, которые полагали, что увеличенная производительность является автоматическим побочным эффектом от секционирования объекта. Секционирование — это всего лишь инструмент, и после разбиения таблицы или индекса на секции вы столкнетесь с одним из трех обстоятельств: приложение, которое использует секционированную таблицу, может функционировать медленнее, быстрее или с прежней скоростью. Более того, если вы просто примените секционирование без понимания, как оно работает и каким образом ваше приложение может воспользоваться им, то с высокой вероятностью негативно повлияете на производительность.

Наконец, мы исследуем очень распространенное применение секций в современном мире: поддержка больших оперативных журналов аудита в OLTP и других действующих системах. Мы посмотрим, как внедрить секционирование и сжатие пространства сегментов, чтобы эффективно хранить оперативный журнал аудита и предоставить возможность помещения в архив старых записей из этого журнала с минимальными трудозатратами.

Обзор секционирования

Секционирование облегчает управление очень большими таблицами и индексами по принципу “разделяй и властвуй”. Оно вводит концепцию *ключа секционирования* (partition key), который используется для разделения данных на основе определенных диапазонов значений, списка специфических значений или значений хеш-функции. Если составить перечень преимуществ секционирования в порядке важности, он будет выглядеть следующим образом.

1. **Повышение доступности данных.** Это характерное свойство применимо ко всем типам систем, будь они по природе OLTP или хранилища данных.
2. **Облегчение администрирования крупных сегментов за счет их устранения из базы данных.** Выполнение административных операций в отношении таблицы размером 100 Гбайт, таких как реорганизация или удаление перемещенных строк либо возвращение “пустого пространства”, оставшегося в таблице после очистки устаревшей информации, будет более обременительным, чем десятикратное выполнение тех же операций на отдельных секциях таблицы по 10 Гбайт каждая. Кроме того, с использованием секционирования мы можем проводить процедуру очистки устаревших данных, вообще не оставляя пустого пространства, что полностью исключает необходимость в реорганизации!
3. **Улучшение производительности некоторых запросов.** Это главным образом полезно в среде крупного хранилища данных, где секционирование можно применять для исключения из процесса анализа значительных диапазонов данных, полностью избегая доступа к ним. Такой прием не особенно подходит для транзакционной системы, поскольку в ней и без того происходит обращение к небольшим объемам данных.
4. **Возможное снижение конкуренции в объемных системах OLTP** за счет разнесения операций модификации по множеству отдельных секций. Если имеется сегмент, подвергающийся высокой конкуренции, то его разбиение на множество сегментов может дать побочный эффект, состоящий в пропорциональном снижении степени этой конкуренции.

Давайте по очереди взглянем на каждое из перечисленных потенциальных преимуществ, связанных с использованием секционирования.

Повышенная доступность

Повышенная доступность происходит от независимости каждой секции. Доступность (или ее нехватка) одиночной секции объекта вовсе не означает, что сам объект является недоступным. Оптимизатору известна имеющаяся схема секционирования, и он соответствующим образом исключит секции, на которые нет ссылок, из плана выполнения запроса. Если в крупном объекте одна секция недоступна, а запрос может исключить ее из рассмотрения, то Oracle успешно обработает запрос.

Чтобы продемонстрировать такую увеличенную доступность, мы подготовим хеш-секционированную таблицу с двумя секциями, каждая из которых находится в отдельном табличном пространстве. Мы создадим таблицу EMP с ключом секционирования по столбцу EMPNO. В данном случае эта структура означает, что для каждой строки, вставленной в таблицу, значение столбца EMPNO хешируется для определе-

ния секции (и, следовательно, табличного пространства), в которую строка будет помещена. Сначала мы создаем два табличных пространства (P1 и P2), а затем секционированную таблицу с двумя секциями (PART_1 и PART_2), по одной секции в каждом табличном пространстве:

```
EOGA@ORA12CR1> create tablespace p1 datafile size 1m autoextend on next 1m;
Tablespace created.
```

Табличное пространство создано.

```
EOGA@ORA12CR1> create tablespace p2 datafile size 1m autoextend on next 1m;
Tablespace created.
```

Табличное пространство создано.

```
EOGA@ORA12CR1> CREATE TABLE emp
2 ( empno int,
3   ename varchar2(20)
4 )
5 PARTITION BY HASH (empno)
6 ( partition part_1 tablespace p1,
7   partition part_2 tablespace p2
8 )
9 /
```

Table created.

Таблица создана.

На заметку! В рассматриваемом примере табличные пространства применяют файлы, управляемые Oracle, с параметром инициализации DB_CREATE_FILE_DEST, установленным в /u01/dbfile/ORA12CR1.

Далее мы вставим в таблицу EMP некоторые данные, после чего, используя имя таблицы с указанием секции, просмотрим содержимое каждой секции:

```
EOGA@ORA12CR1> insert into emp select empno, ename from scott.emp;
14 rows created.
```

14 строк создано.

```
EOGA@ORA12CR1> select * from emp partition(part_1);
```

```
EMPNO ENAME
```

```
-----
7369 SMITH
7499 ALLEN
7654 MARTIN
7698 BLAKE
7782 CLARK
7839 KING
7876 ADAMS
7934 MILLER
```

8 rows selected.

8 строк выбрано.

```
EOGA@ORA12CR1> select * from emp partition(part_2);
```

```
EMPNO ENAME
```

```
-----
7521 WARD
7566 JONES
```

```

7788 SCOTT
7844 TURNER
7900 JAMES
7902 FORD
6 rows selected.
6 строк выбрано.

```

Вы должны заметить, что данные распределены довольно случайным образом. Это является особенностью рассматриваемой ситуации. Применяя хеш-секционирование, мы просим Oracle случайно распределить данные (но надеемся на их равномерное распределение) по множеству секций. Мы не можем управлять тем, в какую именно секцию попадут данные; Oracle решает это на основе значений хеш-ключа. Позже, когда будет рассматриваться секционирование по диапазону и по списку, мы увидим, как управлять тем, какая секция должна принимать данные.

Теперь переведем одно из табличных пространств в отключенный режим (эмулируя, к примеру, аварийный отказ диска), что сделает недоступными данные в находящейся внутри него секции:

```

EODA@ORA12CR1> alter tablespace p1 offline;
Tablespace altered.
Табличное пространство изменено.

```

Далее мы запускаем запрос, затрагивающий все секции, и видим, что он терпит неудачу:

```

EODA@ORA12CR1> select * from emp;
select * from emp
*
ERROR at line 1:
ORA-00376: file 9 cannot be read at this time
ORA-01110: data file 9: '/u01/dbfile/ORA12CR1/datafile/ol_mf_p1_9gck8ndv_.dbf'
ОШИБКА в строке 1:
ORA-00376: файл 9 не может быть прочитан в этот момент
ORA-01110: файл данных 9: '/u01/dbfile/ORA12CR1/datafile/ol_mf_p1_9gck8ndv_.dbf'

```

Однако запрос, который не обращается к отключенному табличному пространству, будет функционировать нормально; Oracle проигнорирует отключенное табличное пространство. Переменная привязки в данном примере используется просто для демонстрации того, что хотя на момент оптимизации запроса Oracle не известно, к каким секциям будет производиться доступ, такое игнорирование может быть осуществлено во время выполнения:

```

EODA@ORA12CR1> variable n number
EODA@ORA12CR1> exec :n := 7844;
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.
EODA@ORA12CR1> select * from emp where empno = :n;

EMPNO ENAME
-----
7844 TURNER

```

Таким образом, когда оптимизатор может исключить секции из плана выполнения запроса, он так и поступает. Это обстоятельство повышает доступность для тех приложений, которые применяют в своих запросах ключ секционирования.

Секции также повышают доступность за счет сокращения времени простоя. Например, если есть таблица размером 100 Гбайт, разбитая на 50 секций по 2 Гбайт, то вы можете восстанавливать ее после ошибок гораздо быстрее. Если повреждена одна из секций 2 Гбайт, то время на ее восстановление определяется тем, сколько нужно для восстановления секции размером 2 Гбайт, а не всей таблицы в 100 Гбайт. Поэтому доступность увеличивается двумя путями.

- Исключение секций оптимизатором означает, что многие пользователи могут вообще не заметить недоступность каких-то данных.
- В случае возникновения ошибки время простоя сокращается, потому что восстановление потребует значительно меньшего объема работ.

Облегчение задач администрирования

Снижение нагрузки по администрированию произрастает из того факта, что выполнение операций в отношении небольших объектов по существу легче, быстрее и менее ресурсоемко, чем выполнение тех же операций над крупным объектом.

Например, пусть в базе данных имеется индекс размером 10 Гбайт. Если такой индекс понадобится перестроить, а он не секционирован, то придется перестраивать все 10 Гбайт как одну единицу работы. Хотя это и правда, что индекс можно перестроить в оперативном режиме, повторное создание полного индекса в 10 Гбайт требует огромного количества ресурсов. Вам нужно будет располагать свободным пространством объемом минимум 10 Гбайт для сохранения старого и нового индексов, временной таблицей журнала транзакций, чтобы записывать изменения, вносимые в базовую таблицу во время перестройки индекса, и т.п. С другой стороны, если сам индекс разбит на десять секций по 1 Гбайт, то вы можете перестроить каждую секцию индекса индивидуально, одну за другой. При этом необходимо иметь всего 10% от объема свободного пространства, которое требовалось ранее. Более того, перестройка отдельной секции индекса выполняется намного быстрее (возможно, в десять раз или около того), поэтому с новым индексом будет объединяться гораздо меньше транзакционных изменений, происходящих во время оперативной перестройки индекса, и т.д.

Кроме того, подумайте, что произойдет в случае отказа системы или программного обеспечения прямо перед завершением перестройки индекса в 10 Гбайт. Пропадет весь объем работ. За счет разбиения задачи на части и секционирования индекса в порции по 1 Гбайт из-за отказа вы потеряете максимум 10% всей работы, требуемой для перестройки индекса.

И последний, но не менее важный момент. Может оказаться так, что необходимо перестроить только 10% всего составного индекса — например, реорганизации должны подвергаться только “более новые” (активные) данные, а все “старые” (относительно статичные) данные остаются незатронутыми.

Еще одним примером может быть ситуация, когда вы обнаруживаете, что 50% строк таблицы являются “перемещенными” строками (за деталями о перемещенных строках обращайтесь в главу 10), и хотите это исправить. Наличие секционированной таблицы облегчает такую задачу. Чтобы устранить проблему с перемещенными строками, обычно требуется перестроить объект — таблицу в данном случае.

При наличии одной таблицы размером 100 Гбайт эту операцию придется выполнять последовательно в одной очень крупной порции с использованием команды

ALTER TABLE MOVE. С другой стороны, если есть 25 секций по 4 Гбайт каждая, то вы можете перестраивать секции по отдельности одну за другой. Более того, если вы делаете это в нерабочие часы и обладаете достаточными ресурсами, то можете даже запускать операторы ALTER TABLE MOVE параллельно в разных сеансах, потенциально сокращая время выполнения всей операции в целом. Практически все, что допустимо делать с несекционированным объектом, можно делать с индивидуальной секцией секционированного объекта. Вы даже можете обнаружить, что перемещенные строки сосредоточены в очень небольшом подмножестве секций, а потому произвести перестройку одной или двух секций вместо всей таблицы.

Рассмотрим краткий пример, демонстрирующий перестройку таблицы с множеством перемещенных строк. Обе таблицы, BIG_TABLE1 и BIG_TABLE2, были созданы из экземпляра таблицы BIG_TABLE, содержащей 10 000 000 строк (сценарий создания BIG_TABLE описан в разделе “Настройка среды” в самом начале книги). BIG_TABLE1 — это обычная несекционированная таблица, тогда как BIG_TABLE2 — хеш-секционированная таблица, состоящая из восьми секций (хеш-секционирование подробно описано в следующем разделе, а пока достаточно знать, что оно довольно равномерно распределяет данные по восьми секциям). В примере создаются два табличных пространства и затем две таблицы в них:

```
EODA@ORA12CR1> create tablespace big1 datafile size 1200m;
Tablespace created.
```

Табличное пространство создано.

```
EODA@ORA12CR1> create tablespace big2 datafile size 1200m;
Tablespace created.
```

Табличное пространство создано.

```
EODA@ORA12CR1> create table big_table1
2 ( ID, OWNER, OBJECT_NAME, SUBOBJECT_NAME,
3   OBJECT_ID, DATA_OBJECT_ID,
4   OBJECT_TYPE, CREATED, LAST_DDL_TIME,
5   TIMESTAMP, STATUS, TEMPORARY,
6   GENERATED, SECONDARY )
7 tablespace big1
8 as
9 select ID, OWNER, OBJECT_NAME, SUBOBJECT_NAME,
10        OBJECT_ID, DATA_OBJECT_ID,
11        OBJECT_TYPE, CREATED, LAST_DDL_TIME,
12        TIMESTAMP, STATUS, TEMPORARY,
13        GENERATED, SECONDARY
14 from big_table;
Table created.
```

Таблица создана.

```
EODA@ORA12CR1> create table big_table2
2 ( ID, OWNER, OBJECT_NAME, SUBOBJECT_NAME,
3   OBJECT_ID, DATA_OBJECT_ID,
4   OBJECT_TYPE, CREATED, LAST_DDL_TIME,
5   TIMESTAMP, STATUS, TEMPORARY,
6   GENERATED, SECONDARY )
7 partition by hash(id)
8 (partition part_1 tablespace big2,
```

```

9  partition part_2 tablespace big2,
10 partition part_3 tablespace big2,
11 partition part_4 tablespace big2,
12 partition part_5 tablespace big2,
13 partition part_6 tablespace big2,
14 partition part_7 tablespace big2,
15 partition part_8 tablespace big2
16 )
17 as
18 select ID, OWNER, OBJECT_NAME, SUBOBJECT_NAME,
19        OBJECT_ID, DATA_OBJECT_ID,
20        OBJECT_TYPE, CREATED, LAST_DDL_TIME,
21        TIMESTAMP, STATUS, TEMPORARY,
22        GENERATED, SECONDARY
23 from big_table;
Table created.
Таблица создана.

```

Теперь каждая из этих таблиц находится в собственном табличном пространстве, так что можно легко запросить словарь данных и просмотреть выделенное и свободное пространство в любом табличном пространстве:

```

EODA@ORA12CR1> select b.tablespace_name,
2      mbytes_alloc,
3      mbytes_free
4  from ( select round(sum(bytes)/1024/1024) mbytes_free,
5          tablespace_name
6        from dba_free_space
7        group by tablespace_name ) a,
8      ( select round(sum(bytes)/1024/1024) mbytes_alloc,
9          tablespace_name
10       from dba_data_files
11       group by tablespace_name ) b
12 where a.tablespace_name (+) = b.tablespace_name
13    and b.tablespace_name in ('BIG1','BIG2')
14 /

```

| TABLESPACE_NAME | MBYTES_ALLOC | MBYTES_FREE |
|-----------------|--------------|-------------|
| BIG2 | 1200 | 175 |
| BIG1 | 1200 | 223 |

Табличные пространства BIG1 и BIG2 имеют размеры 1200 Мбайт, и в каждом из них свободно около 200 Мбайт. Попробуем перестроить первую таблицу, BIG_TABLE1:

```

EODA@ORA12CR1> alter table big_table1 move;
alter table big_table1 move
*
ERROR at line 1:
ORA-01652: unable to extend temp segment by 1024 in tablespace BIG1
ОШИБКА в строке 1:
ORA-01652: не удается расширить временный сегмент на 1024 в табличном
пространстве BIG1

```


Операция не удалась. В табличном пространстве BIG1 должно быть достаточно свободного пространства, чтобы хватило для хранения полной копии BIG_TABLE1 одновременно с ее старой копией. Короче говоря, на короткий период требуется вдвое больше места (может больше, а может меньше — в зависимости от окончательного размера перестроенной таблицы). Теперь попробуем выполнить ту же операцию в отношении таблицы BIG_TABLE2:

```
EODA@ORA12CR1> alter table big_table2 move;
alter table big_table2 move
*
```

ERROR at line 1:

ORA-14511: cannot perform operation on a partitioned object

ОШИБКА в строке 1:

ORA-14511: не удастся выполнить операцию над секционированным объектом

Здесь Oracle сообщает нам о невозможности выполнения операции MOVE над таблицей; эта операция должна выполняться над каждой *секцией* таблицы. Мы можем переместить (следовательно, перестроить и реорганизовать) каждую секцию по очереди:

```
EODA@ORA12CR1> alter table big_table2 move partition part_1;
Table altered.
Таблица изменена.
```

```
EODA@ORA12CR1> alter table big_table2 move partition part_2;
Table altered.
Таблица изменена.
```

```
EODA@ORA12CR1> alter table big_table2 move partition part_3;
Table altered.
Таблица изменена.
```

```
EODA@ORA12CR1> alter table big_table2 move partition part_4;
Table altered.
Таблица изменена.
```

```
EODA@ORA12CR1> alter table big_table2 move partition part_5;
Table altered.
Таблица изменена.
```

```
EODA@ORA12CR1> alter table big_table2 move partition part_6;
Table altered.
Таблица изменена.
```

```
EODA@ORA12CR1> alter table big_table2 move partition part_7;
Table altered.
Таблица изменена.
```

```
EODA@ORA12CR1> alter table big_table2 move partition part_8;
Table altered.
Таблица изменена.
```

Каждому индивидуальному перемещению необходимо свободное пространство, достаточное для хранения копии всего лишь одной восьмой от полного объема данных таблицы! Таким образом, эти команды выполняются успешно при том же объеме свободного пространства, которое было ранее. Требуется существенно меньше временных ресурсов и, более того, если в системе происходит отказ (например, из-за

перебоев электропитания) после перемещения секции PART_4, но перед завершением перемещения секции PART_5, то работа, выполненная до последнего оператора MOVE, не будет утеряна. Когда системы восстановится, первые четыре секции по-прежнему будут перемещенными, и можно продолжать обработку секции PART_5.

Кто-то может посмотреть на код и сказать: “Ого, восемь операторов — слишком много набирать на клавиатуре”, и это правда, особенно если в наличии имеется сто и более секций. К счастью, для этого решения очень легко написать сценарий, который может выглядеть так:

```
EODA@ORA12CR1> begin
  2   for x in ( select partition_name
  3               from user_tab_partitions
  4               where table_name = 'BIG_TABLE2' )
  5   loop
  6       execute immediate
  7       'alter table big_table2 move partition ' ||
  8       x.partition_name;
  9   end loop;
10 end;
11 /
```

PL/SQL procedure successfully completed.

Процедура PL/SQL успешно завершена.

Вся необходимая информация находится в словаре данных Oracle, и на большинстве площадок, где было реализовано секционирование, также имеются наборы хранимых процедур, применение которых облегчает управление большим количеством секций. Вдобавок многие инструменты с графическим пользовательским интерфейсом вроде диспетчера предприятия (Enterprise Manager) также располагают встроенными возможностями выполнения операций подобного рода без необходимости во вводе отдельных команд.

Другим фактором, который следует принять во внимание относительно секций и администрирования, является использование *скользящих окон* данных при их долговременном хранении и архивации. Во многих случаях необходимо предоставлять в оперативном режиме данные, охватывающие N единиц времени. Например, предположим, что нужно обеспечивать доступ в оперативном режиме к данным за последние 12 месяцев или прошедшие 5 лет. Без секций это обычно подразумевает выполнение массовых транзакций вставки, за которыми следуют массовые транзакции удаления. В итоге получается большое количество операций DML и огромный объем сгенерированной информации повторения и отмены. При наличии секций можно просто предпринять следующие действия.

1. Загрузить в отдельную таблицу новые данные для месяцев (лет или любого другого периода).
2. Полностью проиндексировать таблицу. (Указанные шаги могут быть выполнены даже в другом экземпляре и перемещены в эту базу.)
3. Присоединить эту вновь загруженную и проиндексированную таблицу в *конец* секционированной таблицы с применением быстрой DDL-команды ALTER TABLE EXCHANGE PARTITION.
4. Отсоединить самую старую секцию в другом конце секционированной таблицы.

Итак, теперь можно очень легко поддерживать исключительно большие объекты, содержащие чувствительную ко времени информацию. Старые данные могут быть удалены из секционированной таблицы и просто *отброшены*, если они не нужны, либо же помещены в архив где-то в другом месте. Новые данные могут быть загружены в отдельную таблицу, не затрагивая при этом секционированную таблицу до конца загрузки, индексации и т.д. Позже в этой главе мы рассмотрим полный пример скользящего окна.

Короче говоря, секционирование позволяет сделать то, что иначе было бы пугающими, а в некоторых случаях и неосуществимыми операциями, так же легко, как в небольшой базе данных.

Улучшенная производительность операторов

Третье общее (потенциальное) преимущество секционирования связано с улучшением производительности операторов (SELECT, INSERT, UPDATE, DELETE, MERGE). Мы рассмотрим два класса операторов — те, что модифицируют информацию, и те, что только читают информацию — и обсудим, каких преимуществ можно ожидать от секционирования в каждом случае.

Параллельный DML

Операторы, которые модифицируют данные в базе, могут обладать потенциалом выполнения в виде *параллельного DML* (parallel DML — PDML). Во время PDML база данных Oracle использует множество потоков или процессов для выполнения операторов INSERT, UPDATE, DELETE или MERGE вместо единственного последовательного процесса. На многопроцессорной машине с достаточной пропускной способностью ввода-вывода потенциальное увеличение скорости может быть большим для массовых операций DML. В выпусках, предшествовавших Oracle9i, для PDML требовалось секционирование. Если таблицы не были секционированы, то выполнять эти операции параллельно было невозможно. Если же таблицы были секционированы, то объекту назначалась *максимальная* степень параллелизма на основе количества имеющихся в нем секций. Это ограничение большей частью было ослаблено в Oracle9i и последующих версиях, но с двумя заметными исключениями. Если таблица, в отношении которой нужно выполнить PDML, имеет битовый индекс или столбец LOB, то она должна быть секционирована для того, чтобы операции PDML могли выполняться параллельно, и степень параллелизма будет ограничена количеством секций.

Вдобавок следует также отметить, что в версии Oracle 12c появилась возможность выполнять операции PDML со столбцами LOB вида SECUREFILE без необходимости в секционировании. В общем, для применения PDML секционирование больше не требуется.

На заметку! Параллельные операции подробно рассматриваются в главе 14.

Производительность запросов

В области производительности запросов, выполняющих строго чтение (операторы SELECT), секционирование вступает в игру с двумя типами специализированных операций.

- Исключение секций. Некоторые секции данных при обработке запроса не принимаются во внимание. Ранее вы уже видели пример исключения секций.
- Параллельные операции. Примерами могут служить параллельное полное сканирование таблиц и параллельное сканирование диапазонов индекса.

Тем не менее, получаемый выигрыш существенно зависит от типа используемой системы.

Системы OLTP

Вы не должны рассматривать будущие секции как способ значительного улучшения производительности запросов в системе OLTP. На самом деле в традиционной системе OLTP секционирование необходимо применять с осторожностью, чтобы не оказать *негативного* влияния на производительность во время выполнения. В традиционной системе OLTP ожидается, что большинство запросов возвращают управление практически мгновенно, и большинство операций извлечения из базы данных происходят через сканирование очень малых диапазонов индекса. Следовательно, главные выигрыши в производительности от секционирования, перечисленные выше, не вступят в силу. Исключение секций удобно, когда присутствуют полные сканирования крупных объектов, т.к. оно позволяет избежать полного сканирования значительных порций этих объектов. Однако в среде OLTP *полное сканирование крупных объектов не осуществляется* (в противном случае это говорит о наличии серьезной проблемы в проектном решении). Даже когда вы секционируете индексы, любое увеличение производительности, достигаемое сканированием индексов меньшего размера, будет крошечным, если вообще будет. Если некоторые запросы используют индекс и не могут исключить из рассмотрения все секции кроме одной, то вы можете даже обнаружить, что после секционирования запросы выполняются фактически медленнее, поскольку теперь приходится зондировать 5, 10 или 20 мелких индексов вместо одного большого. Мы исследуем это более подробно, когда приступим к рассмотрению доступных типов секционированных индексов.

Что касается параллельных операций, которые обсуждаются в следующей главе, то в системе OLTP выполнять параллельные запросы не нужно. Вы можете оставить применение параллельных операций администратору базы данных для выполнения перестроек, создания индексов, сбора статистики и т.д. Факт состоит в том, что в системе OLTP запросы уже должны отличаться очень быстрым индексным доступом и секционирование не может заметно его ускорить, если рост скорости вообще произойдет. Это вовсе не означает, что вы должны избегать секционирования в системах OLTP; просто не стоит ожидать от него значительного улучшения производительности. Большинство приложений OLTP не способны извлекать выгоду в ситуациях, когда секционирование содействует росту производительности запросов, но вы все же получаете другие преимущества: облегчение администрирования, повышение доступности и сокращение конкуренции.

Системы хранилищ данных

В системе хранилища данных или в системе поддержки принятия решений секционирование является не только замечательным инструментом администрирования, но также и средством ускорения обработки. Например, у вас может быть крупная таблица, в которой необходимо выполнить нерегламентированный запрос.

Вы всегда запускаете этот нерегламентированный запрос ежеквартально, т.к. с каждым кварталом продаж связаны сотни тысяч записей, а всего есть миллионы оперативных записей. Итак, вы хотите выполнить запрос по относительно небольшому срезу всего набора данных, но индексация данных на основе квартала продаж в действительности неосуществима. Такой индекс указывал бы на сотни тысяч записей, а сканирование по диапазону индекса оказалось бы гигантским (более подробно об этом рассказывалось в главе 11). В результате для обработки многих ваших запросов вызывается полное сканирование таблицы, но в конечном итоге приходится сканировать миллионы записей, большинство из которых к запросу не имеют отношения. Используя продуманную схему секционирования, вы можете разделить данные по кварталам, так что при запросе данных за определенный квартал будет выполняться полное сканирование только данных этого квартала. Из всех возможных решений это будет наилучшим.

Кроме того, в среде системы хранилища данных или поддержки принятия решений часто применяются параллельные запросы. Здесь такие операции, как параллельное сканирование диапазона индекса или параллельное быстрое сканирование полного индекса, не только важны, но и полезны. Мы хотим максимизировать использование всех доступных ресурсов, и параллельный запрос является методом достижения этой цели. Значит, в такой среде секционирование имеет хороший шанс ускорить обработку.

Сокращение конкуренции в системе OLTP

Последняя общая область, в которой секционирование обеспечивает преимущества — это потенциальное увеличение степени параллелизма за счет снижения конкуренции в системе OLTP. Секционирование может применяться для распространения операций модификации одиночной таблицы по множеству физических секций. Идея состоит в том, что если есть сегмент, испытывающий высокую конкуренцию, то его превращение в набор из нескольких сегментов может дать побочный эффект в виде пропорционального сокращения такой конкуренции.

Например, вместо единственного табличного сегмента с одиночным индексным сегментом вы могли бы иметь 20 секций таблицы и 20 секций индекса. Это было бы похоже на ситуацию с наличием 20 таблиц вместо одной (и 20 индексов вместо одного), так что во время выполнения операций модификации конкуренция за такой разделяемый ресурс могла бы снизиться.

Схемы секционирования таблиц

В настоящее время доступны девять методов, с помощью которых можно секционировать таблицы в Oracle.

- **Секционирование по диапазонам ключей (range partitioning).** Вы можете указать диапазоны данных, которые должны храниться вместе. Например, все, что имеет отметку времени в пределах января 2014 года, будет сохранено в секции 1, все, что имеет метку в пределах февраля 2014 года — в секции 2 и т.д. Вероятно, это наиболее распространенный механизм секционирования в Oracle.
- **Хеш-секционирование (hash partitioning).** Такой метод секционирования демонстрировался в первом примере настоящей главы. К столбцу (или столб-

цам) применяется хеш-функция, и строка будет помещена в секцию согласно значению, возвращаемому этой хеш-функцией.

- **Секционирование по списку значений ключа (list partitioning).** Вы указываете дискретный набор значений, определяющий данные, которые должны храниться вместе. Например, можно указать, что строки со значением в столбце STATUS из списка ('A', 'M', 'Z') попадают в секцию 1, строки со значением STATUS из списка ('D', 'P', 'Q') — в секцию 2 и т.д.
- **Секционирование по интервалам ключей (interval partitioning).** Очень похоже на секционирование по диапазонам ключей, но с тем отличием, что база данных сама может создавать новые секции по мере поступления данных. При традиционном секционировании по диапазонам ключей администратор базы данных обязан предварительно создать секции для хранения каждого возможного значения данных, как в текущий момент, так и в будущем. Обычно это означает, что администратор базы данных создает секции по расписанию, готовя их под размещение данных для следующих месяцев или следующих недель. При секционировании по интервалам ключей база данных будет самостоятельно создавать секции, когда появятся новые данные, которые не умещаются в существующие секции, на основе правила, указанного администратором базы данных.
- **Секционирование по ссылкам (reference partitioning).** Позволяет дочерней таблице в рамках отношения “родительская—дочерняя”, обеспечиваемом внешним ключом, наследовать схему секционирования родительской таблицы. Это делает возможным эквисекционирование дочерней таблицы по ее родительской таблице, без необходимости в денормализации модели данных. В прошлом таблица могла секционироваться только на основе физически хранимых атрибутов; секционирование по ссылкам по существу позволяет секционировать таблицу на основе атрибутов из родительской таблицы.
- **Секционирование по интервалам ключей и по ссылкам (interval reference partitioning).** Как должно быть понятно по названию, это комбинация секционирования по интервалам ключей и секционирования по ссылкам. Такой тип секционирования доступен, начиная с версии Oracle 12c. Он делает возможным автоматическое добавление секций к родительским и дочерним таблицам, секционированным по ссылкам.
- **Секционирование по виртуальному столбцу (virtual column partitioning).** Позволяет секционировать по выражению, основанному на одном или большем количестве существующих столбцов таблицы. Выражение хранится только как метаданные.
- **Составное секционирование (composite partitioning).** Это комбинация секционирования по диапазонам ключей, хеш-секционирования и секционирования по списку значений ключа. Оно позволяет сначала применить одну схему секционирования к некоторым данным, а затем дополнительно разбить каждую результирующую секцию на подсекции с использованием какой-то другой схемы секционирования.
- **Системное секционирование (system partitioning).** Приложение определяет, в какую секцию будет явно вставляться строка. Такой тип секционирования

имеет ограниченную область применения и в этой главе не рассматривается; он упоминается здесь для предоставления полного списка типов секционирования, поддерживаемых в Oracle. Подробные сведения о системном секционировании можно найти в руководстве разработчика по картриджам данных Oracle (Oracle Database Data Cartridge Developer's Guide).

В последующих разделах мы взглянем на преимущества каждого типа секционирования и на отличия между ними. Мы также обсудим, когда какую схему применять в приложениях разнообразных типов. Этот раздел не задумывался как исчерпывающая демонстрация синтаксиса секционирования и всех доступных опций. Напротив, примеры просты и иллюстративны; они предназначены для того, чтобы предоставить обзор работы секционирования и особенностей различных его типов.

На заметку! Полное описание синтаксиса секционирования можно найти в руководстве по языку SQL для Oracle (*Oracle Database SQL Language Reference*) или в руководстве для администратора баз данных Oracle (*Oracle Database Administrator's Guide*). Кроме того, руководство по очень крупным базам данных и секционированию Oracle (*Oracle Database VLDB and Partitioning Guide*) и руководство по хранилищам данных (*Oracle Database Data Warehousing Guide*) являются великолепными источниками информации по вариантам секционирования и обязательны для изучения каждым, кто планирует внедрять секционирование.

Секционирование по диапазонам ключей

Первой мы рассмотрим таблицу, секционированную по диапазонам ключей. Приведенный ниже оператор CREATE TABLE создает секционированную по диапазонам таблицу с использованием столбца RANGE_KEY_COLUMN. Все данные со значением в столбце RANGE_KEY_COLUMN, *строго меньшим, чем* 01-JAN-2014, будут помещены в секцию PART_1, а все данные со значением в столбце RANGE_KEY_COLUMN, которое *строго меньше, чем* 01-JAN-2015 (и больше или равно 01-JAN-2014), попадут в секцию PART_2. Любые данные, которые не удовлетворяют ни одному из этих условий (например, строка, имеющая в столбце RANGE_KEY_COLUMN значение 01-JAN-2015 или больше), при вставке потерпят неудачу, т.к. не могут быть сопоставлены ни с одной из секций.

```
EODA@ORA12CR1> CREATE TABLE range_example
2 ( range_key_column date NOT NULL,
3   data                varchar2(20)
4 )
5 PARTITION BY RANGE (range_key_column)
6 ( PARTITION part_1 VALUES LESS THAN
7   (to_date('01/01/2014','dd/mm/yyyy')),
8   PARTITION part_2 VALUES LESS THAN
9   (to_date('01/01/2015','dd/mm/yyyy'))
10 )
11 /
```

Table created.

Таблица создана.

На заметку! В операторе CREATE TABLE мы применяем формат даты DD/MM/YYYY, чтобы сделать ее интернациональной. Если использовать формат DD-MON-YYYY, то CREATE TABLE откажется работать и выдаст сообщение ORA-01843: not a valid month (ORA-01843: недопустимый месяц), если аббревиатура января в системе отличается от JAN. На это влияет установка NLS_LANGUAGE. Однако в тексте и операторах INSERT применялась трехсимвольная аббревиатура месяца, чтобы избежать неоднозначности относительно того, какой компонент обозначает день, а какой — месяц.

На рис. 13.1 показано, что Oracle будет инспектировать значение RANGE_KEY_COLUMN и на его основе вставлять строку в одну из двух секций.

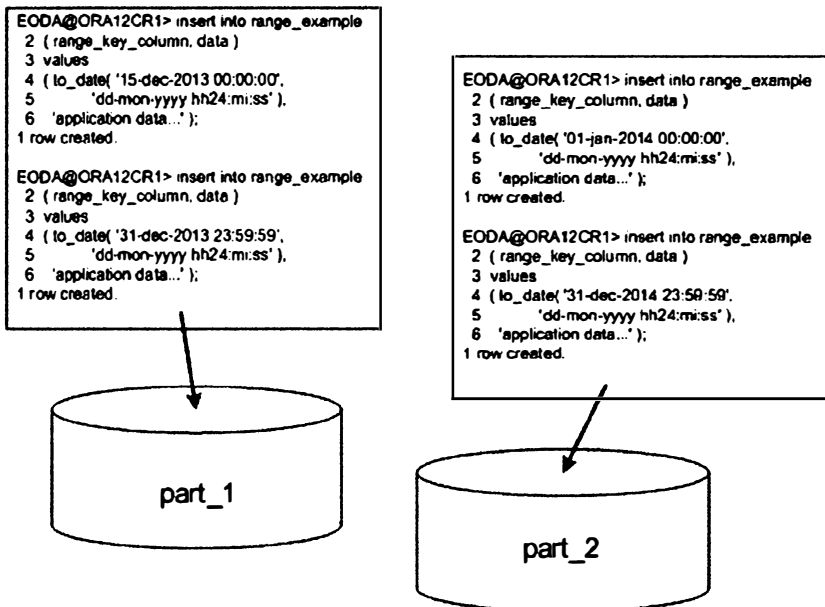


Рис. 13.1. Пример вставки при секционировании по диапазонам

Вставляемые строки были выбраны специально для демонстрации того, что диапазон секций работает по условию “строгое меньше”, а не “меньше или равно”. Сначала мы вставляем значение 15-DEC-2013, которое определенно попадет в секцию PART_1. Мы также вставляем строку с датой/временем на секунду раньше 01-JAN-2014 — эта строка также последует в секцию PART_1, поскольку значение ее столбца RANGE_KEY_COLUMN *меньше* 01-JAN-2014. Однако следующая вставка строки со значением полуночи 01-JAN-2014 отправится в секцию PART_2, т.к. эта дата/время не строго меньше, чем граница диапазона секции для PART_1. Последняя строка очевидным образом относится к секции PART_2, потому что значение ее столбца RANGE_KEY_COLUMN больше или равно границе диапазона секции для PART_1 и меньше границы диапазона секции для PART_2.

Мы можем подтвердить это, выполнив операторы SELECT для индивидуальных секций:


```

EODA@ORA12CR1> select to_char(range_key_column, 'dd-mon-yyyy hh24:mi:ss')
2   from range_example partition (part_1);

TO_CHAR(RANGE_KEY_COLUMN, 'DD-
-----
15-dec-2013 00:00:00
31-dec-2013 23:59:59

EODA@ORA12CR1> select to_char(range_key_column, 'dd-mon-yyyy hh24:mi:ss')
2   from range_example partition (part_2);

TO_CHAR(RANGE_KEY_COLUMN, 'DD-
-----
01-jan-2014 00:00:00
31-dec-2014 23:59:59

```

Вас может интересовать, что произойдет в случае вставки даты, которая выходит за пределы верхней границы. В такой ситуации Oracle сообщит об ошибке:

```

EODA@ORA12CR1> insert into range_example
2   ( range_key_column, data )
3   values
4   ( to_date( '01-jan-2015 00:00:00',
5             'dd-mon-yyyy hh24:mi:ss' ),
6     'application data...' );
insert into range_example
*
ERROR at line 1:
ORA-14400: inserted partition key does not map to any partition
ОШИБКА в строке 1:
ORA-14400: вставляемый ключ секционирования не соответствует ни одной секции

```

Здесь возможны два подхода. Один подход предусматривает использование секционирования по интервалам, которое описано далее в главе. Другой подход предполагает создание секции, перехватывающей все; это мы сейчас и продемонстрируем. Предположим, что вы хотите разносить даты 2013 и 2014 годов по разным секциям, а все остальные даты отправлять в третью секцию. При секционировании по диапазонам это можно сделать с помощью конструкции MAXVALUE, которая выглядит следующим образом:

```

EODA@ORA12CR1> CREATE TABLE range_example
2   ( range_key_column date ,
3     data              varchar2(20)
4   )
5   PARTITION BY RANGE (range_key_column)
6   ( PARTITION part_1 VALUES LESS THAN
7     (to_date('01/01/2014', 'dd/mm/yyyy')),
8     PARTITION part_2 VALUES LESS THAN
9     (to_date('01/01/2015', 'dd/mm/yyyy')),
10    PARTITION part_3 VALUES LESS THAN
11      (MAXVALUE)
12  )
13  /
Table created.
Таблица создана.

```

Теперь вставляемая в таблицу строка попадет в одну из трех секций. Никакие строки не отклоняются, потому что секция PART_3 может принимать строки со значениями RANGE_KEY_COLUMN, которые не попадают в секцию PART_1 или PART_2 (в эту новую секцию будут вставляться даже строки с NULL-значением в столбце RANGE_KEY_COLUMN).

Хеш-секционирование

Когда таблица хеш-секционирована, Oracle будет применять хеш-функцию к ключу секционирования для определения того, в какую из N секций должны быть помещены данные. В Oracle рекомендуют в качестве значения N использовать степень двойки (2, 4, 8, 16 и т.д.), чтобы достичь наилучшего общего распределения, и вскоре вы увидите, что этот совет, безусловно, полезен.

Как работает хеш-секционирование

Хеш-секционирование предназначено для эффективного распределения данных по многим разным устройствам (дискам) или просто для разделения данных на поддающиеся управлению порции. Выбранный для таблицы хеш-ключ должен быть столбцом или набором столбцов, который дает уникальное значение или, по крайней мере, имеет столько отличающихся значений, сколько возможно для обеспечения эффективного распределения строк по секциям. Если вы выберете столбец, имеющий только четыре значения, и организуете две секции, то все строки довольно легко могут попасть *в одну и ту же секцию*, полностью сводя на нет саму цель секционирования!

Давайте создадим хеш-таблицу с двумя секциями и выберем в качестве ключа секционирования столбец по имени HASH_KEY_COLUMN. Хешируя значение в этом столбце, Oracle будет определять секцию, в которой должна сохраняться строка:

```

ODA@ORA12CR1> CREATE TABLE hash_example
2 ( hash_key_column  date,
3   data              varchar2(20)
4 )
5 PARTITION BY HASH (hash_key_column)
6 ( partition part_1 tablespace p1,
7   partition part_2 tablespace p2
8 )
9 /
Table created.
Таблица создана.

```

На рис. 13.2 показано, что Oracle будет инспектировать значение в столбце HASH_KEY_COLUMN, хешировать его и выяснять, в какую из двух секций должна попасть эта строка.

Как упоминалось ранее, хеш-секционирование не предоставляет никакого контроля над тем, в какой именно секции окажется та или иная строка. К значению ключа секционирования Oracle применяет хеш-функцию и определяет, куда отправится строка. Если по какой-то причине вы хотите поместить конкретную строку в секцию PART_1, то *не должны и фактически не сможете* использовать хеш-секционирование. Строка попадает в ту секцию, на которую указывает значение, возвращенное из хеш-функции.

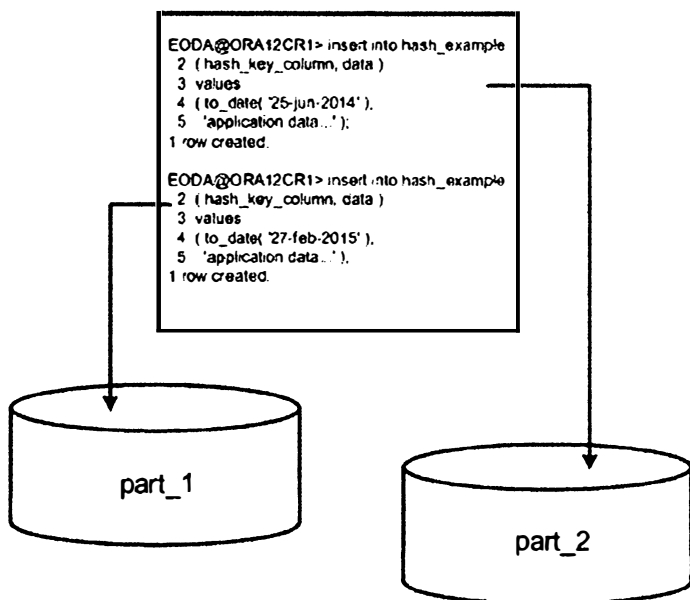


Рис. 13.2. Пример вставки в хеш-секцию

Если вы измените количество хеш-секций, то данные перераспределятся по всем секциям (добавление или удаление секции в таблице с хеш-секционированием вызывает перезапись всех данных, так что каждая строка может оказаться в другой секции).

Хеш-секционирование наиболее удобно, когда есть крупная таблица вроде той, что рассматривалась в разделе “Облегчение задач администрирования”, и вы хотели бы *разделить и властвовать* над ней. Вместо того чтобы управлять одной огромной таблицей, вы желали бы иметь дело с 8 или 16 небольшими таблицами. Хеш-секционирование также полезно для увеличения доступности до определенного уровня, как было продемонстрировано в разделе “Повышенная доступность” ранее в главе; временная потеря одной хеш-секции не мешает доступу ко всем остальным секциям. Некоторых пользователей это может затронуть, но с высокой вероятностью многие пользователи ничего не заметят. Кроме того, единица восстановления становится намного меньше. Вам не придется восстанавливать одну крупную таблицу, а только ее часть. И, наконец, хеш-секционирование удобно в средах с высокой конкуренцией за обновления, как упоминалось в разделе “Сокращение конкуренции в системе OLTP”. Вместо одного горячего сегмента посредством хеш-секционирования можно получить 16 порций, каждая из которых будет принимать модификации.

Использование степени двойки при хеш-секционировании

Ранее я указывал, что количество секций должно выражаться числом, представляющим степень двойки. В истинности этого правила легко убедиться. В целях демонстрации мы подготовим хранимую процедуру для автоматизации создания хеш-секционированной таблицы с N секциями (N послужит параметром). Эта процедура будет конструировать динамический запрос, который извлекает счетчики строк по

секциям и затем отображает их вместе с простой гистограммой. Наконец, процедура откроет запрос и позволит просмотреть результаты. Она начинается с создания хеш-таблицы. Мы будем применять таблицу по имени T.

```

EODA@ORA12CR1> create or replace
2  procedure hash_proc
3      ( p_nhash in number,
4        p_cursor out sys_refcursor )
5  authid current_user
6  as
7      l_text      long;
8      l_template long :=
9          'select $POS$ oc, 'p$POS$' pname, count(*) cnt ' ||
10         'from t partition ( $PNAME$ ) union all ';
11      table_or_view_does_not_exist exception;
12      pragma exception_init( table_or_view_does_not_exist, -942 );
13  begin
14      begin
15          execute immediate 'drop table t';
16      exception when table_or_view_does_not_exist
17          then null;
18      end;
19
20      execute immediate '
21      CREATE TABLE t ( id )
22      partition by hash(id)
23      partitions ' || p_nhash || '
24      as
25      select rownum
26      from all_objects';

```

Затем мы динамически сконструируем запрос для извлечения счетчика строк по секции. Он будет делать это с использованием определенного ранее шаблонного запроса. Для каждой секции мы получим счетчик с указанием имени таблицы, расширенного именем секции, и объединим все счетчики вместе:

```

28      for x in ( select partition_name pname,
29                  PARTITION_POSITION pos
30                  from user_tab_partitions
31                  where table_name = 'T'
32                  order by partition_position )
33  loop
34      l_text := l_text ||
35          replace(
36              replace(l_template,
37                  '$POS$', x.pos),
38                  '$PNAME$', x.pname );
39  end loop;

```

А теперь мы возьмем этот запрос и выберем позицию секции (PNAME) и счетчик строк в этой секции (CNT). С помощью функции RPAD мы построим примитивную, но наглядную гистограмму:

```

41 open p_cursor for
42   'select pname, cnt,
43     substr( rpad(' '*', 30*round( cnt/max(cnt)over(), 2), ' '*'), 1, 30) hg
44     from (' || substr( l_text, 1, length(l_text)-11 ) || ' )
45     order by oc';
46
47 end;
48 /

```

Procedure created.
Процедура создана.

Запустив эту процедуру с входным значением 4, т.е. для четырех хеш-секций, мы можем получить примерно такой вывод:

```

EODA@ORA12CR1> variable x refcursor
EODA@ORA12CR1> set autoprint on
EODA@ORA12CR1> exec hash_proc( 4, :x );
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.

```

| PN | CNT | HG |
|----|-------|-------|
| p1 | 12141 | ***** |
| p2 | 12178 | ***** |
| p3 | 12417 | ***** |
| p4 | 12105 | ***** |

Эта простая гистограмма демонстрирует хорошее равномерное распределение данных по всем четырем секциям. Секции содержат почти одинаковые количества строк. Тем не менее, если мы просто перейдем с четырех на пять секций, то увидим следующую картину:

```

EODA@ORA12CR1> exec hash_proc( 5, :x );
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.

```

| PN | CNT | HG |
|----|-------|-------|
| p1 | 6102 | ***** |
| p2 | 12180 | ***** |
| p3 | 12419 | ***** |
| p4 | 12106 | ***** |
| p5 | 6040 | ***** |

Эта гистограмма показывает, что первая и последняя секции имеют в половину меньше строк, чем остальные. Данные распределены не очень равномерно. Такая тенденция продолжает наблюдаться для шести и семи хеш-секций:

```

EODA@ORA12CR1> exec hash_proc( 6, :x );
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.

```

| PN | CNT | HG |
|----|-------|-------|
| p1 | 6104 | ***** |
| p2 | 6175 | ***** |
| p3 | 12420 | ***** |

```
p4      12106 *****
p5      6040 *****
p6      6009 *****
6 rows selected.
6 строк выбрано.
```

```
EODA@ORA12CR1> exec hash_proc( 7, :x );
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.
```

| PN | CNT | HG |
|----|-------|-------|
| p1 | 6105 | ***** |
| p2 | 6176 | ***** |
| p3 | 6161 | ***** |
| p4 | 12106 | ***** |
| p5 | 6041 | ***** |
| p6 | 6010 | ***** |
| p7 | 6263 | ***** |

```
7 rows selected.
7 строк выбрано.
```

Как только мы возвратимся к количеству хеш-секций, равному степени двойки, мы снова достигнем цели равномерного распределения:

```
EODA@ORA12CR1> exec hash_proc( 8, :x );
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.
```

| PN | CNT | HG |
|----|------|-------|
| p1 | 6106 | ***** |
| p2 | 6178 | ***** |
| p3 | 6163 | ***** |
| p4 | 6019 | ***** |
| p5 | 6042 | ***** |
| p6 | 6010 | ***** |
| p7 | 6264 | ***** |
| p8 | 6089 | ***** |

```
8 rows selected.
8 строк выбрано.
```

Продолжив экспериментирование вплоть до 16 секций, мы увидим те же самые результаты для секций с девятой по пятнадцатую, т.е. смещение данных из краевых секций во внутренние, а по достижении шестнадцатой секции распределение снова выравнивается. Аналогично обстоят дела до секции номер 32, 64 и т.д. Этот пример просто подчеркивает важность применения для количества хеш-секций значения, которое является степенью двойки.

Секционирование по списку значений ключа

Секционирование по списку значений ключа появилось в выпуске Oracle9i Release 1. Оно дает возможность указать, в какой секции будет находиться строка, на основе дискретных списков значений. Часто удобно проводить секционирование по какому-то коду наподобие кода штата или региона.

Например, может понадобиться поместить в одну секцию все записи о людях из штатов Мэн (ME), Нью-Гемпшир (NH), Вермонт (VT) и Массачусетс (MA), т.к. эти штаты расположены рядом друг с другом, а приложение запрашивает данные по географическому региону. Аналогично может потребоваться сгруппировать людей из штатов Коннектикут (CT), Род-Айленд (RI) и Нью-Йорк (NY).

Использовать здесь секционирование по диапазонам не получится, поскольку диапазон для первой секции простирался бы от ME до VT, а для второй — от CT до RI. Эти диапазоны перекрываются (в смысле алфавита). Также нельзя применять хеш-секционирование, потому что в нем невозможно управлять тем, в какую секцию попадет любая взятая строка; этим занимается встроенная хеш-функция, предоставляемая Oracle. Секционирование по списку позволяет легко реализовать следующую специальную схему:

```

EODA@ORA12CR1> create table list_example
2 ( state_cd  varchar2(2),
3   data      varchar2(20)
4 )
5 partition by list(state_cd)
6 ( partition part_1 values ( 'ME', 'NH', 'VT', 'MA' ),
7   partition part_2 values ( 'CT', 'RI', 'NY' )
8 )
9 /

```

Table created.

Таблица создана.

На рис. 13.3 показано, что Oracle будет инспектировать столбец STATE_CD и на основе его значения помещать строку в подходящую секцию.

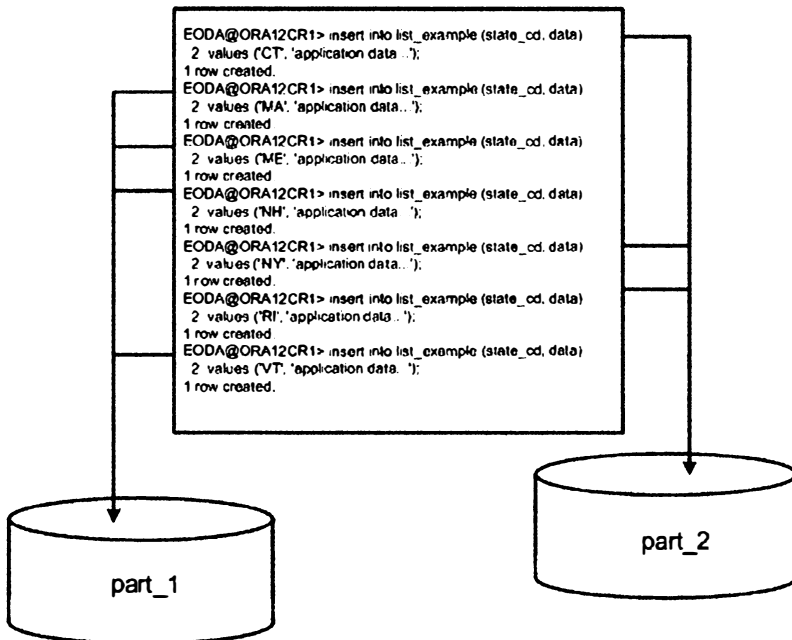


Рис. 13.3. Пример вставки при секционировании по списку

Как и при секционировании по диапазонам, если мы попытаемся вставить значение, которое не указано в списке секций, то Oracle возвратит клиентскому приложению сообщение об ошибке. Другими словами, таблица с секционированием по списку без секции DEFAULT будет неявно налагать ограничение, очень похожее на проверочное ограничение целостности:

```
EOADA@ORA12CR1> insert into list_example values ( 'VA', 'data' );
insert into list_example values ( 'VA', 'data' )
```

*

ERROR at line 1:

ORA-14400: inserted partition key does not map to any partition

ОШИБКА в строке 1:

ORA-14400: вставляемый ключ секционирования не соответствует ни одной секции

Если вы захотите разнести эти семь штатов по отдельным секциям, но все остальные штаты (фактически все строки, которые не содержат один из этих семи кодов штатов) поместить в третью секцию, то придется использовать конструкцию VALUES (DEFAULT). Изменим таблицу, добавив третью секцию (ее можно было бы также указать в операторе CREATE TABLE):

```
EOADA@ORA12CR1> alter table list_example
```

```
2 add partition
```

```
3 part_3 values ( DEFAULT );
```

Table altered.

Таблица изменена.

```
EOADA@ORA12CR1> insert into list_example values ( 'VA', 'data' );
```

1 row created.

1 строка создана.

В третью секцию попадут все значения, которые явно не указаны в списках значений. Здесь необходимо сделать одно предупреждение: после того, как секционированная по списку таблица получит секцию DEFAULT, добавление к ней новых секций станет невозможным:

```
EOADA@ORA12CR1> alter table list_example
```

```
2 add partition
```

```
3 part_4 values( 'CA', 'NM' );
```

```
alter table list_example
```

*

ERROR at line 1:

ORA-14323: cannot add partition when DEFAULT partition exists

ОШИБКА в строке 1:

ORA-14323: невозможно добавить секцию, когда существует секция DEFAULT

В этом случае потребуется удалить секцию DEFAULT, затем добавить секцию PART_4 и, наконец, вернуть секцию DEFAULT на место. Причина в том, что секция DEFAULT может содержать строки со значением ключа секционирования CA или NM — после добавления секции PART_4 они уже не смогут там находиться.

Секционирование по интервалам ключей

Секционирование по интервалам ключей является средством, доступным в Oracle 11g Release 1 и последующих версиях. Оно очень похоже на описанное ранее секционирование по диапазонам. По сути, секционирование по интервалам ключей

начинается с таблицы, секционированной по диапазонам, но с дополнением определения правилом (интервалом), которое позволит базе данных знать, каким образом добавлять секции в будущем.

Цель секционирования по интервалам — создавать новые секции для данных тогда и только тогда, когда существуют данные для конкретной секции и когда эти данные поступают в базу данных. Другими словами, устраняется необходимость в предварительном создании секций для данных, а данные получают возможность самостоятельно создавать секцию во время вставки. Чтобы применить секционирование по интервалам, вы начинаете с таблицы, секционированной по диапазонам, но без секции MAXVALUE, и указываете интервал для добавления к *верхней границе* (самому высокому значению в этой секционированной таблице) и создания нового диапазона. Необходимо иметь таблицу, секционированную по диапазону значений единственного столбца, который разрешает добавлять к нему данные типа NUMBER или INTERVAL (например, таблица, секционированная по столбцу типа VARCHAR2, не может быть секционирована по интервалам; нет ничего такого, что можно было бы добавлять к VARCHAR2). Секционирование по интервалам можно использовать с любой подходящей таблицей, секционированной по диапазонам; т.е. с помощью оператора ALTER TABLE превратить существующую таблицу, секционированную по диапазонам, в секционированную по интервалам, или же создать новую таблицу посредством CREATE TABLE.

Например, предположим, что есть секционированная по диапазону таблица, в которой все, что меньше 01-JAN-2015 (данные за 2014 год и ранее), попадает в секцию P1 — и это все. Таким образом, таблица имеет одну секцию для всех данных за 2014 год и предшествующие годы. Если вы попытаетесь вставить в нее данные за 2015 год, то такая вставка потерпит неудачу, как было показано ранее в разделе, посвященном секционированию по диапазонам. Секционирование по интервалам позволяет создать таблицу и указать и диапазон (строго меньше 01-JAN-2015), и интервал (скажем, продолжительностью 1 месяц), после чего база по мере поступления данных будет создавать секции для месяцев (каждая из них будет хранить данные только за соответствующий месяц). База данных не станет заранее создавать все возможные секции, т.к. это было бы непрактично. Но при появлении каждой строки база данных будет проверять, существует ли секция для месяца, к которому относится строка. При необходимости база данных будет создавать новую секцию.

Ниже приведен пример синтаксиса:

```
EODA@ORA12CR1> create table audit_trail
2 ( ts timestamp,
3   data varchar2(30)
4 )
5 partition by range(ts)
6 interval (numtoyminterval(1,'month'))
7 store in (users, example )
8 (partition p0 values less than
9  (to_date('01-01-1900','dd-mm-yyyy'))
10 )
11 /
```

Table created.

Таблица создана.

На заметку! Может возникнуть вопрос, особенно если вы только что завершили чтение предыдущей главы, посвященной типам данных. Вы видели, что было выполнено секционирование по столбцу типа `TIMESTAMP` и добавлен интервал продолжительностью один месяц. В главе о типах данных было показано, что добавление к значению `TIMESTAMP`, соответствующему 31 января, интервала в один месяц вызывает ошибку, поскольку даты 31 февраля не существует. Возникнет ли та же самая проблема с секционированием по интервалам? Да, возникнет. Если вы попытаетесь применить дату вроде `'29-01-1990'` (любого дня месяца после 28-го будет достаточно), то получите ошибку `ORA-14767: cannot specify this interval with existing high bounds (ORA-14767: невозможно указывать этот интервал при существующих верхних границах)`. База данных не разрешит использовать граничное значение, к которому небезопасно добавлять интервал.

В строках 8 и 9 вы видите схему секционирования по диапазонам для этой таблицы; она начинается с одной пустой секции, которая будет содержать любые данные, предшествующие дате `01-JAN-1900`. Предполагается, что поскольку таблица содержит журнал аудита, эта секция навсегда останется небольшой и пустой. Она является обязательной и называется *переходной* секцией. Все данные, которые строго меньше этого текущего верхнего значения секции, будут подвергнуты традиционному секционированию по диапазонам. Секционирование по интервалам будет использоваться только для данных, которые превышают верхнее значение переходной секции. Запросив словарь данных, можно посмотреть, что было создано до сих пор:

```

EODA@ORA12CR1> select a.partition_name, a.tablespace_name, a.high_value,
2      decode( a.interval, 'YES', b.interval ) interval
3      from user_tab_partitions a, user_part_tables b
4      where a.table_name = 'AUDIT_TRAIL'
5      and a.table_name = b.table_name
6      order by a.partition_position;

```

| PARTITION_ | TABLESPACE | HIGH_VALUE | INTERVAL |
|------------|------------|---------------------------------|----------|
| ----- | ----- | ----- | ----- |
| P0 | USERS | TIMESTAMP' 1900-01-01 00:00:00' | |

К этому моменту мы имеем только одну секцию, и это не секция по интервалу, что видно по пустому столбцу `INTERVAL`. Это просто обычная секция по диапазону, которая будет хранить все, что строго меньше `01-JAN-1900`.

Взглянув на оператор `CREATE TABLE` еще раз, мы можем заметить в строках 6 и 7 новую информацию, специфичную для секционирования по интервалам:

```

6 interval (numtoyminterval(1,'month'))
7 store in (users, example )

```

В строке 6 находится действительная спецификация интервала `NUMTOYMINTERVAL(1, 'MONTH')`. Наша цель заключается в поддержке секций для месяцев, т.е. в создании новой секции для каждого месяца, что представляет собой довольно распространенную задачу. За счет применения даты, к которой можно безопасно добавлять месяц (в главе 12 объяснялось, почему добавление месяца к значению `TIMESTAMP` в ряде случаев может быть чревато ошибками) — первого дня месяца — можно заставить базу данных создавать секции для месяцев на лету по мере поступления данных.

В строке 7 мы видим специфическую конструкцию `STORE IN (USERS, EXAMPLE)`. Она позволяет сообщить базе данных, где создавать новые секции, т.е. какое табличное пространство использовать. Когда база данных выясняет, что за секции нужно создать, она применяет этот список для определения, в каком табличном пространстве создавать каждую секцию. В результате администратор базы данных может управлять максимальным желаемым размером табличного пространства: одно табличное пространство размером 500 Гбайт может не подойти, а 10 табличных пространств по 50 Гбайт окажутся вполне удобными. В таком случае администратору базы данных понадобится настроить 10 табличных пространств и позволить базе данных использовать их все для создания секций.

Давайте теперь вставим строку данных и посмотрим, что произойдет:

```
EODA@ORA12CR1> insert into audit_trail (ts,data) values
2 ( to_timestamp('27-feb-2014','dd-mon-yyyy'), 'xx' );
1 row created.
1 строка создана.
```

```
EODA@ORA12CR1> select a.partition_name, a.tablespace_name, a.high_value,
2 decode( a.interval, 'YES', b.interval ) interval
3 from user_tab_partitions a, user_part_tables b
4 where a.table_name = 'AUDIT_TRAIL'
5 and a.table_name = b.table_name
6 order by a.partition_position;
```

| PARTITION_ | TABLESPACE | HIGH_VALUE | INTERVAL |
|------------|------------|---------------------------------|-----------------------------|
| P0 | USERS | TIMESTAMP' 1900-01-01 00:00:00' | |
| SYS_P1623 | USERS | TIMESTAMP' 2014-03-01 00:00:00' | NUMTOYMINTERVAL(1, 'MONTH') |

Вспомнив материал раздела “Секционирование по диапазонам ключей”, приведенного ранее в главе, можно было бы ожидать, что этот оператор `INSERT` потерпит неудачу. Тем не менее, с учетом реализации секционирования по интервалам он завершится успешно и в действительности создаст новую секцию `SYS_P1623`. Значение `HIGH_VALUE` для этой секции составляет 01-MAR-2014, и если бы применялось секционирование по диапазонам, то оно предполагало бы, что все строки для дат, которые строго меньше 01-MAR-2014 и больше или равны 01-JAN-1900, отправлялись бы в эту секцию, но из-за наличия интервала правила отличаются. Когда интервал установлен, диапазон для этой секции включает все, что больше или равно `HIGH_VALUE-INTERVAL` и строго меньше, чем `HIGH_VALUE`. Поэтому секция имела бы следующий диапазон:

```
EODA@ORA12CR1> select TIMESTAMP' 2014-03-01 00:00:00'-NUMTOYMINTERVAL(1, 'MONTH')
2 greater_than_eq_to,
3 TIMESTAMP' 2014-03-01 00:00:00' strictly_less_than
4 from dual
5 /
```

```
GREATER_THAN_EQ_TO
```

```
STRICTLY_LESS_THAN
```

```
01-FEB-14 12.00.00.0000000000 AM
01-MAR-14 12.00.00.0000000000 AM
```

То есть секция содержит все данные для февраля 2014 года. Если мы вставим еще одну строку для другого месяца, как показано ниже, то заметим, что добавилась новая секция SYS_P1624, которая содержит все данные за июнь 2014 года:

```

EODA@ORA12CR1> insert into audit_trail (ts,data) values
  2  ( to_date('25-jun-2014','dd-mon-yyyy'), 'xx' );
1 row created.
1 строка создана.

```

```

EODA@ORA12CR1> select a.partition_name, a.tablespace_name, a.high_value,
  2      decode( a.interval, 'YES', b.interval ) interval
  3  from user_tab_partitions a, user_part_tables b
  4  where a.table_name = 'AUDIT_TRAIL'
  5  and a.table_name = b.table_name
  6  order by a.partition_position;

```

| PARTITION | TABLESPACE | HIGH_VALUE | INTERVAL |
|-----------|------------|---------------------------------|----------------------------|
| P0 | USERS | TIMESTAMP' 1900-01-01 00:00:00' | |
| SYS_P1623 | USERS | TIMESTAMP' 2014-03-01 00:00:00' | NUMTOYMINTERVAL(1,'MONTH') |
| SYS_P1624 | USERS | TIMESTAMP' 2014-07-01 00:00:00' | NUMTOYMINTERVAL(1,'MONTH') |

Представленный вывод может вызвать вопрос: почему все данные находятся в табличном пространстве USERS? Мы ясно указали, что данные должны быть распределены по табличным пространствам USERS и EXAMPLE, так почему все они оказались только в одном из них? Это связано с тем фактом, что когда база данных выясняет, в какую секцию должны направляться данные, она также определяет и табличное пространство, куда они должны попасть. Поскольку любая секция отстоит на четное количество месяцев от всех других, а мы используем всего два табличных пространства, в итоге одно и то же табличное пространство применяется снова и снова. Если мы будем загружать в таблицу данные только для каждого второго месяца, то это приведет к использованию лишь одного табличного пространства. Зайти в табличное пространство EXAMPLE можно путем добавления строки, которая отстоит от существующих данных на нечетное количество месяцев:

```

EODA@ORA12CR1> insert into audit_trail (ts,data) values
  2  ( to_date('15-mar-2014','dd-mon-yyyy'), 'xx' );
1 row created.
1 строка создана.

```

```

EODA@ORA12CR1> select a.partition_name, a.tablespace_name, a.high_value,
  2      decode( a.interval, 'YES', b.interval ) interval
  3  from user_tab_partitions a, user_part_tables b
  4  where a.table_name = 'AUDIT_TRAIL'
  5  and a.table_name = b.table_name
  6  order by a.partition_position;

```

| PARTITION | TABLESPACE | HIGH_VALUE | INTERVAL |
|-----------|------------|---------------------------------|----------------------------|
| P0 | USERS | TIMESTAMP' 1900-01-01 00:00:00' | |
| SYS_P1623 | USERS | TIMESTAMP' 2014-03-01 00:00:00' | NUMTOYMINTERVAL(1,'MONTH') |
| SYS_P1625 | EXAMPLE | TIMESTAMP' 2014-04-01 00:00:00' | NUMTOYMINTERVAL(1,'MONTH') |
| SYS_P1624 | USERS | TIMESTAMP' 2014-07-01 00:00:00' | NUMTOYMINTERVAL(1,'MONTH') |

Теперь применяется и пространство имен EXAMPLE. Новая секция оказалась между двумя существующими секциями, и она будет содержать данные за март 2014 года.

Может возникнуть вопрос: что произойдет, если в этот момент выполнить откат? Должно быть очевидным, что строки AUDIT_TRAIL, которые были недавно вставлены, исчезнут:

```
EODA@ORA12CR1> select * from audit_trail;
```

| TS | DATA |
|------------------------------|------|
| 27-FEB-14 12.00.00.000000 AM | xx |
| 15-MAR-14 12.00.00.000000 AM | xx |
| 25-JUN-14 12.00.00.000000 AM | xx |

```
EODA@ORA12CR1> rollback;
```

Rollback complete.

Откат выполнен.

```
EODA@ORA12CR1> select * from audit_trail;
```

no rows selected

нет выбранных строк

Однако не сразу ясно, что произойдет с добавленными секциями: они останутся или тоже исчезнут? Быстрый запрос подтвердит, что секции останутся:

```
EODA@ORA12CR1> select a.partition_name, a.tablespace_name, a.high_value,
2      decode( a.interval, 'YES', b.interval ) interval
3   from user_tab_partitions a, user_part_tables b
4   where a.table_name = 'AUDIT_TRAIL'
5   and a.table_name = b.table_name
6   order by a.partition_position;
```

| PARTITION | TABLESPACE | HIGH_VALUE | INTERVAL |
|-----------|------------|---------------------------------|-----------------------------|
| P0 | USERS | TIMESTAMP' 1900-01-01 00:00:00' | |
| SYS_P1623 | USERS | TIMESTAMP' 2014-03-01 00:00:00' | NUMTOYMINTERVAL(1, 'MONTH') |
| SYS_P1625 | EXAMPLE | TIMESTAMP' 2014-04-01 00:00:00' | NUMTOYMINTERVAL(1, 'MONTH') |
| SYS_P1624 | USERS | TIMESTAMP' 2014-07-01 00:00:00' | NUMTOYMINTERVAL(1, 'MONTH') |

Сразу же после создания секции фиксируются и становятся видимыми. Они создаются с использованием *рекурсивной* транзакции — транзакции, выполняемой отдельно и изолированно от любой другой транзакции, которая уже могла быть запущена. Когда мы переходим к вставке строки и база данных обнаруживает, что нужная секция пока не существует, база данных немедленно стартует новую транзакцию, обновляет словарь данных, отражая существование новой секции, и фиксирует работу. Она обязана делать это, иначе возникнет серьезная конкуренция (сериализация) между многими вставками, т.к. другим транзакциям пришлось бы ожидать фиксации, чтобы увидеть новую секцию. Следовательно, такая операция DDL выполняется за пределами существующей транзакции и секции сохраняются.

Вы могли обратить внимание, что база данных назначила имя новой секции — SYS_P1625. Эти имена не упорядочены и не особенно содержательны в привычном для многих людей смысле. Они показывают порядок, в котором секции добавлялись к таблице (хотя нельзя полагаться на то, что это всегда будет так; все может измениться), но не более того. Обычно в таблице, секционированной по диапазонам, ад-

министратор базы данных должен называть секции в соответствии с принятой схемой именования, и в большинстве случаев делать эти имена упорядоченными. Например, данные за февраль должны находиться в секции по имени PART_2014_02 (применяется формат PART_yyyy_mm), данные за март — в секции PART_2014_03 и т.д. При секционировании по интервалам вы не имеете контроля над именованием секций во время их создания, но можете легко переименовать их впоследствии. Скажем, мы могли бы запросить строку HIGH_VALUE и с использованием динамического SQL преобразовать ее в изящно сформатированные и осмысленные имена. Мы можем поступать так потому, что понимаем, каким образом желательно сформатировать имена, но база данных это не известно. Ниже приведен пример:

```

EODA@ORA12CR1> declare
  2     l_str varchar2(4000);
  3 begin
  4     for x in ( select a.partition_name, a.tablespace_name, a.high_value
  5                 from user_tab_partitions a
  6                 where a.table_name = 'AUDIT_TRAIL'
  7                     and a.interval = 'YES'
  8                     and a.partition_name like 'SYS\_P%' escape '\')
  9     loop
 10         execute immediate
 11             'select to_char( ' || x.high_value ||
 12               '-numtodsinterval(1, ''second''), ''PART_"yyyy_mm"' ) from dual'
 13             into l_str;
 14         execute immediate
 15             'alter table audit_trail rename partition "' ||
 16               x.partition_name || '" to "' || l_str || ''';
 17     end loop;
 18 end;
 19 /
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.

```

Итак, здесь мы берем значение HIGH_VALUE и вычитаем из него одну секунду. Мы знаем, что HIGH_VALUE представляет значение *строго меньше чем*, поэтому значение секундой раньше попадает в нужный диапазон. Затем к результирующему значению типа TIMESTAMP мы применяем формат "PART_"yyyy_mm и получаем строку наподобие PART_2014_03 для марта 2014 года. Мы используем эту строку в команде переименования, и теперь словарь данных выглядит следующим образом:

```

EODA@ORA12CR1> select a.partition_name, a.tablespace_name, a.high_value,
  2     decode( a.interval, 'YES', b.interval ) interval
  3 from user_tab_partitions a, user_part_tables b
  4 where a.table_name = 'AUDIT_TRAIL'
  5     and a.table_name = b.table_name
  6 order by a.partition_position;

```

| PARTITION | TABLESPACE | HIGH_VALUE | INTERVAL |
|--------------|------------|---------------------------------|-----------------------------|
| P0 | USERS | TIMESTAMP' 1900-01-01 00:00:00' | |
| PART_2014_02 | USERS | TIMESTAMP' 2014-03-01 00:00:00' | NUMTOYMINTERVAL(1, 'MONTH') |
| PART_2014_03 | EXAMPLE | TIMESTAMP' 2014-04-01 00:00:00' | NUMTOYMINTERVAL(1, 'MONTH') |
| PART_2014_06 | USERS | TIMESTAMP' 2014-07-01 00:00:00' | NUMTOYMINTERVAL(1, 'MONTH') |

Потребуется только время от времени запускать этот сценарий для переименования любых добавленных секций, чтобы поддерживать аккуратное соглашение об именовании. Имейте в виду, что во избежание проблем с атакой внедрением SQL (здесь применяется конкатенация строк, а не переменные привязки, которые в DDL использовать невозможно), мы должны сохранить данный сценарий как анонимный блок или как процедуру с правами вызывающего, если решим сделать его хранимой процедурой. Это предотвратит запуск SQL-кода другими в нашей схеме от нашего имени, что может оказаться опасным.

Секционирование по ссылкам

Секционирование по ссылкам является средством Oracle 11g Release 1 и последующих версий. Оно призвано решить проблему эквисекционирования типа “родительская—дочерняя”, т.е. когда дочерняя таблица должна быть секционирована так, чтобы каждая ее секция находилась в отношении “один к одному” с секцией родительской таблицы. Это важно в ситуациях типа хранилища данных, где требуется держать доступным в оперативном режиме определенный объем данных (таких как информация о заказах (ORDER) за последние пять лет) и нужно гарантировать, что связанные дочерние данные (ORDER_LINE_ITEMS) также доступны. В этом классическом примере таблица ORDERS обычно должна иметь столбец ORDER_DATE, что облегчает секционирование по месяцам и, как следствие, позволяет поддерживать в оперативном режиме данные за последние пять лет. С течением времени необходимо лишь просто делать доступной для загрузки секцию, соответствующую следующему месяцу, и удалять самую старую секцию. Однако, взглянув на таблицу ORDER_LINE_ITEMS, вы заметите, что есть одна проблема. Таблица не имеет столбца ORDER_DATE, и в ней нет ничего, на чем можно было бы провести секционирование; следовательно, это не способствует упрощению удаления старой информации или загрузки новой.

В прошлом, до появления секционирования по ссылкам, разработчики вынуждены были производить денормализацию данных, в сущности, копируя атрибут ORDER_DATE из родительской таблицы ORDERS в дочернюю таблицу ORDER_LINE_ITEMS. Это порождало типичные проблемы избыточности данных, приводящие к росту накладных расходов при хранении, увеличению потребления ресурсов при загрузке данных, потребностям в каскадном обновлении (если вы модифицируете родительскую таблицу, то должны обеспечить обновление всех копий родительских данных) и т.д. Вдобавок, если в базе данных включены ограничения внешнего ключа (что и должно делаться), то вы обнаружите, что возможность усечения или удаления старых секций в родительской таблице утеряна. Давайте в качестве примера создадим обычные таблицы ORDERS и ORDER_LINE_ITEM, начав с ORDERS:

```
EODA@ORA12CR1> create table orders
2 (
3   order#      number primary key,
4   order_date  date,
5   data        varchar2(30)
6 )
7 enable row movement
8 PARTITION BY RANGE (order_date)
9 (
```

```

10 PARTITION part_2014 VALUES LESS THAN (to_date('01-01-2015','dd-mm-yyyy')) ,
11 PARTITION part_2015 VALUES LESS THAN (to_date('01-01-2016','dd-mm-yyyy'))
12 )
13 /

```

Table created.

Таблица создана.

```

EODA@ORA12CR1> insert into orders values
  2 ( 1, to_date( '01-jun-2014', 'dd-mon-yyyy' ), 'xxx' );
1 row created.
1 строка создана.

```

```

EODA@ORA12CR1> insert into orders values
  2 ( 2, to_date( '01-jun-2015', 'dd-mon-yyyy' ), 'xxx' );
1 row created.
1 строка создана.

```

А теперь создадим таблицу ORDER_LINE_ITEMS с частью данных, которые указывают на таблицу ORDERS:

```

EODA@ORA12CR1> create table order_line_items
  2 (
  3   order#      number,
  4   line#       number,
  5   order_date  date, -- manually copied from ORDERS!
  6   data        varchar2(30),
  7   constraint c1_pk primary key(order#,line#),
  8   constraint c1_fk_p foreign key(order#) references orders
  9 )
10 enable row movement
11 PARTITION BY RANGE (order_date)
12 (
13   PARTITION part_2014 VALUES LESS THAN (to_date('01-01-2015','dd-mm-yyyy')) ,
14   PARTITION part_2015 VALUES LESS THAN (to_date('01-01-2016','dd-mm-yyyy'))
15 )
16 /

```

Table created.

Таблица создана.

```

EODA@ORA12CR1> insert into order_line_items values
  2 ( 1, 1, to_date( '01-jun-2014', 'dd-mon-yyyy' ), 'yyy' );
1 row created.
1 строка создана.

```

```

EODA@ORA12CR1> insert into order_line_items values
  2 ( 2, 1, to_date( '01-jun-2015', 'dd-mon-yyyy' ), 'yyy' );
1 row created.
1 строка создана.

```

Если бы нам пришлось удалить секцию таблицы ORDER_LINE_ITEMS, содержащую данные для 2014 года, то мы знаем, что соответствующую секцию таблицы ORDERS для 2014 года также можно было бы удалить, не нарушая ограничение ссылочной целостности. Хотя мы об этом знаем, базе данных такой факт не известен:

```

EODA@ORA12CR1> alter table order_line_items drop partition part_2014;
Table altered.
Таблица изменена.

```



```

EODA@ORA12CR1> alter table orders drop partition part_2014;
alter table orders drop partition part_2014
*
ERROR at line 1:
ORA-02266: unique/primary keys in table referenced by enabled foreign keys
ОШИБКА в строке 1:
ORA-02266: в таблице есть уникальный/первичный ключ, на который ссылаются
внешние ключи

```

Таким образом, подход с денормализацией данных не только является громоздким, ресурсоемким и потенциально нарушающим целостность данных, он также препятствует выполнению того, что часто должно делаться в секционированных таблицах — очистке устаревшей информации.

Именно здесь в игру вступает секционирование по ссылкам. При таком секционировании дочерняя таблица наследует схему секционирования своей родительской таблицы без необходимости в денормализации ключа секционирования и позволяет базе данных знать, что дочерняя таблица эквисекционирована с родительской таблицей. То есть мы имеем возможность удалять или усекать секцию родительской таблицы, когда усекаем или удаляем соответствующую секцию дочерней таблицы.

Ниже приведен простой синтаксис, воссоздающий предыдущий пример. Мы повторно применим существующую родительскую таблицу ORDERS и просто выполним ее усечение:

```

EODA@ORA12CR1> drop table order_line_items cascade constraints;
Table dropped.
Таблица удалена.

EODA@ORA12CR1> truncate table orders;
Table truncated.
Таблица усечена.

EODA@ORA12CR1> insert into orders values
  2 ( 1, to_date( '01-jun-2014', 'dd-mon-yyyy' ), 'xxx' );
1 row created.
1 строка создана.

EODA@ORA12CR1> insert into orders values
  2 ( 2, to_date( '01-jun-2015', 'dd-mon-yyyy' ), 'xxx' );
1 row created.
1 строка создана.

```

Теперь создадим новую дочернюю таблицу:

```

EODA@ORA12CR1> create table order_line_items
  2 (
  3   order#      number,
  4   line#       number,
  5   data        varchar2(30),
  6   constraint c1_pk primary key(order#,line#),
  7   constraint c1_fk_p foreign key(order#) references orders
  8 )
  9 enable row movement
 10 partition by reference(c1_fk_p)
 11 /
Table created.
Таблица создана.

```

```

EODA@ORA12CR1> insert into order_line_items values ( 1, 1, 'zzz' );
1 row created.
1 строка создана.

EODA@ORA12CR1> insert into order_line_items values ( 2, 1, 'zzz' );
1 row created.
1 строка создана.

```

Магия скрыта в строке 10 оператора CREATE TABLE. В ней конструкция PARTITION BY RANGE заменена конструкцией PARTITION BY REFERENCE.

На заметку! Если при работе в Oracle 11g Release 1 сообщается об ошибке ORA-14652: r eference partitioning foreign key is not supported (ORA-14652: внешний ключ секционирования по ссылкам не поддерживается), то причина в том, что этот выпуск Oracle требует обязательного наличия ограничения NOT NULL на каждом столбце внешнего ключа. Поскольку ORDER# — часть первичного ключа, мы знаем, что она не равна NULL, но выпуск Oracle 11g Release 1 это не распознает. В таком случае понадобится оп-ределить столбцы внешнего ключа как NOT NULL.

Это позволяет именовать ограничение внешнего ключа так, чтобы выяснить, ка-кой будет схема секционирования. Здесь мы видим, что внешний ключ указывает на таблицу ORDERS: база данных читает структуру таблицы ORDERS и определяет, что она имеет две секции, а потому дочерняя таблица также будет иметь две секции. В сущности, если запросить словарь данных прямо сейчас, то можно увидеть, что эти две таблицы имеют одну и ту же структуру секционирования:

```

EODA@ORA12CR1> select table_name, partition_name
2   from user_tab_partitions
3   where table_name in ( 'ORDERS', 'ORDER_LINE_ITEMS' )
4   order by table_name, partition_name
5   /

```

| TABLE_NAME | PARTITION_NAME |
|------------------|----------------|
| ORDERS | PART_2014 |
| ORDERS | PART_2015 |
| ORDER_LINE_ITEMS | PART_2014 |
| ORDER_LINE_ITEMS | PART_2015 |

Более того, т.к. базе данных известно, что эти две таблицы связаны, мы можем удалить секцию родительской таблицы и заставить базу данных автоматически очис-тить связанные секции дочерней таблицы (поскольку дочерняя таблица наследует структуру секционирования у родительской таблицы, любые изменения в структуре секционирования родительской таблицы распространяются каскадным образом):

```

EODA@ORA12CR1> alter table orders drop partition part_2014 update global
indexes;
Table altered.
Таблица изменена.

EODA@ORA12CR1> select table_name, partition_name
2   from user_tab_partitions
3   where table_name in ( 'ORDERS', 'ORDER_LINE_ITEMS' )
4   order by table_name, partition_name
5   /

```

| TABLE_NAME | PARTITION_NAME |
|------------------|----------------|
| ORDERS | PART_2015 |
| ORDER_LINE_ITEMS | PART_2015 |

Таким образом, операция DROP, выполнение которой ранее не допускалось, теперь разрешена, и она автоматически распространяется каскадным образом на дочернюю таблицу. Кроме того, если добавить секцию, как показано ниже, можно заметить, что эта операция также распространяется каскадным образом; отношение “один к одному” между секциями родительской и дочерней таблицей сохраняется:

```

EODA@ORA12CR1> alter table orders add partition
  2 part_2016 values less than
  3 (to_date( '01-01-2017', 'dd-mm-yyyy' ));
Table altered.

```

Таблица изменена.

```

EODA@ORA12CR1> select table_name, partition_name
  2 from user_tab_partitions
  3 where table_name in ( 'ORDERS', 'ORDER_LINE_ITEMS' )
  4 order by table_name, partition_name
  5 /

```

| TABLE_NAME | PARTITION_NAME |
|------------------|----------------|
| ORDERS | PART_2015 |
| ORDERS | PART_2016 |
| ORDER_LINE_ITEMS | PART_2015 |
| ORDER_LINE_ITEMS | PART_2016 |

В приведенном выше операторе CREATE TABLE имеется часть, которая пока еще не обсуждалась — ENABLE ROW MOVEMENT. Эта конструкция была добавлена в Oracle8i, и она детально рассматривается ниже в отдельном разделе. Выражаясь кратко, ее синтаксис позволяет оператору UPDATE модифицировать значение ключа секции так, что это вызовет перемещение строки из текущей секции в какую-то другую. В версиях, предшествующих Oracle8i, такая операция не допускалась; ключи секций можно было изменять, но это не приводило к перемещению строк в другие секции.

Поскольку родительская таблица изначально определена как допускающая перемещение строк, мы обязаны определить все дочерние таблицы (а также их дочерние таблицы и т.д.) с той же самой возможностью, чтобы в случае использования секционирования по ссылкам при перемещении родительской строки дочерние строки также должны перемещаться. Например:

```

EODA@ORA12CR1> select '2015', count(*) from order_line_items
partition(part_2015)
  2 union all
  3 select '2016', count(*) from order_line_items partition(part_2016);
'201  COUNT(*)
----
2015          1
2016          0

```

Можно заметить, что прямо сейчас данные в дочерней таблице ORDER_LINE_ITEMS находятся внутри секции 2015. Выполнив простое обновление родительской таблицы ORDERS, можно удостовериться, что данные переместились и в дочерней таблице:

```
EODA@ORA12CR1> update orders set order_date = add_months(order_date,12);
1 row updated.
1 строка обновлена.

EODA@ORA12CR1> select '2015', count(*) from order_line_items partition(part_2015)
2 union all
3 select '2016', count(*) from order_line_items partition(part_2016);
'201  COUNT(*)
----
2015      0
2016      1
```

Обновление родительской таблицы каскадным образом распространилось на дочернюю таблицу и вызвало перемещение строки (или при необходимости нескольких строк).

Подводя итоги, секционирование по ссылкам устраняет необходимость в денормализации данных при секционировании родительских и дочерних таблиц. Кроме того, удаление родительской секции приводит к автоматическому удалению связанной дочерней секции. Такие свойства очень полезны в средах хранилищ данных.

Секционирование по интервалам ключей и по ссылкам

До выхода версии Oracle 12c комбинация секционирования по интервалам ключей и секционирования по ссылкам не поддерживалась. Например, если создать родительскую таблицу, секционированную по интервалам, в Oracle 11g:

```
EODA@ORA11GR2> create table orders
2 (order#          number primary key,
3  order_date      timestamp,
4  data            varchar2(30))
5 PARTITION BY RANGE (order_date)
6 INTERVAL (numtoyminterval(1,'year'))
7 (PARTITION part_2014 VALUES LESS THAN (to_date('01-01-2015','dd-mm-yyyy')) ,
8  PARTITION part_2015 VALUES LESS THAN (to_date('01-01-2016','dd-mm-yyyy')));
Table created.
Таблица создана.
```

А затем попытаться создать дочернюю таблицу, секционированную по ссылкам, то возникнет ошибка:

```
EODA@ORA11GR2> create table order_line_items
2 ( order#          number,
3  line#            number,
4  data            varchar2(30),
5  constraint c1_pk primary key(order#,line#),
6  constraint c1_fk_p foreign key(order#) references orders)
7 partition by reference(c1_fk_p);
create table order_line_items
*
ERROR at line 1:
```

ORA-14659: Partitioning method of the parent table is not supported

ОШИБКА в строке 1:

ORA-14659: метод секционирования родительской таблицы не поддерживается

В версии Oracle 12c эта проблема исчезла, т.к. здесь можно сочетать секционирование по интервалам ключей и секционирование по ссылкам. Предшествующий код успешно выполнится в базе данных Oracle 12c, создав необходимую дочернюю таблицу:

```

EODA@ORA12CR1> create table order_line_items
2  ( order#          number,
3    line#           number,
4    data            varchar2(30),
5    constraint c1_pk primary key(order#,line#),
6    constraint c1_fk_p foreign key(order#) references orders)
7  partition by reference(c1_fk_p);

```

Table created.

Таблица создана.

Чтобы увидеть секционирование по интервалам и по ссылкам в действии, давайте вставим какие-нибудь данные. Для начала мы вставим строки, которые будут умещаться в существующие секции по диапазонам:

```

EODA@ORA12CR1> insert into orders values (1, to_date( '01-jun-2014',
❖'dd-mon-yyyy' ), 'xxx');

```

1 row created.

1 строка создана.

```

EODA@ORA12CR1> insert into orders values (2, to_date( '01-jun-2015',
❖'dd-mon-yyyy' ), 'xxx');

```

1 row created.

1 строка создана.

```

EODA@ORA12CR1> insert into order_line_items values( 1, 1, 'yyy' );

```

1 row created.

1 строка создана.

```

EODA@ORA12CR1> insert into order_line_items values( 2, 1, 'yyy' );

```

1 row created.

1 строка создана.

Все эти строки попадают в секции, указанные во время создания таблиц. Следующий запрос отображает текущие секции:

```

EODA@ORA12CR1> select table_name, partition_name from user_tab_partitions
2  where table_name in ( 'ORDERS', 'ORDER_LINE_ITEMS' )
3  order by table_name, partition_name;

```

| TABLE_NAME | PARTITION_NAME |
|------------------|----------------|
| ORDERS | PART_2014 |
| ORDERS | PART_2015 |
| ORDER_LINE_ITEMS | PART_2014 |
| ORDER_LINE_ITEMS | PART_2015 |

Далее мы вставим строки, которые не умещаются ни в одну секцию по диапазону из числа имеющихся; вследствие этого Oracle автоматически создает секции для хранения вновь вставленных строк:

```

EODA@ORA12CR1> insert into orders values (3, to_date( '01-jun-2016',
❏ 'dd-mon-yyyy' ), 'xxx');
1 row created.
1 строка создана.

EODA@ORA12CR1> insert into order_line_items values (3, 1, 'zzz' );
1 row created.
1 строка создана.

```

Приведенный ниже запрос показывает, что были автоматически созданы две интервальных секции, одна для родительской таблицы и одна для дочерней:

```

EODA@ORA12CR1> select a.table_name, a.partition_name, a.high_value,
2 decode( a.interval, 'YES', b.interval ) interval
3 from user_tab_partitions a, user_part_tables b
4 where a.table_name IN ('ORDERS', 'ORDER_LINE_ITEMS')
5 and a.table_name = b.table_name
6 order by a.table_name;

```

| PARTITION_ | TABLESPACE | HIGH_VALUE | INTERVAL |
|----------------------------|------------|---------------------------------|----------|
| ORDERS | PART_2014 | TIMESTAMP' 2015-01-01 00:00:00' | |
| ORDERS | PART_2015 | TIMESTAMP' 2016-01-01 00:00:00' | |
| ORDERS | SYS_P1626 | TIMESTAMP' 2017-01-01 00:00:00' | |
| NUMTOYMINTERVAL(1, 'YEAR') | | | |
| ORDER_LINE_ITEMS | PART_2014 | | |
| ORDER_LINE_ITEMS | PART_2015 | | |
| ORDER_LINE_ITEMS | SYS_P1626 | | YES |

В выводе видно, что созданы две секции по имени SYS_P1626, и секция родительской таблицы имеет значение HIGH_VALUE, равное 2017-01-01. При желании можете переименовать секции посредством команды ALTER TABLE:

```

EODA@ORA12CR1> alter table orders rename partition sys_p1626 to part_2016;
Table altered.
Таблица изменена.

EODA@ORA12CR1> alter table order_line_items rename partition sys_p1626
❏ to part_2016;
Table altered.
Таблица изменена.

```

Совет. Пример автоматизации переименования секций с помощью PL/SQL рассматривался в разделе “Секционирование по интервалам ключей” ранее в этой главе.

Секционирование по виртуальному столбцу

Секционирование по виртуальному столбцу позволяет разбивать на секции, основываясь на SQL-выражении. Такой тип секционирования удобен, когда столбец таблицы перегружен множеством бизнес-значений, и необходимо разбивать на секции по какой-то части этого столбца. Например, пусть в таблице имеется столбец RESERVATION_CODE:

```

EODA@ORA12CR1> create table res(reservation_code varchar2(30));
Table created.
Таблица создана.

```

Первый символ значения в столбце `RESERVATION_CODE` определяет регион, из которого поступил запрос на бронирование. Для целей этого примера предположим, что первый символ A или C соответствует региону NE, первый символ B — региону SW, а первый символ D — региону NW.

Вставим в таблицу тестовые данные:

```
EODA@ORA12CR1> insert into res (reservation_code)
  2  select chr(64+(round(dbms_random.value(1,4)))) || level
  3  from dual connect by level < 100000;

EODA@ORA12CR1> select * from res;
```

Ниже показан неполный фрагмент вывода:

```
RESERVATION_CODE
-----
C1
D2
...
A72827
B72828
```

В этом сценарии нам известно, что первый символ представляет регион, и мы хотим иметь возможность провести секционирование по списку на столбце `RESERVATION_CODE`. Секционирование по виртуальному столбцу позволяет применять SQL-функцию к столбцу и секционировать по списку на основе первого символа. Вот как будет выглядеть определение таблицы с секционированием по виртуальному столбцу:

```
EODA@ORA12CR1> create table res(
  reservation_code varchar2(30),
  region as
    (decode(substr(reservation_code,1,1), 'A', 'NE'
                                     , 'C', 'NE'
                                     , 'B', 'SW'
                                     , 'D', 'NW')
    )
)
partition by list (region)
(partition p1 values('NE'),
 partition p2 values('SW'),
 partition p3 values('NW'));
Table created.
Таблица создана.
```

Таким образом, секционирование по виртуальному столбцу часто подходит в ситуации, когда есть бизнес-требование секционировать по частям данных в столбце или по комбинации данных из разных столбцов (особенно при отсутствии очевидного способа выполнить секционирование по списку или диапазону). Выражение, лежащее в основе виртуального столбца, может быть довольно сложным и возвращать подмножество последовательности столбцов, комбинацию значений в столбцах и т.п.

Составное секционирование

И, наконец, мы рассмотрим ряд примеров составного секционирования, которое представляет собой смесь секционирования по диапазонам, хеш-секционирования и/или секционирования по списку. Методы, посредством которых можно комбинировать секционирование, т.е. типы схем секционирования, которые допускается смешивать, варьируются в зависимости от выпуска. В табл. 13.1 показано то, что доступно в каждом крупном выпуске. Сверху вниз перечислены разрешенные схемы секционирования верхнего уровня, а слева направо — схемы добавочного секционирования внутри секций.

Таблица 13.1. Схемы составного секционирования, поддерживаемые в разных версиях Oracle

| | Секционирование по диапазонам | Секционирование по списку | Хеш-секционирование |
|-------------------------------|-------------------------------|---------------------------|----------------------|
| Секционирование по диапазонам | Oracle 11g Release 1 | Oracle9i Release 2 | Oracle9i Release 1 |
| Секционирование по списку | Oracle 11g Release 1 | Oracle 11g Release 1 | Oracle 11g Release 1 |
| Хеш-секционирование | Oracle 11g Release 2 | Oracle 11g Release 2 | Oracle 11g Release 2 |

Итак, например, в Oracle9i Release 2 и последующих версиях таблицу можно секционировать по диапазонам и затем внутри каждой секции диапазона секционировать по спискам или хеш-значениям. В Oracle 11g Release 1 и последующих версиях количество составных схем возрастает с двух до шести. А в Oracle 11g Release 2 и последующих версиях на выбор доступно девять составных схем.

Интересно отметить, что когда вы используете составное секционирование, сегменты секций отсутствуют — есть только сегменты подсекций. При составном секционировании сами секции не имеют сегментов (подобно тому, как секционированная таблица не имеет сегмента). Данные физически хранятся в сегментах подсекций, и секция становится *логическим контейнером*, или контейнером, который указывает на действительные подсекции.

В нашем примере мы взглянем на составное секционирование по диапазонам и по хеш-значениям. Для секционирования по диапазонам мы применяем набор столбцов, который отличается от набора, используемого при хеш-секционировании. Это не является обязательным; мы могли бы выбрать для обоих видов секционирования один и тот же набор столбцов:

```

EODA@ORA12CR1> CREATE TABLE composite_example
2  ( range_key_column  date,
3    hash_key_column   int,
4    data               varchar2(20)
5  )
6  PARTITION BY RANGE (range_key_column)
7  subpartition by hash(hash_key_column) subpartitions 2
8  (
9    PARTITION part_1
10   VALUES LESS THAN(to_date('01/01/2014','dd/mm/yyyy'))
11   (subpartition part_1_sub_1,
12     subpartition part_1_sub_2
13   ),

```



```

14 PARTITION part_2
15     VALUES LESS THAN(to_date('01/01/2015', 'dd/mm/yyyy'))
16     (subpartition part_2_sub_1,
17      subpartition part_2_sub_2
18     )
19 )
20 /

```

Table created.

Таблица создана.

При составном секционировании по диапазонам и по хеш-значениям Oracle сначала применяет правила секционирования по диапазонам, чтобы выяснить, в какой диапазон попадают данные. Затем с использованием хеш-функции решается, в какую физическую секцию в конечном итоге должны быть помещены данные. Процесс представлен на рис. 13.4.

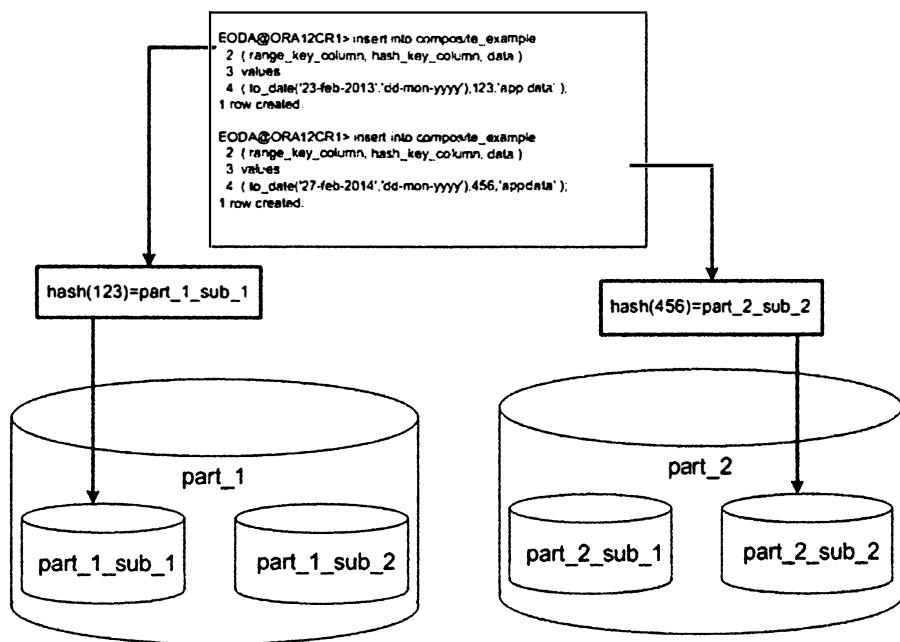


Рис. 13.4. Пример составного секционирования по диапазонам и по хеш-значениям

Таким образом, составное секционирование дает возможность разбивать данные по диапазону и, когда диапазон оказывается слишком большим или может быть полезно дальнейшее исключение секций, дополнительно разбить данные по хеш-значениям или списку. Следует отметить, что секции по диапазонам не обязаны иметь одинаковые количества подсекций. Предположим для примера, что вы секционировали таблицу по диапазону значений столбца даты в качестве поддержки процедуры очистки данных (быстрого и легкого удаления всех устаревших данных). В секции для 2013 года и предыдущих лет имеются равные объемы данных с нечетными и четными номерами кодов в столбце CODE_KEY_COLUMN. Но вы знаете, что после этого периода количество записей, ассоциированных с нечетными номерами кодов, вдвое

больше, и хотите иметь больше подсекций для записей с нечетными номерами кодов. Достичь этого довольно легко, просто определив дополнительные подсекции:

```

EODA@ORA12CR1> CREATE TABLE composite_range_list_example
 2 ( range_key_column date,
 3   code_key_column int,
 4   data varchar2(20)
 5 )
 6 PARTITION BY RANGE (range_key_column)
 7 subpartition by list(code_key_column)
 8 (
 9   PARTITION part_1
10     VALUES LESS THAN(to_date('01/01/2014','dd/mm/yyyy'))
11     (subpartition part_1_sub_1 values( 1, 3, 5, 7 ),
12      subpartition part_1_sub_2 values( 2, 4, 6, 8 )
13     ),
14   PARTITION part_2
15     VALUES LESS THAN(to_date('01/01/2015','dd/mm/yyyy'))
16     (subpartition part_2_sub_1 values ( 1, 3 ),
17      subpartition part_2_sub_2 values ( 5, 7 ),
18      subpartition part_2_sub_3 values ( 2, 4, 6, 8 )
19     )
20 )
21 /

```

Table created.

Таблица создана.

Здесь в итоге создаются пять секций: две подсекции в секции PART_1 и три — в секции PART_2.

Перемещение строк

Вас может интересовать, что произойдет, если значение столбца, применяемого для определения секции, модифицировано в любой из предшествующих схем секционирования. Существуют два возможных случая.

- Модификация может не вызвать необходимости использования другой секции; строка будет по-прежнему принадлежать той же секции, где она была. Это поддерживается во всех случаях.
- Модификация может привести к перемещению строки между секциями. Это поддерживается, если перемещение строк включено для данной таблицы; в противном случае возникнет ошибка.

Мы легко можем наблюдать такое поведение. В примере, который рассматривался в разделе “Секционирование по диапазонам ключей”, мы вставили пару строк в секцию PART_1 таблицы RANGE_EXAMPLE:

```

EODA@ORA12CR1> CREATE TABLE range_example
 2 ( range_key_column date ,
 3   data varchar2(20)
 4 )
 5 PARTITION BY RANGE (range_key_column)
 6 ( PARTITION part_1 VALUES LESS THAN
 7   (to_date('01/01/2014','dd/mm/yyyy')) ,

```

```

8   PARTITION part_2 VALUES LESS THAN
9     (to_date('01/01/2015','dd/mm/yyyy'))
10 )
11 /

```

Table created.

Таблица создана.

```

EODA@ORA12CR1> insert into range_example

```

```

2  ( range_key_column, data )
3  values
4  ( to_date( '15-dec-2013 00:00:00',
5            'dd-mon-yyyy hh24:mi:ss' ),
6    'application data...' );

```

1 row created.

1 строка создана.

```

EODA@ORA12CR1> insert into range_example

```

```

2  ( range_key_column, data )
3  values
4  ( to_date( '01-jan-2014 00:00:00',
5            'dd-mon-yyyy hh24:mi:ss' )-1/24/60/60,
6    'application data...' );

```

1 row created.

1 строка создана.

```

EODA@ORA12CR1> select * from range_example partition(part_1);
RANGE_KEY DATA

```

```

-----
15-DEC-13 application data...
31-DEC-13 application data...

```

Мы берем одну из строк и обновляем значение в ее столбце RANGE_KEY_COLUMN так, что она остается в секции PART_1:

```

EODA@ORA12CR1> update range_example
2   set range_key_column = trunc(range_key_column)
3   where range_key_column =
4     to_date( '31-dec-2013 23:59:59',
5             'dd-mon-yyyy hh24:mi:ss' );

```

1 row updated.

1 строка обновлена.

Как и ожидалось, операция завершается успешно: строка остается в секции PART_1. Затем мы устанавливаем для столбца RANGE_KEY_COLUMN в строке значение, которое приведет к тому, что строка станет принадлежать секции PART_2:

```

EODA@ORA12CR1> update range_example
2   set range_key_column = to_date('01-jan-2014','dd-mon-yyyy')
3   where range_key_column = to_date('31-dec-2013','dd-mon-yyyy');
update range_example
*

```

ERROR at line 1:

ORA-14402: updating partition key column would cause a partition change
ОШИБКА в строке 1:

ORA-14402: обновление столбца ключа секционирования приведет к изменению секции

Немедленно возникает ошибка, т.к. перемещение строк между секциями не было явно разрешено. В Oracle8i и последующих версиях для этой таблицы можно включить перемещение строк из секции в секцию.

Однако, выполняя такое действие, вы должны помнить о тонком побочном эффекте, а именно — в результате этого обновления изменится ROWID строки:

```
EODA@ORA12CR1> select rowid
2   from range_example
3   where range_key_column = to_date('31-dec-2013', 'dd-mon-yyyy');
```

ROWID

AAAtzXAAGAAAa06AAB

```
EODA@ORA12CR1> alter table range_example enable row movement;
Table altered.
```

Таблица изменена.

```
EODA@ORA12CR1> update range_example
2   set range_key_column = to_date('01-jan-2014', 'dd-mon-yyyy')
3   where range_key_column = to_date('31-dec-2013', 'dd-mon-yyyy');
```

1 row updated.

1 строка обновлена.

```
EODA@ORA12CR1> select rowid
2   from range_example
3   where range_key_column = to_date('01-jan-2014', 'dd-mon-yyyy');
```

ROWID

AAAtzYAAGAAAae6AAA

При условии понимания, что при таком обновлении будет изменяться ROWID строки, включение перемещения строк позволит вам обновлять ключи секционирования.

На заметку! Существуют и другие ситуации, при которых ROWID может изменяться в результате обновления. Это может произойти из-за обновления первичного ключа индекс-таблицы. Универсальный идентификатор ROWID в такой строке также изменится. Команда FLASHBACK TABLE в Oracle 10g и последующих версиях тоже может изменять ROWID строк, как и команда ALTER TABLE SHRINK в тех же версиях.

Вы должны понимать, что перемещение строки внутренне выполняется так, как если бы вы фактически удалили ее и вставили заново. При этом обновляется каждый индекс таблицы, удаляются его старые записи и вставляются новые. То есть выполняется физическая работа операторов DELETE и INSERT. Однако Oracle трактует это как обновление, даже несмотря на то, что производится удаление и вставка строки — т.е. триггеры INSERT и DELETE не запускаются, а инициируются только триггеры UPDATE. К тому же дочерние таблицы, которые могут помешать удалению из-за ограничений внешнего ключа, не делают этого. Тем не менее, следует подготовиться к тому, что будет выполняться дополнительная работа; перемещение строки является намного более затратным в плане ресурсов, чем обычное обновление. Следовательно, проектное решение, допускающее частую модификацию ключей секционирования, которая приводит к перемещению строк между секциями, нельзя считать удачным.

Заключительные соображения по поводу схем секционирования таблиц

В общем случае секционирование по диапазонам полезно, когда есть данные, которые логически разделены по какому-то значению (значениям). Классическим примером могут служить данные, основанные на времени — “Квартал продаж”, “Финансовый год” или “Месяц”. Секционирование по диапазонам позволяет во многих случаях воспользоваться преимуществом исключения секций, в том числе применение точного равенства и диапазонов (меньше, больше, между и т.д.).

Хеш-секционирование подходит для данных, не имеющих естественных диапазонов значений, по которым их можно было бы разделить. Например, если нужно загрузить таблицу результатов переписи населения, то в них может не оказаться атрибута, по которому имело бы смысл секционировать по диапазонам. Однако вас все равно могут интересовать улучшения администрирования, производительности и доступности, предлагаемые секционированием. Для этого понадобится просто указать уникальный или почти уникальный набор столбцов для хеширования. В результате будет обеспечено равномерное распределение данных между многими секциями. Хеш-секционированные объекты используют преимущества исключения секций, когда в критериях выборки применяется явное равенство или конструкция `IN (значение, значение, . . .)`, но не в случае использования диапазонов данных.

Секционирование по списку подходит для данных, у которых имеется столбец с дискретными наборами значений, и секционирование по такому столбцу имеет смысл в соответствии со способом его применения в приложении (например, оно позволяет легко исключать секции в запросах). Классическим примером может быть код штата или региона или любые другие атрибуты, выраженные с помощью кода.

Секционирование по интервалам ключей расширяет секционирование по диапазонам и позволяет автоматически добавлять секции, когда вставляемые в таблицу данные не вписываются в существующие секции. Значительное улучшение связано с меньшим объемом работ по обслуживанию (поскольку администратору базы данных не придется обязательно отслеживать диапазоны и вручную добавлять секции).

Секционирование по ссылкам облегчает реализацию секционированных таблиц, которые связаны через ограничения ссылочной целостности. Оно позволяет дочерней таблице логически секционироваться в той же манере, что и родительская таблица, без необходимости в дублировании столбцов родительской таблицы внутри дочерней таблицы.

Секционирование по интервалам ключей и по ссылкам предоставляет возможность комбинировать схему секционирования по интервалам ключей со схемой секционирования по ссылкам. Это средство появилось в версии Oracle 12c и полезно, когда нужно использовать схемы секционирования по интервалам ключей и по ссылкам в тандеме.

Секционирование по виртуальному столбцу позволяет разбивать на секции с применением виртуального столбца в качестве ключа. Это средство предлагает гибкость в секционировании по подстроке, полученной из значения обычного столбца (или по любому другому SQL-выражению). Такой вид секционирования удобен, когда использовать существующий столбец для ключа нереально, но можно разбивать на секции, основываясь на части значения, которое содержится в каком-то столбце.

Составное секционирование полезно, когда имеются логические принципы, по которым можно провести секционирование по диапазонам, но результирующие секции диапазонов все равно остаются слишком большими, чтобы ими можно было эффективно управлять. Вы можете применить секционирование по диапазонам, списку или хеш-значениям, а затем дополнительно разделить каждый диапазон с помощью хеш-функции либо использовать секционирование по списку или даже по диапазонам. Это позволяет в любой крупной секции разнести запросы ввода-вывода по множеству устройств. Кроме того, теперь можно достичь исключения секций на трех уровнях. Если вы запускаете запрос по ключу секционирования, то Oracle может исключить из рассмотрения любые секции, которые не удовлетворяют указанному критерию. Если вы добавите к запросу ключ подсекции, то Oracle сможет исключить другие подсекции в пределах этой секции. Если же вы запрашиваете только по ключу подсекции (не применяя ключ секции), то Oracle будет обращаться лишь к подходящим хеш- или списковым подсекциям из каждой секции.

Поэтому если есть что-то, на основании чего имеет смысл секционировать данные по диапазонам, то вы должны проводить его поверх хеш-секционирования или секционирования по списку. Хеш-секционирование и секционирование по списку приносят много заметных преимуществ, но они не так удобны, как секционирование по диапазонам значений, когда дело доходит до исключения секций при обработке запросов. Использование хеш-секционирования или секционирования по списку внутри диапазонных секций целесообразно, когда результирующие секции диапазонов оказываются слишком большими для управления или когда нужно задействовать все возможности PDML либо параллельного сканирования индексов в единственной диапазонной секции.

Секционирование индексов

Индексы, подобно таблицам, можно секционировать. Для этого доступны два метода.

- Эквисекционирование индекса с таблицей. Также называется локальным индексом. Для каждой табличной секции создается секция индекса, которая индексирует только эту табличную секцию. Все записи в отдельно взятой секции индекса указывают на единственную табличную секцию, а все строки в одной табличной секции представлены в единственной секции индекса.
- Секционирование индекса по диапазону или хеш-значению. Также называется глобально секционированным индексом. Здесь индекс секционируется по диапазону или дополнительно по хеш-значению (в Oracle 10g и последующих версиях), и единственная секция индекса может указывать на любую (и все) секции таблицы.

На рис. 13.5 демонстрируется разница между локальным и глобальным индексами.

Обратите внимание, что в случае глобально секционированного индекса количество индексных секций может отличаться от количества табличных секций.

Поскольку глобальные индексы могут секционироваться только по диапазону или хеш-значению, вы должны применять локальные индексы, если нужен индекс с секционированием по списку или составным секционированием. Локальный индекс будет секционирован по той же схеме, что и таблица, к которой он относится.

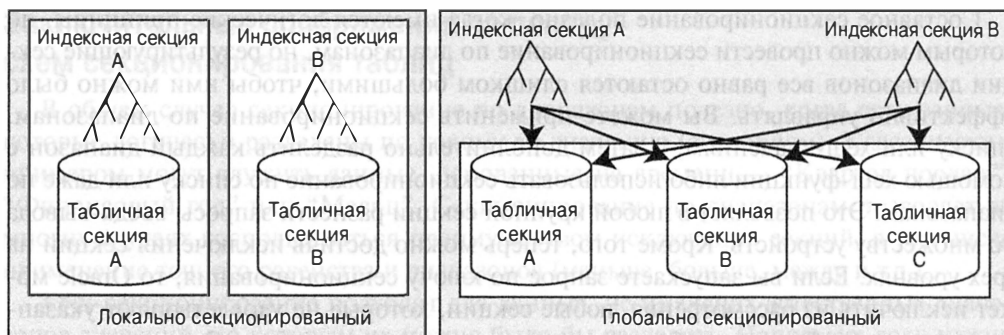


Рис. 13.5. Локальные и глобальные индексные секции

Сравнение локальных и глобальных индексов

Согласно моему опыту, большинство реализаций секционирования в системах хранилищ данных используют локальные индексы. В системе OLTP более распространены глобальные индексы, и вскоре мы увидим почему. Дело в том, что это связано с необходимостью выполнения исключения секций в индексных структурах, чтобы поддерживать время отклика после секционирования таким же, каким оно было до секционирования.

На заметку! За последние несколько лет все чаще можно наблюдать применение локальных индексов в системах OLTP, т.к. эти системы быстро растут в размерах.

Локальные индексы обладают определенными свойствами, которые делают их лучшим выбором для большинства реализаций хранилищ данных. Они поддерживают среды с более высокой доступностью (меньшим временем простоя), поскольку проблемы будут изолированы одним диапазоном или хеш-значением данных. С другой стороны, т.к. глобальный индекс способен указывать на множество табличных секций, он может стать точкой отказа, делая все секции недоступными некоторым запросам.

Локальные индексы являются более гибкими, когда речь идет об операциях обслуживания секций. Если администратор базы данных решит переместить табличную секцию, то перестройки или обслуживания потребует только ассоциированная с ней секция локального индекса. В глобальном индексе все индексные секции должны перестраиваться или обслуживаться в реальном времени. То же самое справедливо и в отношении реализаций скользящих окон, когда старые данные изымаются из секции, а взамен поступают новые данные. Локальным индексам перестройка не понадобится, но все глобальные индексы будут либо перестраиваться, либо обслуживаться во время выполнения операции с секциями. В некоторых случаях Oracle может извлечь преимущество из того факта, что индекс локально секционирован с таблицей, и на основе этого разрабатывать оптимизированные планы запросов. В случае глобальных индексов такое отношение между секцией индекса и секцией таблицы отсутствует.

Локальные индексы также облегчают выполнение операции восстановления секций на определенный момент времени. Если по какой-то причине одиночная секция должна быть восстановлена в состоянии на более ранний момент времени, чем остальная таблица, все локально секционированные индексы могут быть восстановлены на тот же момент времени. Все глобальные индексы для этого объекта потребуется перестроить. Это вовсе не означает, что стоит избегать глобальных индексов; на самом деле они жизненно важны по причинам, связанным с производительностью, о чем вы скоро узнаете — просто необходимо помнить о последствиях их использования.

Локальные индексы

В Oracle делается различие между следующими двумя типами локальных индексов.

- Локальные префиксные индексы (local prefixed index). Это индексы, в которых ключи секционирования находятся в головной части определения индекса. Например, если таблица секционирована по диапазонам с применением столбца `LOAD_DATE`, то локальный префиксный индекс должен иметь `LOAD_DATE` в качестве первого столбца в своем списке столбцов.
- Локальные непрефиксные индексы (local nonprefixed index). Эти индексы не содержат ключ секционирования в головной части своего списка столбцов. Индекс может содержать, а может и не содержать столбцов ключа секционирования.

Оба типа индексов обладают преимуществом исключения секций, оба могут поддерживать уникальность (если в непрефиксном индексе присутствует ключ секционирования) и т.д. Фактически запрос, использующий локальный префиксный индекс, всегда *позволяет* исключать индексные секции при выполнении запросов, в то время как запрос, в котором применяется локальный префиксный индекс, может делать это не всегда. Вот почему локальные непрефиксные индексы некоторые считают медленными — они не *навязывают* исключение секций (хотя и поддерживают его).

В локальном префиксном индексе нет ничего по своему существу лучшего по сравнению с локальным непрефиксным индексом, когда индекс используется как первоначальный путь к таблице в запросе. Я имею в виду, что если запрос может начинаться со сканирования индекса в качестве первого шага, то между префиксным и непрефиксным индексами особой разницы нет.

Поведение исключения секций

Для запроса, который начинается с доступа к индексу, наличие возможности исключения секций из рассмотрения на самом деле всецело зависит от предиката внутри запроса. Продемонстрировать это поможет небольшой пример. Приведенный ниже код создает таблицу `PARTITIONED_TABLE`, которая секционирована по диапазонам на числовом столбце `A`, так что значения меньше 2 попадают в секцию `PART_1`, а значения меньше 3 — в секцию `PART_2`:

```

EODA@ORA12CR1> CREATE TABLE partitioned_table
2  ( a int,
3    b int,
4    data char(20)
5  )

```



```

6 PARTITION BY RANGE (a)
7 (
8 PARTITION part_1 VALUES LESS THAN(2) tablespace p1,
9 PARTITION part_2 VALUES LESS THAN(3) tablespace p2
10 )
11 /

```

Table created.

Таблица создана.

Затем мы создаем локальный префиксный индекс, LOCAL_PREFIXED, и локальный непрефиксный индекс, LOCAL_NONPREFIXED. Обратите внимание, что непрефиксный индекс не содержит столбец А в головной части своего определения, что и делает его непрефиксным:

```
EODA@ORA12CR1> create index local_prefixed on partitioned_table (a,b) local;
```

Index created.

Индекс создан.

```
EODA@ORA12CR1> create index local_nonprefixed on partitioned_table (b) local;
```

Index created.

Индекс создан.

Далее мы вставим данные в одну секцию и соберем статистику:

```
EODA@ORA12CR1> insert into partitioned_table
```

```
2 select mod(rownum-1,2)+1, rownum, 'x'
```

```
3 from dual connect by level <= 70000;
```

70000 rows created.

70000 строк создано.

```
EODA@ORA12CR1> begin
```

```
2 dbms_stats.gather_table_stats
```

```
3 ( user,
```

```
4 'PARTITIONED_TABLE',
```

```
5 cascade=>TRUE );
```

```
6 end;
```

```
7 /
```

PL/SQL procedure successfully completed.

Процедура PL/SQL успешно завершена.

Теперь переведем в автономный режим табличное пространство P2, в котором находится секция PART_2 для таблицы и индекса:

```
EODA@ORA12CR1> alter tablespace p2 offline;
```

Tablespace altered.

Табличное пространство изменено.

Перевод табличного пространства P2 в автономный режим предотвратит доступ Oracle к специфическим индексным секциям. Это будет имитировать аварийный отказ носителя, который вызвал их недоступность. Теперь выполним запрос к таблице, чтобы посмотреть, какие индексные секции понадобятся разным запросам. Первый запрос реализован так, чтобы обеспечить применение локального префиксного индекса:

```
EODA@ORA12CR1> select * from partitioned_table where a = 1 and b = 1;
```

| A | B | DATA |
|---|---|------|
| 1 | 1 | x |

Этот запрос завершился успешно, и мы можем увидеть почему, просмотрев план выполнения. Мы будем использовать встроенный пакет DBMS_XPLAN, который позволит увидеть, к каким секциям обращается этот запрос. Столбцы PSTART (начало секции) и PSTOP (конец секции) в выводе точно показывают, какие секции должны находиться в оперативном режиме и быть доступными, чтобы запрос смог нормально выполниться:

```
EODA@ORA12CR1> explain plan for select * from partitioned_table
where a = 1 and b = 1;
Explained.
Объяснено.
```

Теперь с помощью функции DBMS_XPLAN.DISPLAY отобразим базовый план выполнения и сведения о секционировании:

```
EODA@ORA12CR1> select * from table(dbms_xplan.display(null,null,
where 'BASIC +PARTITION'));
```

| Id | Operation | Name | Pstart | Pstop |
|----|---|-------------------|--------|-------|
| 0 | SELECT STATEMENT | | | |
| 1 | PARTITION RANGE SINGLE | | 1 | 1 |
| 2 | TABLE ACCESS BY LOCAL INDEX ROWID BATCHED | PARTITIONED_TABLE | 1 | 1 |
| 3 | INDEX RANGE SCAN | LOCAL_PREFIXED | 1 | 1 |

Итак, запрос, который применяет LOCAL_PREFIXED, выполняется успешно. Оптимизатор смог исключить из рассмотрения секцию PART_2 индекса LOCAL_PREFIXED, т.к. в запросе было указано A=1; в плане ясно видно, что в столбцах PSTART и PSTOP находится значение 1. Исключение секции произошло. Тем не менее, следующий запрос не заработает:

```
EODA@ORA12CR1> select * from partitioned_table where b = 1;
ERROR:
ORA-00376: file 10 cannot be read at this time
ORA-01110: data file 10: '/u01/dbfile/ORA12CR1/datafile/o1_mf_p2_9hstdql2_.dbf'
no rows selected
ORA-00376: файл 10 не может быть прочитан в этот момент
ORA-01110: файл данных 10: '/u01/dbfile/ORA12CR1/datafile/o1_mf_p2_9hstdql2_.dbf'
нет выбранных строк
```

Используя тот же прием, выясним причину:

```
EODA@ORA12CR1> explain plan for select * from partitioned_table where b = 1;
Explained.
Объяснено.
```

```
EODA@ORA12CR1> select * from table(dbms_xplan.display(null,null,
where 'BASIC +PARTITION'));
```

| Id | Operation | Name | Pstart | Pstop |
|----|---|-------------------|--------|-------|
| 0 | SELECT STATEMENT | | | |
| 1 | PARTITION RANGE ALL | | 1 | 2 |
| 2 | TABLE ACCESS BY LOCAL INDEX ROWID BATCHED | PARTITIONED_TABLE | 1 | 2 |
| 3 | INDEX RANGE SCAN | LOCAL_NONPREFIXED | 1 | 2 |

Здесь оптимизатор *не* смог исключить из рассмотрения секцию PART_2 индекса LOCAL_NONPREFIXED — необходимо было просмотреть обе индексных секции, PART_1 и PART_2, на предмет наличия в них данных, для которых B=1. В этом и состоит проблема с производительностью локальных непрефиксных индексов: они не применяют ключ секционирования в предикате, как это делает префиксный индекс. Дело не в том, что префиксные индексы лучше; просто для их использования вы должны применять запрос, который разрешает исключение секций.

Если удалить индекс LOCAL_PREFIXED и перезапустить успешный исходный запрос:

```
EOA@ORA12CR1> drop index local_prefixed;
```

```
Index dropped.
```

Индекс удален.

```
EOA@ORA12CR1> select * from partitioned_table where a = 1 and b = 1;
```

```

      A      B DATA
-----
      1      1   x

```

то он выполнится нормально, но при этом использует тот же самый индекс, который ранее приводил к неудаче. План показывает, что Oracle теперь удалось задействовать исключение секций — в предикате A=1 было достаточно информации для базы данных, чтобы исключить из рассмотрения секцию PART_2:

```
EOA@ORA12CR1> explain plan for select * from partitioned_table
```

```
⌘ where a = 1 and b = 1;
```

```
Explained.
```

Объяснено.

```
EOA@ORA12CR1> select * from table(dbms_xplan.display(null,null,
```

```
⌘ 'BASIC +PARTITION'));
```

| Id | Operation | Name | Pstart | Pstop |
|----|---|-------------------|--------|-------|
| 0 | SELECT STATEMENT | | | |
| 1 | PARTITION RANGE SINGLE | | 1 | 1 |
| 2 | TABLE ACCESS BY LOCAL INDEX ROWID BATCHED | PARTITIONED_TABLE | 1 | 1 |
| 3 | INDEX RANGE SCAN | LOCAL_NONPREFIXED | 1 | 1 |

Обратите внимание на значения 1 и 1 в столбцах PSTART и PSTOP. Они подтверждают, что оптимизатор способен выполнять исключение секций даже для непрефиксных локальных индексов.

Если вы часто запускаете следующие запросы к показанной выше таблице, можете обдумать применение локального непрефиксного индекса на (b, a):

```
select ... from partitioned_table where a = :a and b = :b;
```

```
select ... from partitioned_table where b = :b;
```

Такой индекс был бы удобным для обоих запросов подобного рода. Локальный префиксный индекс на (a, b) оказался бы полезным только для первого запроса.

Суть здесь в том, что не стоит остерегаться непрефиксных индексов или считать их основными источниками снижения производительности. При наличии множества запросов, которые могут выиграть от непрефиксного индекса, как было показано выше, вы должны рассмотреть возможность его использования. Главное — удостове-

вериться, что запросы содержат предикаты, которые позволяют исключать индексные секции всякий раз, когда это возможно. Применение локальных префиксных индексов учитывает это соображение, а использование локальных непрефиксных индексов — нет. Принимайте также во внимание и то, каким образом индекс будет эксплуатироваться. Если он задействуется на первом шаге в плане запроса, то особых отличий между упомянутыми двумя типами индексов не будет.

Локальные индексы и ограничения уникальности

Чтобы обеспечить уникальность, а это подразумевает ограничение UNIQUE или ограничения PRIMARY KEY, ключ секционирования *должен входить в само ограничение*, если вы хотите применять это ограничение посредством локального индекса. По моему мнению, это является крупнейшим недостатком локального индекса. Уникальность в Oracle обеспечивается только внутри одной индексной секции — и никогда между секциями. Это подразумевает, например, что вы не можете провести секционирование по диапазонам на поле TIMESTAMP и иметь первичный ключ на поле ID, который поддерживается с использованием локально секционированного индекса. Взамен для обеспечения уникальности Oracle задействует *глобальный индекс*.

В следующем примере мы создадим таблицу, которая секционирована по диапазонам на столбце с именем TIMESTAMP, но имеет первичный ключ на столбце ID. Это можно сделать, запустив приведенный ниже оператор CREATE TABLE в схеме, которая не владеет какими-то другими объектами, что позволит увидеть, какие объекты создаются, просмотрев все принадлежащие пользователю сегменты:

```
EODA@ORA12CR1> CREATE TABLE partitioned
2 ( timestamp date,
3   id          int,
4   constraint partitioned_pk primary key(id)
5 )
6 PARTITION BY RANGE (timestamp)
7 (
8   PARTITION part_1 VALUES LESS THAN
9   ( to_date('01/01/2014','dd/mm/yyyy') ) ,
10  PARTITION part_2 VALUES LESS THAN
11  ( to_date('01/01/2015','dd/mm/yyyy') )
12 )
13 /
```

Table created.

Таблица создана.

Вставим некоторые данные, чтобы вызвать создание сегментов:

```
EODA@ORA12CR1> insert into partitioned values(to_date('01/01/2013',
  1 'dd/mm/yyyy'),1);
1 row created.
1 строка создана.
```

```
EODA@ORA12CR1> insert into partitioned values(to_date('01/01/2014',
  1 'dd/mm/yyyy'),2);
1 row created.
1 строка создана.
```

Предполагая, что другие объекты в схеме отсутствуют, мы получим следующий вывод:

```
EODA@ORA12CR1> select segment_name, partition_name,
  segment_type from user_segments;
```

| SEGMENT_NAME | PARTITION_NAME | SEGMENT_TYPE |
|----------------|----------------|-----------------|
| PARTITIONED | PART_1 | TABLE PARTITION |
| PARTITIONED | PART_2 | TABLE PARTITION |
| PARTITIONED_PK | | INDEX |

Индекс PARTITIONED_PK даже не секционирован, не говоря уже о локальном секционировании, и как мы увидим, он не может быть локально секционирован. Даже если мы попытаемся обмануть Oracle, явно указав, что первичный ключ может быть обеспечен неunikальным индексом так же, как уникальным, то обнаружим, что такой подход все равно не заработает:

```
EODA@ORA12CR1> CREATE TABLE partitioned
```

```
2 ( timestamp date,
3   id      int
4 )
5 PARTITION BY RANGE (timestamp)
6 (
7   PARTITION part_1 VALUES LESS THAN
8   ( to_date('01-jan-2014','dd-mon-yyyy') ) ,
9   PARTITION part_2 VALUES LESS THAN
10  ( to_date('01-jan-2015','dd-mon-yyyy') )
11 )
12 /
```

Table created.

Таблица создана.

```
EODA@ORA12CR1> create index partitioned_idx on partitioned(id) local;
```

Index created.

Индекс создан.

Вставим данные, чтобы создались сегменты:

```
EODA@ORA12CR1> insert into partitioned values(to_date('01/01/2013',
'dd/mm/yyyy'),1);
```

1 row created.

1 строка создана.

```
EODA@ORA12CR1> insert into partitioned values(to_date('01/01/2014',
'dd/mm/yyyy'),2);
```

1 row created

1 строка создана.

```
EODA@ORA12CR1> select segment_name, partition_name, segment_type
  2   from user_segments;
```

| SEGMENT_NAME | PARTITION_NAME | SEGMENT_TYPE |
|-----------------|----------------|-----------------|
| PARTITIONED | PART_1 | TABLE PARTITION |
| PARTITIONED | PART_2 | TABLE PARTITION |
| PARTITIONED_IDX | PART_1 | INDEX PARTITION |
| PARTITIONED_IDX | PART_2 | INDEX PARTITION |

```

EODA@ORA12CR1> alter table partitioned
2 add constraint
3 partitioned_pk
4 primary key(id)
5 /
alter table partitioned
*
ERROR at line 1:
ORA-01408: such column list already indexed
ОШИБКА в строке 1:
ORA-01408: такой список столбцов уже проиндексирован

```

Здесь Oracle пытается создать глобальный индекс на столбце ID, но обнаруживает, что это невозможно, потому что индекс уже существует. Предыдущие операторы заработали бы, если бы созданный нами индекс не был секционирован, т.к. Oracle приходится его применять, чтобы обеспечить ограничение.

Причина того, что уникальность не может быть обеспечена, если только ключ секционирования не является частью ограничения, двояка. Если бы в Oracle это было разрешено, большинство преимуществ секций исчезло бы. Доступность и масштабируемость были бы утрачены, поскольку каждая секция должна была быть всегда *доступной* и *сканироваться*, чтобы выполнялись любые вставки и обновления. Чем больше имелось бы секций, тем менее доступными становились бы данные. Чем больше появлялось бы секций таблицы, тем больше индексных секций приходилось бы сканировать и тем менее масштабируемыми становились бы секции. Вместо предоставления доступности и масштабируемости в действительности то и другое снизилось бы.

Кроме того, Oracle понадобилось бы эффективно *сериализовать* вставки и обновления в эту таблицу на транспортном уровне. Причина в том, что если мы добавляем строку с ID=1 в секцию PART_1, то Oracle пришлось бы *предотвратить* вставку кем-то другой строки с ID=1 в секцию PART_2. Единственный способ сделать это заключался бы в предотвращении модификации индексной секции PART_2 другими клиентами, т.к. нет ничего, что можно было бы действительно заблокировать в этой секции.

В системе OLTP ограничения уникальности должны обеспечиваться системой (т.е. Oracle) для поддержания целостности данных. Это подразумевает, что логическая модель вашего приложения будет оказывать влияние на физическое проектное решение. Ограничения уникальности либо будут управлять лежащей в основе схемой секционирования таблицы, определяя выбор ключей секционирования, либо склонять вас к использованию *глобальных* индексов. В следующем разделе глобальные индексы рассматриваются более подробно.

Глобальные индексы

Глобальные индексы секционируются с применением схемы, которая отличается от той, что используется в таблице, к которой они относятся. Таблица может быть разбита на десять секций по столбцу TIMESTAMP, а глобальный индекс для этой таблицы — на пять секций по столбцу REGION. В отличие от локальных индексов, существует только один класс глобальных индексов — *глобальный префиксный индекс*. Отсутствует поддержка глобального индекса, ключ которого не начинается с ключа секционирования *для этого индекса*. Из этого вытекает, что атрибуты, используемые для секционирования индекса, будут находиться в головной части самого ключа индекса.

Основываясь на предыдущем примере, рассмотрим краткий пример применения глобального индекса. Он покажет, что глобальный секционированный индекс может использоваться для обеспечения уникальности первичному ключу, поэтому можно иметь секционированные индексы, которые поддерживают уникальность, но не содержат ключ секционирования таблицы. Ниже создается таблица, секционированная по столбцу `TIMESTAMP`, которая имеет индекс, секционированный по столбцу `ID`:

```

ODA@ORA12CR1> CREATE TABLE partitioned
2 ( timestamp      date,
3   id              int
4 )
5 PARTITION BY RANGE (timestamp)
6 (
7   PARTITION part_1 VALUES LESS THAN
8   ( to_date('01-jan-2014','dd-mon-yyyy') ) ,
9   PARTITION part_2 VALUES LESS THAN
10  ( to_date('01-jan-2015','dd-mon-yyyy') )
11 )
12 /

```

Table created.

```

ODA@ORA12CR1> create index partitioned_index
2 on partitioned(id)
3 GLOBAL
4 partition by range(id)
5 (
6   partition part_1 values less than(1000),
7   partition part_2 values less than (MAXVALUE)
8 )
9 /

```

Index created.

Обратите внимание на применение `MAXVALUE` в этом индексе. Константа `MAXVALUE` может использоваться в любой секционированной по диапазонам таблице, а также в индексе. Она представляет *бесконечную верхнюю границу* диапазона. В примерах, приведенных до сих пор, применялись жесткие верхние границы диапазонов (значения меньше *некоторого значения*). Однако глобальный индекс подчиняется требованию того, что самая высокая (последняя) секция должна иметь границу, значение которой равно `MAXVALUE`. Это гарантирует, что все строки таблицы могут размещаться в индексе.

Теперь для завершения примера добавим к таблице первичный ключ:

```

ODA@ORA12CR1> alter table partitioned add constraint
2 partitioned_pk
3 primary key(id)
4 /

```

Table altered.

Таблица изменена.

Данный код не позволяет увидеть, что Oracle использует созданный нами индекс для поддержания первичного ключа, так что удостоверимся в этом, просто попробовав удалить индекс:

```

EODA@ORA12CR1> drop index partitioned_index;
drop index partitioned_index
*
ERROR at line 1:
ORA-02429: cannot drop index used for enforcement of unique/primary key
ОШИБКА в строке 1:
ORA-02429: невозможно удалить индекс, используемый для принудительного
применения уникального/первичного ключа

```

Чтобы показать, что Oracle не позволит создать глобальный *непрефиксный* индекс, достаточно попытаться выполнить следующий оператор:

```

EODA@ORA12CR1> create index partitioned_index2
2 on partitioned(timestamp,id)
3 GLOBAL
4 partition by range(id)
5 (
6 partition part_1 values less than(1000),
7 partition part_2 values less than (MAXVALUE)
8 )
9 /
partition by range(id)
*
ERROR at line 4:
ORA-14038: GLOBAL partitioned index must be prefixed
ОШИБКА в строке 4:
ORA-14038: ГЛОБАЛЬНЫЙ секционированный индекс должен быть префиксным

```

Это сообщение об ошибке предельно ясно. Глобальный индекс *должен* быть префиксным. Итак, когда же необходимо использовать глобальный индекс? Мы взглянем на два типа систем — хранилище данных и OLTP — и посмотрим, когда могут применяться глобальные индексы.

Хранилища данных и глобальные индексы

В прошлом хранилища данных и глобальные индексы практически взаимно исключали друг друга. Хранилище данных подразумевает определенные характеристики, такие как крупный объем входящих и исходящих данных. Многие хранилища данных реализуют подход скользящих окон для управления данными — т.е. удаляют самую старую секцию таблицы и добавляют новую секцию для вновь загружаемых данных. Ранее (в Oracle8i и предшествующих версиях) эти системы должны были избегать использования глобальных индексов по очень весомой причине: недостаток доступности. Это было вызвано тем, что большинство операций с секциями (вроде удаления старой секции) делали глобальные индексы *недействительными*, что предотвращало их применение вплоть до перестройки. Это серьезно снижало доступность.

В последующих разделах мы посмотрим, что понимается под *скользящим окном данных*, и оценим *потенциальное* влияние на него глобального индекса. Я подчеркиваю слово “потенциальное”, потому что речь также пойдет о том, как обойти эту проблему, и что такой обход может за собой повлечь.

Скользящие окна и индексы

В следующем примере реализуется классическое скользящее окно данных. Во многих реализациях данные добавляются в хранилище с течением времени, и наиболее старые данные из него изымаются. Часто эти данные секционируются по диапазону на основе какого-то атрибута типа даты, так что самые старые данные хранятся вместе в одной секции, и подобным же образом вновь загруженные данные попадают в новую секцию. Ежемесячный процесс загрузки включает в себя перечисленные ниже действия.

- Отсоединение старых данных. Наиболее старая секция либо удаляется, либо обменивается с пустой таблицей (данные этой секции перемещаются в таблицу), чтобы позволить архивирование старых данных.
- Загрузка и индексация новых данных. Новые данные загружаются в рабочую таблицу, индексируются и проверяются.
- Присоединение новых данных. После того как новые данные загружены и обработаны, таблица, в которой они находятся, обменивается с пустой секцией в секционированной таблице. Это приводит к переносу новых данных, загруженных в таблицу, в секцию крупной секционированной таблицы.

Такой процесс повторяется ежемесячно или с частотой выполнения процесса загрузки; он может происходить ежедневно или еженедельно. Здесь мы реализуем этот очень типичный процесс, чтобы продемонстрировать влияние глобальных секционированных индексов и показать возможности увеличения доступности данных во время выполнения операций с секциями, которые позволяют реализовать скользящее окно данных и поддерживать непрерывную доступность данных.

В этом примере мы будем обрабатывать ежегодные данные, располагая загруженной информацией о финансовых годах 2014 и 2015. Таблица будет секционирована по столбцу `TIMESTAMP` и иметь два индекса — локально секционированный индекс по столбцу `ID` и глобальный (не секционированный) индекс по столбцу `TIMESTAMP`:

```
EODA@ORA12CR1> CREATE TABLE partitioned
 2 ( timestamp date,
 3   id          int
 4 )
 5 PARTITION BY RANGE (timestamp)
 6 (
 7   PARTITION fy_2014 VALUES LESS THAN
 8   ( to_date('01-jan-2015','dd-mon-yyyy') ) ,
 9   PARTITION fy_2015 VALUES LESS THAN
10   ( to_date('01-jan-2016','dd-mon-yyyy') )
11 )
12 /
Table created.
```

```
EODA@ORA12CR1> insert into partitioned partition(fy_2014)
 2 select to_date('31-dec-2014','dd-mon-yyyy')-mod(rownum,360), rownum
 3 from dual connect by level <= 70000
 4 /
70000 rows created.
```

```

EODA@ORA12CR1> insert into partitioned partition(fy_2015)
2  select to_date('31-dec-2015','dd-mon-yyyy')-mod(rownum,360), rownum
3  from dual connect by level <= 70000
4  /
70000 rows created.
70000 строк создано.

EODA@ORA12CR1> create index partitioned_idx_local
2  on partitioned(id)
3  LOCAL
4  /
Index created.
Индекс создан.

EODA@ORA12CR1> create index partitioned_idx_global
2  on partitioned(timestamp)
3  GLOBAL
4  /
Index created.
Индекс создан.

```

Приведенный код настраивает таблицу хранилища. Данные секционированы по финансовому году и в оперативном режиме находится часть данных, относящихся к последним двум годам. Таблица имеет два индекса: LOCAL и GLOBAL. Близится конец года, и мы хотим выполнить следующие действия.

1. Удалить данные для самого старого финансового года. Мы не хотим потерять эти данные навсегда, а просто желаем их вывести из эксплуатации и заархивировать.
2. Добавить данные для нового финансового года. Их загрузка, трансформация и индексация займет некоторое время. Мы хотели бы, чтобы эта работа по возможности не повлияла на доступность текущих данных.

Первый шаг заключается в подготовке пустой таблицы для финансового года 2014, которая выглядит точно так же, как секционированная таблица. Мы будем использовать эту таблицу для обмена с секцией FY_2014 секционированной таблицы, перенеся в таблицу данные из указанной секции и затем опустошив секцию в секционированной таблице. В результате самые старые данные в секционированной таблице после обмена будут удалены:

```

EODA@ORA12CR1> create table fy_2014 ( timestamp date, id int );
Table created.
Таблица создана.

EODA@ORA12CR1> create index fy_2014_idx on fy_2014(id);
Index created.
Индекс создан.

```

Мы проделаем то же самое в отношении новых данных, подлежащих загрузке. Создадим и загрузим данными таблицу со структурой, аналогичной структуре существующей секционированной таблицы (но без секционирования):

```

EODA@ORA12CR1> create table fy_2016 ( timestamp date, id int );
Table created.
Таблица создана.

```

```

EODA@ORA12CR1> insert into fy_2016
2  select to_date('31-dec-2016','dd-mon-yyyy')-mod(rownum,360), rownum
3  from dual connect by level <= 70000
4  /
70000 rows created.
70000 строк создано.

EODA@ORA12CR1> create index fy_2016_idx on fy_2016(id) nologging;
Index created.
Индекс создан.

```

Мы опустошим текущую полную секцию и создадим таблицу, заполненную данными из секции FY_2014. Кроме того, мы завершим всю работу, необходимую для подготовки данных FY_2016. Это включает верификацию данных, их трансформацию — словом, все сложные задачи, которые должны выполняться для обеспечения готовности данных.

Теперь можно приступить к обновлению актуальных данных с применением обмена секциями:

```

EODA@ORA12CR1> alter table partitioned
2  exchange partition fy_2014
3  with table fy_2014
4  including indexes
5  without validation
6  /
Table altered.
Таблица изменена.

EODA@ORA12CR1> alter table partitioned drop partition fy_2014;
Table altered.
Таблица изменена.

```

Это и все, что понадобилось для исключения устаревших данных. Мы превратили секцию в полную таблицу, а пустую таблицу — в секцию. Это было простым обновлением словаря данных. Никакого существенного объема ввода-вывода не потребовалось. Теперь можно экспортировать таблицу FY_2014 (возможно, используя переносимое табличное пространство) из нашей базы данных для последующего архивирования. При необходимости ее можно будет быстро присоединить повторно.

Далее мы хотим внести новые данные:

```

EODA@ORA12CR1> alter table partitioned
2  add partition fy_2016
3  values less than ( to_date('01-jan-2017','dd-mon-yyyy') )
4  /
Table altered.
Таблица изменена.

EODA@ORA12CR1> alter table partitioned
2  exchange partition fy_2016
3  with table fy_2016
4  including indexes
5  without validation
6  /
Table altered.
Таблица изменена.

```

И снова все происходит мгновенно, т.к. достигается простыми обновлениями словаря данных — добиться этого позволяет конструкция `WITHOUT VALIDATION`. В случае применения такой конструкции база данных будет доверять помещаемым в эту секцию данным, считая их допустимыми. Добавление пустой секции занимает очень короткое время. Затем мы производим обмен новой пустой секции с полной таблицей, а полной таблицы — с пустой секцией, и эта операция также выполняется быстро. Новые данные находятся в оперативном режиме.

Тем не менее, взглянув на индексы, мы увидим вот что:

```
EODA@ORA12CR1> select index_name, status from user_indexes;
```

| INDEX_NAME | STATUS |
|------------------------|----------|
| PARTITIONED_IDX_LOCAL | N/A |
| PARTITIONED_IDX_GLOBAL | UNUSABLE |
| FY_2014_IDX | VALID |
| FY_2016_IDX | VALID |

После такой операции глобальный индекс, естественно, становится непригодным. Поскольку каждая секция индекса может указывать на любую секцию таблицы, и мы только что избавились от одной секции и добавили другую секцию, индекс стал недействительным. Он содержит записи, которые указывают на удаленную секцию. Он *не* имеет записей, которые указывали бы на только что добавленную секцию. Любой запрос, который использует этот индекс, потерпит неудачу и не выполнится, а если пропустить непригодные индексы, то производительность запроса пострадает из-за того, что индекс не будет задействован:

```
EODA@ORA12CR1> select /*+ index( partitioned PARTITIONED_IDX_GLOBAL ) */ count(*)
  2 from partitioned
  3 where timestamp between to_date( '01-mar-2016', 'dd-mon-yyyy' )
  4 and to_date( '31-mar-2016', 'dd-mon-yyyy' );
select /*+ index( partitioned PARTITIONED_IDX_GLOBAL ) */ count(*)
*
```

ERROR at line 1:

ORA-01502: index 'EODA.PARTITIONED_IDX_GLOBAL' or partition of such index is in unusable state

ОШИБКА в строке 1:

ORA-01502: индекс EODA.PARTITIONED_IDX_GLOBAL или секция этого индекса находится в непригодном состоянии

```
EODA@ORA12CR1> explain plan for select count(*)
```

```
  2 from partitioned
  3 where timestamp between to_date( '01-mar-2016', 'dd-mon-yyyy' )
  4 and to_date( '31-mar-2016', 'dd-mon-yyyy' );
```

```
EODA@ORA12CR1> select * from table(dbms_xplan.display(null,null,
  5 'BASIC +PARTITION'));
```

| Id | Operation | Name | Pstart | Pstop |
|----|------------------------|-------------|--------|-------|
| 0 | SELECT STATEMENT | | | |
| 1 | SORT AGGREGATE | | | |
| 2 | PARTITION RANGE SINGLE | | 2 | 2 |
| 3 | TABLE ACCESS FULL | PARTITIONED | 2 | 2 |

Таким образом, в отношении глобальных индексов после выполнения такой операции с секциями на выбор доступны следующие действия.

- Пропустить индекс либо прозрачным образом, как поступает Oracle в приведенном примере, либо путем установки параметра сеанса `SKIP_UNUSABLE_INDEXES=TRUE` в Oracle9i (в Oracle 10g он установлен в `TRUE` по умолчанию). Но тогда утрачивается рост производительности, обеспечиваемый индексом.
- Позволить запросам получить ошибку, как это будет с запросами без установки `SKIP_UNUSABLE_INDEXES` в `FALSE` в Oracle9i и предшествующих версиях, или с запросами, которые явно применяют подсказку в версии Oracle 10g. Чтобы снова нормально извлекать данные, индекс понадобится перестроить.

Процесс скользящего окна, при котором до сих пор практически отсутствовало время простоя, теперь будет требовать очень длительного времени на завершение, пока мы перестраиваем глобальный индекс. Производительность запросов во время выполнения, которые полагаются на эти индексы, в этот период пострадает — они либо вообще не запустятся, либо запустятся, но не воспользуются преимуществами индекса. Все данные должны сканироваться, а индекс полностью воссоздаваться из табличных данных. Если таблица имеет размер в сотни гигабайтов, это потребует значительных ресурсов.

Оперативное обслуживание глобальных индексов

Начиная с Oracle9i, к обслуживанию секций была добавлена еще одна функция: возможность *обслуживания глобальных индексов* во время проведения операции над секциями с применением конструкции `UPDATE GLOBAL INDEXES`. Это означает, что когда вы удаляете секцию, разделяете секцию или выполняете другую необходимую операцию над секцией, Oracle внесет любые необходимые модификации в глобальный индекс, чтобы поддержать его в актуальном состоянии. Поскольку большинство операций над секциями делают глобальный индекс недействительным, такое средство может оказаться благотворным для систем, которым необходимо предоставлять непрерывный доступ к данным. Вы обнаружите, что жертвуете высокой скоростью выполнения операций над секциями, но с сопутствующим периодом недоступности, немедленно возникающим в результате перестройки индексов, в пользу замедления этих операций при стопроцентной доступности данных. Короче говоря, если у вас есть хранилище данных, для которого простои неприемлемы, но которое должно поддерживать распространенные приемы поступления и изъятия данных, то описанное средство подойдет наилучшим образом, однако при этом следует понимать последствия.

Вернемся к предыдущему примеру, и там, где это уместно, используем в операциях над секциями конструкцию `UPDATE GLOBAL INDEXES` (здесь нет нужды в применении оператора `ADD PARTITION`, т.к. вновь добавленная секция не будет содержать строк). Мы увидим, что индексы полностью допустимы и пригодны к использованию *и во время, и после* операции:

```

EODA@ORA12CR1> alter table partitioned
2  exchange partition fy_2014
3  with table fy_2014
4  including indexes
5  without validation
6  UPDATE GLOBAL INDEXES
7  /
Table altered.

```

```
EODA@ORA12CR1> alter table partitioned drop partition fy_2014
```

```
2 update global indexes;
```

```
Table altered.
```

Таблица изменена.

```
EODA@ORA12CR1> alter table partitioned
```

```
2 add partition fy_2016
```

```
3 values less than ( to_date('01-jan-2017','dd-mon-yyyy') )
```

```
4 /
```

```
Table altered.
```

Таблица изменена.

```
EODA@ORA12CR1> alter table partitioned
```

```
2 exchange partition fy_2016
```

```
3 with table fy_2016
```

```
4 including indexes
```

```
5 without validation
```

```
6 UPDATE GLOBAL INDEXES
```

```
7 /
```

```
Table altered.
```

Таблица изменена.

В приведенном ниже выводе состояние N/A (Не применяется) для индекса PARTITIONED_IDX_LOCAL просто означает, что состояния ассоциированы с *секциями* индекса, а не самим индексом. Не имеет смысла говорить, допустим ли секционированный индекс; это всего лишь контейнер, который логически содержит в себе сами секции индекса:

```
EODA@ORA12CR1> select index_name, status from user_indexes;
```

| INDEX_NAME | STATUS |
|------------------------|--------|
| PARTITIONED_IDX_LOCAL | N/A |
| PARTITIONED_IDX_GLOBAL | VALID |
| FY_2014_IDX | VALID |
| FY_2016_IDX | VALID |

```
EODA@ORA12CR1> explain plan for select count(*)
```

```
2 from partitioned
```

```
3 where timestamp between to_date( '01-mar-2016', 'dd-mon-yyyy' )
```

```
4 and to_date( '31-mar-2016', 'dd-mon-yyyy' );
```

```
EODA@ORA12CR1> select * from table(dbms_xplan.display(null,null,  
  ↳ 'BASIC +PARTITION'));
```

| Id | Operation | Name |
|----|------------------|------------------------|
| 0 | SELECT STATEMENT | |
| 1 | SORT AGGREGATE | |
| 2 | INDEX RANGE SCAN | PARTITIONED_IDX_GLOBAL |

Но здесь присутствует компромисс: мы выполняем логический эквивалент операций DELETE и INSERT на структурах глобального индекса. Когда мы уничтожаем секцию, то должны удалить все записи глобального индекса, которые могут указывать на эту секцию. Когда мы производим обмен таблицы с секцией, то должны удалить все записи глобального индекса, указывающие на исходные данные, а затем

вставить записи, которые указывают на новые данные. Таким образом, объем работы, выполняемой командами ALTER, значительно возрастает.

Согласно соображениям по обслуживанию глобальных индексов, следует ожидать, что подход без обслуживания индексов потребует меньше ресурсов базы данных и потому выполняется быстрее, но влечет за собой заметный период простоя. Второй подход, предусматривающий обслуживание индексов, будет потреблять больше ресурсов и, возможно, займет больше времени, но зато он позволит исключить простой. С точки зрения конечных пользователей их работа никогда не прекращается. Обработка может чуть замедлиться (т.к. происходит конкуренция с ними за ресурсы), но пользователи *продолжат работать безо всяких остановок*.

Подход с перестройкой индекса почти наверняка будет работать быстрее, принимая во внимание затраченное время и время процессора. Этот факт может вызвать у многих администраторов баз данных вывод: “Я не хочу использовать конструкцию UPDATE GLOBAL INDEXES — она медленная”. Тем не менее, это слишком упрощенный взгляд. Важно помнить, что хотя операция в целом занимает больше времени, обработка в системе не обязательно прерывается. Конечно, вам, как администратору базы данных, скорее всего, придется дольше смотреть на экран, но действительно важная работа, происходящая в системе, по-прежнему будет продолжаться. Вам необходимо оценить, насколько оправдан такой компромисс в вашей системе. Если в вашем распоряжении есть восьмичасовое окно на обслуживание по ночам, в рамках которого загружаются новые данные, то, конечно же, применяйте подход с перестройкой индексов, если он имеет смысл. Однако если вы обязаны обеспечивать непрерывную доступность, то возможность обслуживания глобальных индексов во время операций будет решающей.

Еще один аспект, который следует учитывать — объем информации повторения (redo), генерируемой при каждом подходе. Вы обнаружите, что при наличии конструкции UPDATE GLOBAL INDEXES генерируется заметно больше информации повторения (из-за обслуживания индекса), и с добавлением к таблице дополнительных глобальных индексов ее объем будет только расти. Информация повторения, сгенерированная при обработке UPDATE GLOBAL INDEXES, неизбежна и не может быть отключена через NOLOGGING, т.к. обслуживание глобальных индексов — это не полная перестройка их структуры, а скорее инкрементный процесс. Кроме того, поскольку обслуживается действующая структура индекса, для этого должна генерироваться информация повторения — в случае отказа операции с секцией вы должны быть готовы к возвращению индекса в состояние, в котором он находился до операции. И помните, что информация отмены (undo) также защищается информацией повторения, поэтому часть информации повторения, которую мы здесь видим, сгенерирована обновлениями индекса, а часть предназначена для отката. Добавьте один или два дополнительных глобальных индекса — и вы получите ожидаемый рост этих показателей.

Итак, конструкция UPDATE GLOBAL INDEXES позволяет достичь компромисса между доступностью и потреблением ресурсов. Если вам необходимо обеспечить непрерывную доступность, то этот вариант как раз подойдет. Но вы должны понимать последствия и влияние на другие компоненты системы. В частности, со временем во многих хранилищах данных ловко внедряется использование прямых операций массовой обработки, которые обходят генерацию данных отмены, а когда это разрешено, то и генерацию данных повторения.

Применение конструкции USING GLOBAL INDEXES не позволяет обойти ни то, ни другое. Прежде чем использовать это средство, вам придется исследовать правила, применяемые для регулирования размеров данных отмены и повторения, и убедиться в том, что оно сможет работать в системе.

Асинхронное обслуживание глобальных индексов

Как было показано в предыдущем разделе, в Oracle9i и последующих версиях с помощью конструкции UPDATE GLOBAL INDEXES вы можете обслуживать глобальные индексы во время удаления или усечения секций. Тем не менее, вы узнали, что в этом случае такие операции будут выполняться дольше и потреблять больше ресурсов.

Начиная с версии Oracle 12c, при удалении или усечении табличных секций Oracle откладывает удаление записей глобального индекса, связанных с удаляемыми или усекаемыми секциями. Такая процедура называется *асинхронным обслуживанием глобальных индексов*. Обслуживание глобального индекса переносится на будущее время, а глобальный индекс остается пригодным для эксплуатации. Идея заключается в том, что это улучшает производительность удаления/усечения секций, одновременно поддерживая любые глобальные индексы в рабочем состоянии. Фактическая очистка записей индексов производится позже (асинхронно) либо администратором базы данных, либо автоматически запланированным заданием Oracle. Дело вовсе не в том, что при этом выполняется меньше работы, а в том, что очистка записей индекса отвязывается от оператора DROP/TRUNCATE.

Продемонстрируем асинхронное обслуживание глобального индекса на небольшом примере. Для его подготовки мы создадим таблицу в базе данных Oracle 11g, заполним ее тестовыми данными и построим глобальный индекс:

```
EODA@ORA11GR2> CREATE TABLE partitioned
  2   ( timestamp   date,
  3     id          int
  4   )
  5   PARTITION BY RANGE (timestamp)
  6   (PARTITION fy_2014 VALUES LESS THAN
  7   (to_date('01-jan-2015','dd-mon-yyyy')),
  8   PARTITION fy_2015 VALUES LESS THAN
  9   ( to_date('01-jan-2016','dd-mon-yyyy')));

EODA@ORA11GR2> insert into partitioned partition(fy_2014)
  2   select to_date('31-dec-2014','dd-mon-yyyy')-mod(rownum,364), rownum
  3   from dual connect by level < 100000;
99999 rows created.
99999 строк создано.

EODA@ORA11GR2> insert into partitioned partition(fy_2015)
  2   select to_date('31-dec-2015','dd-mon-yyyy')-mod(rownum,364), rownum
  3   from dual connect by level < 100000;
99999 rows created.
99999 строк создано.

EODA@ORA11GR2> create index partitioned_idx_global
  2   on partitioned(timestamp)
  3   GLOBAL;
Index created.
Индекс создан.
```


Далее мы запустим запрос, извлекающий значения статистических показателей redo size (размер данных повторения) и db block gets (количество извлечений блоков базы данных) для текущего сеанса:

```
EODA@ORA11GR2> col r1 new_value r2
EODA@ORA11GR2> col b1 new_value b2
EODA@ORA11GR2> select * from
  2 (select b.value r1
  3   from v$statname a, v$mystat b
  4   where a.statistic# = b.statistic#
  5   and a.name = 'redo size'),
  6 (select b.value b1
  7   from v$statname a, v$mystat b
  8   where a.statistic# = b.statistic#
  9   and a.name = 'db block gets');
```

| R1 | B1 |
|---------|------|
| 4816712 | 4512 |

Затем удалим секцию, указав конструкцию UPDATE GLOBAL INDEXES:

```
EODA@ORA11GR2> alter table partitioned drop partition fy_2014
❖update global indexes;
Table altered.
Таблица изменена.
```

Теперь подсчитаем объем сгенерированной информации повторения и количество блоков, к которым производился доступ:

```
EODA@ORA11GR2> select * from
  2 (select b.value - &r2 redo_gen
  3   from v$statname a, v$mystat b
  4   where a.statistic# = b.statistic#
  5   and a.name = 'redo size'),
  6 (select b.value - &b2 db_block_gets
  7   from v$statname a, v$mystat b
  8   where a.statistic# = b.statistic#
  9   and a.name = 'db block gets');
```

```
old 2: (select b.value - &r2 redo_gen
new 2: (select b.value -      4816712 redo_gen
old 6: (select b.value - &b2 db_block_gets
new 6: (select b.value -      4512 db_block_gets

REDO_GEN DB_BLOCK_GETS
-----
2459820      1495
```

Если тот же самый код выполнить в базе данных Oracle 12c, то мы получим следующие результаты во время удаления секции при указанной конструкции UPDATE GLOBAL INDEXES:

| REDO_GEN | DB_BLOCK_GETS |
|----------|---------------|
| 9872 | 43 |

По сравнению с Oracle 11g в базе данных Oracle 12c была сгенерирована только небольшая доля информации повторения и значительно меньшее число обращений к блокам. Причина в том, что Oracle не производит немедленное обслуживание индекса, удаляя из него записи, которые относятся к отброшенной секции. Вместо этого такие записи помечаются как *висячие* и будут позже очищены Oracle. В существовании висячих записей можно убедиться посредством приведенного ниже запроса:

```
EODA@ORA12CR1> select index_name, orphaned_entries, status from user_indexes
2 where table_name='PARTITIONED';
```

| INDEX_NAME | ORP | STATUS |
|------------------------|-----|--------|
| PARTITIONED_IDX_GLOBAL | YES | VALID |

Как происходит очистка висячих записей? В версии Oracle 12c предусмотрено автоматически запланированное задание PMO_DEFERRED_GIDX_MAINT_JOB, которое запускается во время ночного окна обслуживания. Если вы не хотите ждать, пока запустится это задание, можете очистить висячие записи вручную:

```
EODA@ORA12CR1> exec dbms_part.cleanup_gidx;
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.
```

При таком подходе вы можете выполнять операции вроде удаления и усечения секций и по-прежнему поддерживать глобальные индексы в рабочем состоянии безо всяких накладных расходов по очистке индексных записей как части операции удаления/усечения.

Совет. Дополнительные сведения об асинхронном обслуживании глобальных индексов ищите в примечании MOS под номером 1482264.1.

Системы OLTP и глобальные индексы

Для системы OLTP характерно частое выполнение множества небольших транзакций чтения и записи. В общем случае первостепенным является быстрый доступ к строке (или строкам). Жизненно важна и целостность данных. Также очень важна доступность.

В системах OLTP глобальные индексы оправданы во многих ситуациях. Табличные данные могут секционироваться только по одному ключу — одному набору столбцов. Однако вы можете нуждаться в доступе к данным многими отличающимися способами. Скажем, данные в таблице EMPLOYEE (сотрудники) могут быть секционированы по столбцу LOCATION (местоположение), но вам по-прежнему необходим быстрый доступ к данным EMPLOYEE по следующим столбцам.

- DEPARTMENT. Отделы географически рассредоточены. Между отделом и местоположением какие-либо отношения отсутствуют.
- EMPLOYEE_ID. Хотя идентификатор сотрудника будет определять местоположение, вы не хотите осуществлять поиск по EMPLOYEE_ID и LOCATION, поэтому исключение секций не может происходить для индексных секций. Кроме того, сам столбец EMPLOYEE_ID должен быть уникальным.
- JOB_TITLE. Между должностью (JOB_TITLE) и местоположением какие-либо отношения отсутствуют. В любом местоположении могут присутствовать все должности.

Потребность в доступе к данным EMPLOYEE по многим отличающимся ключам возникает в разных местах приложения, и первостепенное значение имеет скорость. В хранилище данных мы могли бы просто использовать локально секционированные индексы на этих ключах и применять параллельные сканирования диапазонов по индексам для быстрого сбора крупного объема данных. В таких случаях применять исключение индексных секций не обязательно. Тем не менее, в системе OLTP такое исключение необходимо использовать. Параллельные запросы для этих систем не подходят; мы должны предоставить соответствующие индексы. Следовательно, нам придется обеспечить применение глобальных индексов на определенных полях.

Ниже перечислены цели, которых мы хотим добиться:

- быстрый доступ;
- целостность данных;
- доступность.

Достичь указанных целей в системе OLTP могут помочь глобальные индексы. Скорее всего, мы не будем организовывать скользящие окна, отложив этот вопрос на некоторое время. Мы не будем разделять секции (если только не располагаем запланированным временем простоя), не станем перемещать данные и т.п. В целом операции, выполняемые в хранилище данных, в действующей системе OLTP не производятся.

Рассмотрим небольшой пример, в котором показано, каким образом можно достичь трех перечисленных выше целей посредством глобальных индексов. Я собираюсь использовать простые *односекционные* глобальные индексы, но результаты не будут отличаться от тех, которые были бы получены с глобальными индексами, имеющими множество секций (за исключением того факта, что с добавлением индексных секций доступность и управляемость должны *возрасти*). Мы начнем с создания табличных пространств P1, P2, P3 и P4, а затем создадим таблицу, секционированную по диапазонам на основе местоположения (LOC) в соответствии с правилами, которые предусматривают помещение всех значений LOC меньше 'C' в секцию P1, всех значений LOC меньше 'D' — в секцию P2 и т.д.:

```
EODA@ORA12CR1> create tablespace p1 datafile size 1m autoextend on next 1m;
Tablespace created.
```

Табличное пространство создано.

```
EODA@ORA12CR1> create tablespace p2 datafile size 1m autoextend on next 1m;
Tablespace created.
```

Табличное пространство создано.

```
EODA@ORA12CR1> create tablespace p3 datafile size 1m autoextend on next 1m;
Tablespace created.
```

Табличное пространство создано.

```
EODA@ORA12CR1> create tablespace p4 datafile size 1m autoextend on next 1m;
Tablespace created.
```

Табличное пространство создано.

```
EODA@ORA12CR1> create table emp
2  (EMPNO          NUMBER(4) NOT NULL,
3   ENAME          VARCHAR2(10),
4   JOB            VARCHAR2(9),
5   MGR            NUMBER(4),
```

```

6  HIREDATE      DATE,
7  SAL           NUMBER(7,2),
8  COMM          NUMBER(7,2),
9  DEPTNO        NUMBER(2) NOT NULL,
10 LOC           VARCHAR2(13) NOT NULL
11 )
12 partition by range(loc)
13 (
14 partition p1 values less than('C') tablespace p1,
15 partition p2 values less than('D') tablespace p2,
16 partition p3 values less than('N') tablespace p3,
17 partition p4 values less than('Z') tablespace p4
18 )
19 /

```

Table created.

Таблица создана.

Изменим таблицу, добавив ограничение на столбце первичного ключа:

```

EODA@ORA12CR1> alter table emp add constraint emp_pk
2 primary key(empno)
3 /

```

Table altered.

Таблица изменена.

Побочным эффектом будет появление уникального индекса на столбце EMPNO. Это показывает, что мы можем поддерживать и обеспечивать соблюдение целостности данных, достигая одной из наших целей. Наконец, создадим еще два глобальных индекса на столбцах DEPTNO и JOB, чтобы содействовать быстрому доступу к записям по этим атрибутам:

```

EODA@ORA12CR1> create index emp_job_idx on emp(job)
2 GLOBAL
3 /

```

Index created.

Индекс создан.

```

EODA@ORA12CR1> create index emp_dept_idx on emp(deptno)
2 GLOBAL
3 /

```

Index created.

Индекс создан.

```

EODA@ORA12CR1> insert into emp
2 select e.*, d.loc
3 from scott.emp e, scott.dept d
4 where e.deptno = d.deptno
5 /

```

14 rows created.

14 строк создано.

Давайте посмотрим, что находится в каждой секции:

```

EODA@ORA12CR1> break on pname skip 1
EODA@ORA12CR1> select 'p1' pname, empno, job, loc from emp partition(p1)
2 union all
3 select 'p2' pname, empno, job, loc from emp partition(p2)

```

```
4 union all
5 select 'p3' pname, empno, job, loc from emp partition(p3)
6 union all
7 select 'p4' pname, empno, job, loc from emp partition(p4)
8 /
```

| PN | EMPNO | JOB | LOC |
|----|-------|-----------|----------|
| p2 | 7499 | SALESMAN | CHICAGO |
| | 7521 | SALESMAN | CHICAGO |
| | 7654 | SALESMAN | CHICAGO |
| | 7698 | MANAGER | CHICAGO |
| | 7844 | SALESMAN | CHICAGO |
| | 7900 | CLERK | CHICAGO |
| p3 | 7369 | CLERK | DALLAS |
| | 7566 | MANAGER | DALLAS |
| | 7788 | ANALYST | DALLAS |
| | 7876 | CLERK | DALLAS |
| | 7902 | ANALYST | DALLAS |
| p4 | 7782 | MANAGER | NEW YORK |
| | 7839 | PRESIDENT | NEW YORK |
| | 7934 | CLERK | NEW YORK |

14 rows selected.
14 строк выбрано.

Вывод показывает распределение данных по местоположению в индивидуальных секциях. Теперь мы можем просмотреть некоторые планы запросов, чтобы увидеть, чего следует ожидать в отношении производительности:

```
EODA@ORA12CR1> variable x varchar2(30);
EODA@ORA12CR1> begin
2   dbms_stats.set_table_stats
3   ( user, 'EMP', numrows=>100000, numblks => 10000 );
4 end;
5 /
```

PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.

```
EODA@ORA12CR1> explain plan for select empno, job, loc from emp
where empno = :x;
Explained.
Объяснено.
```

```
EODA@ORA12CR1> select * from table(dbms_xplan.display(null,null,
where empno = :x;
Explained.
Объяснено.
```

| Id | Operation | Name | Pstart | Pstop |
|----|------------------------------------|--------|--------|-------|
| 0 | SELECT STATEMENT | | | |
| 1 | TABLE ACCESS BY GLOBAL INDEX ROWID | EMP | ROWID | ROWID |
| 2 | INDEX UNIQUE SCAN | EMP_PK | | |

В плане присутствует операция INDEX UNIQUE SCAN по несекционированному индексу EMP_PK, который был создан для поддержки первичного ключа. Затем

идет операция TABLE ACCESS BY GLOBAL INDEX ROWID с PSTART и PSTOP, равными ROWID/ROWID, т.е., когда мы получаем идентификатор ROWID из индекса, он в точности сообщит, какую индексную секцию нужно прочитать, чтобы получить эту строку. Такой индексный доступ будет эффективным в той же степени, что и для несекционированной таблицы, и потребует того же самого объема ввода-вывода. Это лишь простое одиночное сканирование уникального индекса, за которым следует “получение строки по идентификатору ROWID”. Теперь давайте взглянем на еще один глобальный индекс — по столбцу JOB:

```
EOA@ORA12CR1> explain plan for select empno, job, loc from emp where job = :x;
Explained.
```

Объяснено.

```
EOA@ORA12CR1> select * from table(dbms_xplan.display);
```

| Id | Operation | Name | Pstart | Pstop |
|----|--|-------------|--------|-------|
| 0 | SELECT STATEMENT | | | |
| 1 | TABLE ACCESS BY GLOBAL INDEX ROWID BATCHED | EMP | ROWID | ROWID |
| 2 | INDEX RANGE SCAN | EMP_JOB_IDX | | |

Без сомнений мы видим похожий результат для операции INDEX RANGE SCAN. Наши индексы задействованы и могут предоставить высокоскоростной OLTP-доступ к лежащим в основе данным. Если бы они были секционированы, то должны были быть префиксными и обеспечивать исключение индексных секций; отсюда следует их масштабируемость — в том смысле, что мы можем их секционировать и наблюдать идентичное поведение. Вскоре мы посмотрим, что произошло бы в случае применения только индексов LOCAL.

Наконец, давайте обратимся к области, связанной с доступностью. В документации по Oracle утверждается, что глобально секционированные индексы обуславливают менее высокую доступность данных, чем локально секционированные индексы. Я не вполне согласен с такой поверхностной характеристикой. Я уверен, что в системе OLTP они предлагают такую же высокую доступность, как и локально секционированные индексы. Рассмотрим следующие действия:

```
EOA@ORA12CR1> alter tablespace p1 offline;
```

```
Tablespace altered.
```

Табличное пространство изменено.

```
EOA@ORA12CR1> alter tablespace p2 offline;
```

```
Tablespace altered.
```

Табличное пространство изменено.

```
EOA@ORA12CR1> alter tablespace p3 offline;
```

```
Tablespace altered.
```

Табличное пространство изменено.

```
EOA@ORA12CR1> select empno, job, loc from emp where empno = 7782;
```

| EMPNO | JOB | LOC |
|-------|---------|----------|
| 7782 | MANAGER | NEW YORK |

Несмотря на то что здесь большая часть данных в таблице является недоступной, мы по-прежнему можем обращаться к любой порции доступных данных через

индекс. До тех пор, пока данные с интересующим номером EMPNO находятся в доступном табличном пространстве, и доступен индекс GLOBAL, этот индекс работает на нас. С другой стороны, если бы в предыдущем случае мы *действительно* использовали высокодоступный локальный индекс, то вообще могли бы не иметь доступа к данным! Это побочный эффект того факта, что данные секционируются по LOC, но должны запрашиваться по EMPNO. Нам пришлось бы зондировать каждую секцию локального индекса и получать отказ в индексных секциях, которые не были доступными в тот момент.

Тем не менее, другие типы запросов не будут (и не могут) функционировать в этот момент времени:

```

EODA@ORA12CR1> select empno, job, loc from emp where job = 'CLERK';
select empno, job, loc from emp where job = 'CLERK'
*
ERROR at line 1:
ORA-00376: file 10 cannot be read at this time
ORA-01110: data file 10: '/u01/dbfile/ORA12CR1/datafile/ol_mf_p2_9hx10fqv_.dbf'
ORA-00376: файл 10 не может быть прочитан в этот момент
ORA-01110: файл данных 10: '/u01/dbfile/ORA12CR1/datafile/ol_mf_p2_9hx10fqv_.dbf'

```

Данные CLERK присутствуют во всех секциях, и тот факт, что три табличных пространства отключены, оказывает влияние. Это неизбежно, если только не провести секционирование по столбцу JOB, но тогда мы получаем те же самые проблемы с запросами, которые извлекают данные по LOC. Такая ситуация будет возникать всегда, когда необходим доступ к данным по множеству разных *ключей*. База данных Oracle будет предоставлять данные всякий раз, когда сможет.

Однако обратите внимание, что если ответ на запрос можно получить из индекса, избегая операции TABLE ACCESS BY ROWID, то факт недоступности табличных данных теряет свою важность:

```

EODA@ORA12CR1> select count(*) from emp where job = 'CLERK';
COUNT (*)
-----
4

```

Так как в этом случае Oracle не нуждается в таблице, факт пребывания большинства секций в автономном режиме не влияет на выполнение запроса (разумеется, предполагая, что индекс не расположен в одном из автономных табличных пространств). Поскольку такой вид оптимизации (т.е. формирование ответа на запрос с применением только индекса) распространен в системах OLTP, найдется немало приложений, которые никак не затронет то, что данные находятся в автономном режиме. Все, что потребуется предпринять — это сделать автономные данные доступными как можно быстрее (восстановить их).

Частичные индексы

Начиная с версии Oracle 12c, вы можете создавать либо локальные, либо глобальные индексы на подмножестве разделов в таблице. Это может понадобиться, если вы имеете заранее созданные секции и пока что не располагаете данными для диапозонных секций, которые соответствуют будущим датам — идея в том, что вы будете строить индекс после загрузки данных в секции (в какой-то момент в будущем).

Чтобы настроить использование частичного индекса, сначала необходимо указать конструкцию INDEXING ON|OFF для каждой секции в таблице. В следующем примере индексация секции PART_1 включена, а индексация секции PART_2 отключена:

```

EODA@ORA12CR1> CREATE TABLE p_table (a int)
2 PARTITION BY RANGE (a)
3 (PARTITION part_1 VALUES LESS THAN(1000) INDEXING ON,
4 PARTITION part_2 VALUES LESS THAN(2000) INDEXING OFF);
Table created.
Таблица создана.

```

Затем создается частичный локальный индекс:

```

EODA@ORA12CR1> create index pil on p_table(a) local indexing partial;
Index created.
Индекс создан.

```

В этом сценарии конструкция INDEXING PARTIAL инструктирует Oracle о том, что должны быть построены и сделаны пригодными к употреблению только такие индексные секции, которые указывают на секции в таблице, определенные с конструкцией INDEXING ON. В данном случае будет создана одна индексная секция с записями, указывающими на данные в табличной секции PART_1:

```

EODA@ORA12CR1> select a.index_name, a.partition_name, a.status
2 from user_ind_partitions a, user_indexes b
3 where b.table_name = 'P_TABLE'
4 and a.index_name = b.index_name;

```

| INDEX_NAME | ARTITION_NAME | STATUS |
|------------|---------------|----------|
| PI1 | PART_2 | UNUSABLE |
| PI1 | PART_1 | USABLE |

Далее мы вставим тестовые данные, сгенерируем статистику, включим средство AUTOTRACE и запустим запрос, который должен определить местоположение данных в секции PART_1:

```

EODA@ORA12CR1> insert into p_table select rownum from dual connect
by level < 2000;
1999 rows created.
1999 строк создано.

EODA@ORA12CR1> exec dbms_stats.gather_table_stats(user, 'P_TABLE');
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.

EODA@ORA12CR1> explain plan for select * from p_table where a = 20;
Explained.
Объяснено.

EODA@ORA12CR1> select * from table(dbms_xplan.display(null,null,
by 'BASIC +PARTITION'));

```

| Id | Operation | Name | Pstart | Pstop |
|----|------------------------|------|--------|-------|
| 0 | SELECT STATEMENT | | | |
| 1 | PARTITION RANGE SINGLE | | 1 | 1 |
| 2 | INDEX RANGE SCAN | PI1 | 1 | 1 |

Как и ожидалось, оптимизатор оказался способным сгенерировать план выполнения, который задействует индекс. После этого выдается запрос, извлекающий данные из секции, которая была определена с конструкцией INDEXING OFF:

```
EOGA@ORA12CR1> explain plan for select * from p_table where a = 1500;
Explained.
Объяснено.
```

```
EOGA@ORA12CR1> select * from table(dbms_xplan.display(null,null,
  ↳ 'BASIC +PARTITION'));
```

| Id | Operation | Name | Pstart | Pstop |
|----|------------------------|---------|--------|-------|
| 0 | SELECT STATEMENT | | | |
| 1 | PARTITION RANGE SINGLE | | 2 | 2 |
| 2 | TABLE ACCESS FULL | P_TABLE | 2 | 2 |

Вывод показывает, что для секции PART_2 требовалось полное сканирование таблицы, т.к. нет рабочего индекса с записями, указывающими на данные в PART_2. Мы можем сообщить Oracle о необходимости создания индексных записей, которые указывают на данные в секции PART_2, путем перестройки секции индекса, связанной с секцией PART_2:

```
EOGA@ORA12CR1> alter index pil rebuild partition part_2;
Index altered.
Индекс изменен.
```

Повторное выполнение предыдущего запроса говорит о том, что оптимизатор теперь применяет локальный секционированный индекс, указывающий на табличную секцию PART_2:

| Id | Operation | Name | Pstart | Pstop |
|----|------------------------|------|--------|-------|
| 0 | SELECT STATEMENT | | | |
| 1 | PARTITION RANGE SINGLE | | 2 | 2 |
| 2 | INDEX RANGE SCAN | PI1 | 2 | 2 |

Таким образом, частичные индексы позволяют отключать индекс на время загрузки данных в секцию таблицы (в итоге увеличивая скорость загрузки), а затем позже перестроить частичный индекс, чтобы сделать его доступным.

Еще раз о секционировании и производительности

Много раз мне приходилось слышать, как люди говорят: “Мы очень разочарованы в секционировании. Мы секционировали нашу самую большую таблицу, и она стала работать намного медленнее. Поэтому называть секционирование средством увеличения производительности — явное преувеличение”.

Секционирование может подействовать на общую производительность запросов одним из трех способов.

- Ускорить выполнение запросов.
- Вообще не повлиять на производительность запросов.

- Существенно замедлить выполнение запросов и задействовать во много раз больше ресурсов, чем в случае реализации без секционирования.

При понимании всех запросов, отправляемых к данным, первый результат очень часто достижим в хранилище данных. Секционирование может положительно повлиять на запросы, которые регулярно производят полное сканирование таблиц базы данных, исключая крупные фрагменты данных из рассмотрения. Предположим, что имеется таблица с одним миллиардом строк. В таблице есть атрибут — отметка времени. Ваш запрос собирается извлечь из этой таблицы данные за один год (всего в ней хранятся данные за 10 лет). Для извлечения этих данных запрос использует полное сканирование таблицы. Если таблица будет секционирована по атрибуту отметки времени (скажем, секция на месяц), то придется выполнять полное сканирование только одной десятой части данных (при условии их равномерного распределения по годам). Исключение секций позволит не рассматривать остальные 90% данных. Ваш запрос, вероятно, станет выполняться быстрее.

Теперь возьмем похожую таблицу в системе OLTP. В приложениях такого типа вы никогда не будете извлекать 10% из миллиарда строк таблицы. Следовательно, значительное увеличение скорости, наблюдаемое в хранилище данных, в транзакционной системе попросту недостижимо. Здесь не делается работа того же рода, что в хранилищах данных, поэтому ожидать аналогичных улучшений производительности нереально. Таким образом, в системе OLTP первый из трех перечисленных выше результатов в основном недостижим, и вы не будете применять секционирование преимущественно для повышения производительности. А вот увеличение доступности вполне достижимо. Облегчение администрирования — весьма вероятно. Но в системе OLTP придется сильно постараться, чтобы достичь второго из трех перечисленных результатов: вообще не оказать влияния на производительность запросов — ни отрицательного, ни положительного. В очень многих случаях именно это и будет вашей целью — секционирование без влияния на время выполнения запросов.

Мне приходилось много раз наблюдать следующую картину. Команда разработчиков имеет дело с таблицей средних размеров, скажем, 100 миллионов строк. Количество 100 миллионов для них *звучит* как невероятно большое число (а пять или десять лет тому назад так и было, но время меняет все). Разработчики решают секционировать данные, но осмотр данных не выявляет логических атрибутов, наличие которых делало бы осмысленным секционирование по диапазонам. Подходящие для такого секционирования атрибуты отсутствуют. Аналогично не имеет смысла и секционирование по списку. В таблице не обнаружено ничего, что позволило бы по нему правильно секционировать. Таким образом, разработчики выбирают вариант хеш-секционирования по первичному ключу, который по стечению обстоятельств заполняется значением последовательности Oracle. Первичный ключ выглядит безупречным, он уникален и легко хешируется, а многие запросы имеют вид `SELECT * FROM T WHERE PRIMARY_KEY = :X`.

Но проблема в том, что к этой таблице выполняется и множество других запросов, которые имеют не такую форму. В целях иллюстрации предположим, что обсуждаемая таблица — это на самом деле представление словаря данных `ALL_OBJECTS`, и в то время как внутренне многие запросы имеют форму `WHERE OBJECT_ID = :X`, конечные пользователи часто отправляют запросы следующих видов:

- показать подробные сведения о таблице EMP из схемы SCOTT (where owner = :o and object_type = :t and object_name = :n);
- показать все таблицы, которыми владеет схема SCOTT (where owner = :o and object_type = :t);
- показать все объекты, которыми владеет схема SCOTT (where owner = :o).

Для поддержки этих запросов предусмотрен индекс на (OWNER, OBJECT_TYPE, OBJECT_NAME). Но вы также читали, что локальные индексы являются более доступными, и хотите обеспечить более высокую доступность для системы, поэтому реализуете их. В конечном итоге вы воссоздадите свою таблицу с 16 хеш-секциями примерно так:

```
EODA@ORA12CR1> create table t
 2 ( OWNER, OBJECT_NAME, SUBOBJECT_NAME, OBJECT_ID, DATA_OBJECT_ID,
 3   OBJECT_TYPE, CREATED, LAST_DDL_TIME, TIMESTAMP, STATUS,
 4   TEMPORARY, GENERATED, SECONDARY )
 5 partition by hash(object_id)
 6 partitions 16
 7 as
 8 select OWNER, OBJECT_NAME, SUBOBJECT_NAME, OBJECT_ID, DATA_OBJECT_ID,
 9   OBJECT_TYPE, CREATED, LAST_DDL_TIME, TIMESTAMP, STATUS,
10   TEMPORARY, GENERATED, SECONDARY
11 from all_objects;
Table created.
Таблица создана.

EODA@ORA12CR1> create index t_idx
 2 on t(owner,object_type,object_name)
 3 LOCAL
 4 /
Index created.
Индекс создан.

EODA@ORA12CR1> exec dbms_stats.gather_table_stats( user, 'T' );
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.
```

Затем вы запускаете типичные запросы OLTP, которые, как вам известно, будут часто выполняться в системе:

```
variable o varchar2(30)
variable t varchar2(30)
variable n varchar2(30)

exec :o := 'SCOTT'; :t := 'TABLE'; :n := 'EMP';

select *
  from t
 where owner = :o
    and object_type = :t
    and object_name = :n
/

select *
  from t
 where owner = :o
    and object_type = :t
/
```

```

select *
  from t
 where owner = :o
/

```

Однако когда вы запускаете эти запросы с SQL_TRACE=TRUE и просматриваете результирующий отчет TKPROF, то замечаете следующие характеристики производительности:

```

select * from t where owner = :o and object_type = :t and object_name = :n
call      count      cpu  elapsed      disk  query  current    rows
-----
total         4      0.00    0.01         0     34         0         1
...
Rows (1st)  Rows (avg)  Rows (max)  Row Source Operation
-----
          1          1          1  PARTITION HASH ALL PARTITION: 1 16
⌚(cr=34 pr=0 pw=0 time=95...
          1          1          1  TABLE ACCESS BY LOCAL INDEX ROWID
⌚BATCHED T PARTITION: ...
          1          1          1  INDEX RANGE SCAN T_IDX PARTITION: 1 16
⌚(cr=33 pr=0 pw=0...

```

Вы сравниваете это с показателями, полученными для той же самой таблицы, но *без секционирования*, и обнаруживаете вот что:

```

select * from t where owner = :o and object_type = :t and object_name = :n
call      count      cpu  elapsed      disk  query  current    rows
-----
total         4      0.00    0.00         0         4         0         1
...
Rows (1st)  Rows (avg)  Rows (max)  Row Source Operation
-----
          1          1          1  TABLE ACCESS BY INDEX ROWID BATCHED T
⌚(cr=4 pr=0 pw=0...
          1          1          1  INDEX RANGE SCAN T_IDX (cr=3 pr=0 pw=0
⌚time=14 us cost=1...

```

Вы тут же можете сделать (ошибочный) вывод о том, что секционирование послужило причиной почти восьмикратного увеличения объема ввода-вывода: 4 извлечения в режиме запроса требуются при отсутствии секций и 34 — при их наличии. Если в вашей системе ранее присутствовала проблема с высоким количеством операций согласованного чтения (логических операций ввода-вывода), то теперь ситуация еще хуже. Если такой проблемы не было, то сейчас она может появиться. То же самое можно наблюдать с другими двумя запросами. В следующем отчете первая итоговая строка (total) касается секционированной таблицы, а вторая — несекционированной:

```

select * from t where owner = :o and object_type = :t
call      count      cpu  elapsed      disk  query  current    rows
-----
total         5      0.00    0.00         0     49         0     20
total         5      0.00    0.00         0     11         0     20

```

```
select * from t where owner = :o
```

| call | count | cpu | elapsed | disk | query | current | rows |
|-------|-------|------|---------|------|-------|---------|------|
| total | 5 | 0.00 | 0.00 | 0 | 665 | 0 | 26 |
| total | 5 | 0.00 | 0.00 | 0 | 628 | 0 | 26 |

Каждый запрос возвращает один и тот же ответ, но потребляет значительно больше операций ввода-вывода, что не очень хорошо. В чем основная причина? В схеме секционирования индекса. В показанном ранее плане обратите внимание на секции, перечисленные в последней строке — с 1 по 16:

```

1          1          1 PARTITION HASH ALL PARTITION: 1 16 (cr=34 pr=0
↳pw=0 time=95...
1          1          1 TABLE ACCESS BY LOCAL INDEX ROWID BATCHED T
↳PARTITION: ...
1          1          1 INDEX RANGE SCAN T_IDX PARTITION: 1 16
↳(cr=33 pr=0 pw=0...
```

Данный запрос должен просматривать *каждую* секцию индекса, потому что записи для SCOTT вполне могут оказаться в *каждой* секции индекса и, скорее всего, так и будет. Индекс логически хеш-секционирован по столбцу OBJECT_ID; любой запрос, который использует этот индекс, но также не ссылается на столбец OBJECT_ID в предикате, должен просматривать *каждую* секцию индекса! Каково же решение? Оно заключается в глобальном секционировании индекса. Взяв в качестве примера предыдущий случай, мы можем выбрать хеш-секционирование индекса:

На заметку! Хеш-секционирование индексов — это средство, появившееся в Oracle 10g, и в версии Oracle9i оно недоступно. Существует ряд соображений относительно сканирования диапазонов, которые должны быть приняты во внимание при работе с хеш-секционированными индексами; мы обсудим их позже в этом разделе.

```

EODA@ORA12CR1> create index t_idx
2 on t(owner,object_type,object_name)
3 global
4 partition by hash(owner)
5 partitions 16
6 /
Index created.
Индекс создан.
```

Во многом подобно хеш-секционированным таблицам, которые мы исследовали ранее, Oracle берет значение OWNER, хеширует его в секцию между 1 и 16 и помещает туда запись индекса. Если мы снова просмотрим информацию TKPROF для этих трех запросов, то обнаружим, что теперь картина намного ближе к той, которую мы наблюдали ранее с несекционированной таблицей — т.е. мы не оказываем отрицательного влияния на работу запросов:

| call | count | cpu | elapsed | disk | query | current | rows |
|-------|-------|------|---------|------|-------|---------|------|
| total | 4 | 0.00 | 0.00 | 0 | 4 | 0 | 1 |
| total | 5 | 0.00 | 0.00 | 0 | 11 | 0 | 20 |
| total | 5 | 0.00 | 0.00 | 0 | 628 | 0 | 26 |

Тем не менее, следует отметить, что хеш-секционированный индекс не может подвергаться сканированию диапазонов; в общем случае он больше всего подходит для условий с точным равенством (явное сравнение или вхождение в список конструкции IN). Если вы укажете в запросе конструкцию WHERE OWNER > :X, применяя предыдущий индекс, то выполнить простое сканирование диапазона с исключением секций не получится, а придется просматривать все 16 хеш-секций.

Использование конструкции ORDER BY

Этот пример невольно напомнил о несвязанном, но очень важном факте. При просмотре хеш-секционированных индексов мы столкнулись с еще одним случаем, когда применение индекса для получения данных *не* приводит к их автоматическому извлечению в отсортированном виде. Многие предполагают, что если в плане выполнения запроса отражено использование индекса при извлечении данных, то данные будут получены отсортированными. *Это никогда не было правдой.* Единственный способ извлечь данные отсортированными в любом порядке предусматривает применение в запросе конструкции ORDER BY. Если запрос не содержит конструкцию ORDER BY, то нельзя выдвигать какие-либо предположения, касающиеся порядка поступления данных.

Продemonстрируем сказанное на небольшом примере. Мы создадим маленькую таблицу как копию ALL_USERS с хеш-секционированным индексом, имеющим четыре секции по столбцу USER_ID:

```
EODA@ORA12CR1> create table t
2 as
3 select *
4   from all_users
5 /
```

Table created.

Таблица создана.

```
EODA@ORA12CR1> create index t_idx
2 on t(user_id)
3 global
4 partition by hash(user_id)
5 partitions 4
6 /
```

Index created.

Индекс создан.

Теперь запустим запрос в отношении этой таблицы и воспользуемся подсказкой, чтобы вынудить Oracle задействовать индекс. Обратите внимание на упорядочение (точнее на его отсутствие) данных:

```
EODA@ORA12CR1> set autotrace on explain
EODA@ORA12CR1> select /*+ index( t t_idx ) */ user_id
2   from t
3  where user_id > 0
4  /
```

```
USER_ID
-----
      13
...
      97
      22
...
     104
       8
```

```
...
      93
      7
...
      96
43 rows selected.
43 строки выбрано.
```

| Id | Operation | Name | Rows | Bytes | Cost (%CPU) | Time | Pstart | Pstop |
|-----|--------------------|-------|------|-------|-------------|----------|--------|-------|
| 0 | SELECT STATEMENT | | 43 | 172 | 4 (0) | 00:00:01 | | |
| 1 | PARTITION HASH ALL | | 43 | 172 | 4 (0) | 00:00:01 | 1 | 4 |
| * 2 | INDEX RANGE SCAN | T_IDX | 43 | 172 | 4 (0) | 00:00:01 | 1 | 4 |

```
EODA@ORA12CR1> set autotrace off
```

Таким образом, хотя Oracle применяет индекс при сканировании диапазона, данные очевидным образом не отсортированы. Фактически вы можете заметить в них некоторый шаблон. Здесь присутствуют четыре группы отсортированных результатов: многоточием заменены увеличивающиеся значения; между строками с USER_ID, равным 34 и 81, в выводе значения растут. Затем появляется строка с USER_ID = 19. Это объясняется тем, что Oracle возвращает "отсортированные данные" по очереди из каждой из четырех хеш-секций.

Просто примите к сведению: если запрос не имеет конструкции ORDER BY, то нет причин ожидать, что возвращенные данные будут каким-либо образом упорядочены. (И еще: конструкция GROUP BY также не обязана сортировать данные! Замены конструкции ORDER BY не существует.)

Значит ли это, что секционирование вообще не оказывает никакого влияния на производительность системы OLTP в позитивном смысле? Не совсем — просто мы должны взглянуть на вопрос по-другому. В целом секционирование не будет положительно влиять на производительность извлечения данных в системе OLTP; взамен необходимо позаботиться о том, чтобы избежать его отрицательного влияния на извлечение данных. Но в отношении модификации данных секционирование может обеспечить заметные преимущества в средах с высокой степенью параллелизма.

Вернемся к довольно простому предыдущему примеру с одной таблицей и единственным индексом и добавим к таблице первичный ключ. Без секционирования имеется одна таблица: в нее поступают все производимые вставки. Возможно, в таблице возникает конкуренция за списки свободных блоков. Кроме того, индекс по первичному ключу на столбце OBJECT_ID может оказаться правосторонним, как обсуждалось в главе 11. Предположительно он должен заполняться последовательностью; таким образом, все вставки будут осуществляться после крайнего правого блока, приводя к событиям ожидания занятых буферов. К тому же имеется единственная индексная структура T_IDX, за которую клиенты могут соперничать. Итак, пока что есть много одиночных элементов.

Перейдем к секционированию. Вы проводите хеш-секционирование таблицы по столбцу OBJECT_ID, разбивая ее на 16 секций. Теперь для состязаний имеется 16 таблиц, и к каждой таблице одновременно обращается одна шестнадцатая часть всех пользователей. Вы локально секционируете индекс по первичному ключу с использованием столбца OBJECT_ID, получая 16 секций. Вы располагаете 16 правосторонними индексными структурами, каждая из которых принимает одну шестнадцатую часть первоначальной нагрузки. И так далее. То есть вы можете применять секци-

онирование в среде с высокой степенью параллелизма для снижения конкуренции подобно тому, как в главе 11 использовался индекс по реверсированному ключу для сокращения событий ожидания занятых буферов. Однако вы должны помнить, что обработка секционированных данных требует большего количества ресурсов процессора, чем обработка несекционированных данных. Более высокая нагрузка на процессор связана с поиском места для размещения данных.

Следовательно, как и в случае любого другого средства, прежде чем применять в системе секционирование для увеличения производительности, убедитесь, что хорошо понимаете то, в чем *нуждается* система. Если система в текущий момент ограничена по ресурсам процессора, но его потребление не вызвано конкуренцией и ожиданиями освобождения зашлюк, то введение секций может только усугубить проблему, а не устранить ее!

Удобство средств обслуживания

В начале главы я указал, что преследуемой целью было предоставление практического руководства по реализации приложений с секционированием, не заостряя при этом особого внимания на администрировании. Тем не менее, в версии Oracle 12c появилось несколько новых средств администрирования, которые заслуживают специального обсуждения, в частности:

- множественные операции обслуживания секций;
- каскадный обмен;
- каскадное удаление.

Эти средства оказывают положительное влияние на удобство обслуживания, целостность данных и производительность. Следовательно, при внедрении секционирования важно иметь о них представление.

Множественные операции обслуживания секций

Это средство облегчает администрирование секций и в некоторых сценариях сокращает количество ресурсов базы данных, требуемых для выполнения операций обслуживания. До выхода Oracle 12c при проведении операций обслуживания секций, таких как добавление секции, было разрешено работать только с одной секцией за раз. Например, взгляните на следующую таблицу, секционированную по диапазонам:

```
EODA@ORA12CR1> create table p_table
  2 (a int)
  3 partition by range (a)
  4 (partition p1 values less than (1000),
  5  partition p2 values less than (2000));
Table created.
Таблица создана.
```

В версиях, предшествующих Oracle 12c, добавление двух секций к таблице делалось с помощью двух отдельных операторов SQL:

```
EODA@ORA12CR1> alter table p_table add partition p3 values less than (3000);
Table altered.
Таблица изменена.
```



```

EODA@ORA12CR1> alter table p_table add partition p4 values less than (4000);
Table altered.

```

Таблица изменена.

Начиная с Oracle 12c, в одном операторе можно выполнять множество операций над секциями. Показанный выше код можно переписать следующим образом:

```

EODA@ORA12CR1> alter table p_table add
  2 partition p3 values less than (3000),
  3 partition p4 values less than (4000);
Table altered.

```

Таблица изменена.

На заметку! В дополнение к добавлению секций множественные операции обслуживания секций могут применяться к удалению, объединению, разделению и усечению.

Выполнение множественных операций обслуживания секций в одном операторе DDL особенно полезно при разделении секций, поэтому оно заслуживает более подробного обсуждения. Подумайте о том, что произойдет в среде Oracle 11g, когда необходимо разделить годовую секцию P2014 на четыре квартальные секции: Q1, Q2, Q3 и Q4. Вам пришлось бы разделять секцию P2014 с помощью трех отдельных операторов DDL; каждая операция требует сканирования всех строк в разделяемой секции, после чего Oracle определяет, в какую секцию должна отправиться каждая строка, и затем осуществляет собственно вставку строки. Многократное разделение приводит к потреблению намного большего числа ресурсов, чем в случае, если разделение на множество секций можно было бы выполнить как одну операцию. Прдемонстрируем это на небольшом примере. Начнем с создания таблицы и загрузки в нее данных:

```

EODA@ORA12CR1> CREATE TABLE sales(
  2 sales_id int,
  3 s_date date)
  4 PARTITION BY RANGE (s_date)
  5 (PARTITION P2014 VALUES LESS THAN (to_date('01-jan-2015','dd-mon-yyyy')));
Table created.

```

Таблица создана.

```

EODA@ORA12CR1> insert into sales
  2 select level, to_date('01-jan-2014','dd-mon-yyyy') +
  3 ceil(dbms_random.value(1,364))
  4 from dual connect by level < 100000;
99999 rows created.
99999 строк создано.

```

Далее создадим служебную функцию, которая поможет измерить ресурсы, потребляемые во время выполнения операции:

```

EODA@ORA12CR1> create or replace function get_stat_val( p_name in
varchar2 ) return number
  2 as
  3 l_val number;
  4 begin
  5 select b.value
  6 into l_val

```

```

7      from v$statname a, v$mystat b
8      where a.statistic# = b.statistic#
9      and a.name = p_name;
10     return l_val;
11 end;
12 /

```

Function created.

Функция создана.

Теперь реализуем метод разделения секции на несколько подсекций, принятый в версиях до Oracle 12c, и измерим объем информации повторения, сгенерированной сеансом. Используя функцию GET_STAT_VAL, мы получаем текущее значение для статистического показателя redo size (размер данных повторения):

```

EODA@ORA12CR1> var r1 number
EODA@ORA12CR1> exec :r1 := get_stat_val('redo size');
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.

```

С помощью пакета DBMS_UTILITY мы записываем текущее время процессора:

```

EODA@ORA12CR1> var c1 number
EODA@ORA12CR1> exec :c1 := dbms_utility.get_cpu_time;
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.

```

После этого мы разделим секцию P2014 на четыре подсекции с применением трех отдельных операторов DDL согласно синтаксису для версий, предшествующих Oracle 12c:

```

EODA@ORA12CR1> alter table sales split partition
2  P2014 at (to_date('01-apr-2014','dd-mon-yyyy'))
3  into (partition Q1, partition Q2);
Table altered.
Таблица изменена.

```

```

EODA@ORA12CR1> alter table sales split partition
2  Q2 at (to_date('01-jul-2014','dd-mon-yyyy'))
3  into (partition Q2, partition Q3);
Table altered.
Таблица изменена.

```

```

EODA@ORA12CR1> alter table sales split partition
2  Q3 at (to_date('01-oct-2014','dd-mon-yyyy'))
3  into (partition Q3, partition Q4);
Table altered.
Таблица изменена.

```

Отообразим разницу в размерах информации повторения и времени процессора:

```

EODA@ORA12CR1> exec dbms_output.put_line(get_stat_val('redo size') - :r1);
4747712
EODA@ORA12CR1> exec dbms_output.put_line(dbms_utility.get_cpu_time - :c1);
16

```

Значительный объем данных повторения был сгенерирован из-за нескольких операций разделения, которые в результате привели к выполнению многочисленных

операторов вставки, т.к. Oracle нужно было разделять секцию много раз и повторно вставлять строки.

А теперь запустим точно такой же тест, но воспользуемся новым синтаксисом версии Oracle 12c, чтобы разделить секцию P2014 на четыре подсекции в одном операторе DDL (код для повторного создания и заполнения таблицы здесь не показан):

```

EODA@ORA12CR1> var r1 number
EODA@ORA12CR1> exec :r1 := get_stat_val('redo size');
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.

EODA@ORA12CR1> var c1 number
EODA@ORA12CR1> exec :c1 := dbms_utility.get_cpu_time;
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.

EODA@ORA12CR1> alter table sales split partition P2014
  2 into (partition Q1 values less than (to_date('01-apr-2014','dd-mon-yyyy')),
  3      partition Q2 values less than (to_date('01-jul-2014','dd-mon-yyyy')),
  4      partition Q3 values less than (to_date('01-oct-2014','dd-mon-yyyy')),
  5      partition Q4);
Table altered.
Таблица изменена.

EODA@ORA12CR1> exec dbms_output.put_line(get_stat_val('redo size') - :r1);
2099288

EODA@ORA12CR1> exec dbms_output.put_line(dbms_utility.get_cpu_time - :c1);
6

```

Объем информации повторения, сгенерированной единственным оператором DDL, почти вдвое меньше объема такой информации в предыдущем случае, а потребление времени процессора более чем в два раза ниже. В зависимости от количества разделяемых секций и включения одновременного обновления индексов объем данных повторения и время процессора могут оказаться значительно меньше, чем при разнесении операций обслуживания по нескольким операторам.

Каскадное усечение

Начиная с версии Oracle 12c, вы можете усекать родительские и дочерние таблицы вместе в одном атомарном операторе DDL. Хотя происходит каскадное усечение, любые запросы, выдаваемые в отношении комбинации родительской и дочерней таблиц, будут всегда иметь дело с согласованным по чтению представлением данных. Это значит, что данные в родительской и дочерней таблицах будут видны либо в состоянии, когда обе таблицы заполнены, либо в состоянии, когда обе таблицы усечены.

Функциональность каскадного усечения инициируется посредством оператора TRUNCATE...CASCADE для родительской таблицы. Чтобы каскадное усечение состоялось, любые дочерние таблицы должны быть определены с ограничением внешнего ключа ON DELETE CASCADE. Что каскадному усечению придется делать с секционированием? В секционированной по ссылкам таблице вы можете выполнить усечение секции родительской таблицы и распространить его каскадным образом на секцию дочерней таблицы в одной транзакции.

Давайте рассмотрим пример. Для применения функциональности TRUNCATE ... CASCADE к таблицам, секционированным по ссылкам, мы создадим сначала таблицу ORDERS:

```
EODA@ORA12CR1> create table orders
2 (
3   order#          number primary key,
4   order_date      date,
5   data            varchar2(30)
6 )
7 PARTITION BY RANGE (order_date)
8 (
9   PARTITION part_2014 VALUES LESS THAN (to_date('01-01-2015','dd-mm-yyyy')),
10  PARTITION part_2015 VALUES LESS THAN (to_date('01-01-2016','dd-mm-yyyy'))
11 )
12 /
```

Table created.

Таблица создана.

```
EODA@ORA12CR1> insert into orders values
2 ( 1, to_date( '01-jun-2014', 'dd-mon-yyyy' ), 'xyz' );
1 row created.
1 строка создана.
```

```
EODA@ORA12CR1> insert into orders values
2 ( 2, to_date( '01-jun-2015', 'dd-mon-yyyy' ), 'xyz' );
1 row created.
1 строка создана.
```

Следующей мы создадим таблицу ORDER_LINE_ITEMS, в которой используется конструкция ON DELETE CASCADE в ограничении внешнего ключа:

```
EODA@ORA12CR1> create table order_line_items
2 (
3   order#    number,
4   line#     number,
5   data      varchar2(30),
6   constraint cl_pk primary key(order#,line#),
7   constraint cl_fk_p foreign key(order#) references orders on delete cascade
8 ) partition by reference(cl_fk_p)
9 /
```

```
EODA@ORA12CR1> insert into order_line_items values ( 1, 1, 'yyy' );
1 row created.
1 строка создана.
```

```
EODA@ORA12CR1> insert into order_line_items values ( 2, 1, 'yyy' );
1 row created.
1 строка создана.
```

А теперь мы можем выдать оператор с конструкцией TRUNCATE...CASCADE, который выполнит усечение секции родительской таблицы и секции дочерней таблицы в виде одной транзакции:

```
EODA@ORA12CR1> alter table orders truncate partition PART_2014 cascade;
Table truncated.
Таблица усечена.
```

Другими словами, функциональность TRUNCATE ... CASCADE не позволяет приложениям видеть усеченную дочернюю таблицу до того, как будет усечена родительская таблица. Можно также выполнить усечение всех секций в родительской и дочерней таблицах:

```
EODA@ORA12CR1> truncate table orders cascade;
Table truncated.
Таблица усечена.
```

Имеет смысл повторить ради ясности: возможность каскадного усечения родительской и дочерней таблиц не является исключительно функцией секционирования. Это средство также применимо к несекционированным родительским и дочерним таблицам. Оно позволяет использовать один оператор DDL для запуска операций усечения и также гарантирует, что приложение базы данных всегда имеет дело с согласованным представлением родительских и дочерних секций.

Каскадный обмен

До выхода версии Oracle 12c при обмене секциями для таблицы, секционированной по ссылкам, применялась примерно такая последовательность действий.

1. Создание родительской таблицы и загрузка в нее данных.
2. Создание родительской секции в секционированной по ссылкам таблице.
3. Обмен секциями в родительской таблице с указанием конструкции UPDATE GLOBAL INDEXES.
4. Создание дочерней таблицы с ограничением внешнего ключа, которое указывает на секционированную по ссылкам родительскую таблицу.
5. Загрузка данных в дочернюю таблицу.
6. Обмен секциями в дочерней таблице.

В перечисленных выше действиях можно заметить, что у пользователей, обращающихся к базе данных, есть потенциальная возможность видеть строки в родительской таблице без соответствующих им строк в дочерней таблице. В версиях, предшествующих Oracle 12c, не существовало способа обойти такое поведение.

В Oracle 12c появилась возможность производить обмен комбинации родительских и дочерних таблиц, секционированных по ссылкам, в одном атомарном операторе DDL. Продемонстрируем это на небольшом примере. Создадим родительскую и дочернюю таблицы, секционированные по ссылкам:

```
EODA@ORA12CR1> create table orders
2  ( order#          number primary key,
3    order_date      date,
4    data            varchar2(30))
5  PARTITION BY RANGE (order_date)
6  (PARTITION part_2014 VALUES LESS THAN (to_date('01-01-2015','dd-mm-yyyy')) ,
7   PARTITION part_2015 VALUES LESS THAN (to_date('01-01-2016','dd-mm-yyyy')));
Table created.

EODA@ORA12CR1> insert into orders values (1, to_date( '01-jun-2014',
❖ 'dd-mon-yyyy' ), 'xyz');
1 row created.
1 строка создана.
```

```

EODA@ORA12CR1> insert into orders values (2, to_date( '01-jun-2015',
    ❸ 'dd-mon-yyyy' ), 'xyz');
1 row created.

```

1 строка создана.

```

EODA@ORA12CR1> create table order_line_items
2      (order#      number,
3        line#      number,
4        data       varchar2(30),
5        constraint c1_pk primary key(order#,line#),
6        constraint c1_fk_p foreign key(order#) references orders
7      ) partition by reference(c1_fk_p);

```

Table created.

Таблица создана.

```

EODA@ORA12CR1> insert into order_line_items values ( 1, 1, 'yyy' );
1 row created.

```

1 строка создана.

```

EODA@ORA12CR1> insert into order_line_items values ( 2, 1, 'yyy' );
1 row created.

```

1 строка создана.

Добавим пустую секцию в секционированную по ссылкам таблицу:

```

EODA@ORA12CR1> alter table orders add partition part_2016
2      values less than (to_date('01-01-2017','dd-mm-yyyy'));

```

Table altered.

Таблица изменена.

Родительская и дочерняя таблицы созданы, и данные в них загружены. Создадим таблицы, которые подлежат обмену с пустыми секциями в таблице, секционированной по ссылкам:

```

EODA@ORA12CR1> create table part_2016
2      ( order#      number primary key,
3        order_date  date,
4        data        varchar2(30));

```

Table created.

Таблица создана.

```

EODA@ORA12CR1> insert into part_2016 values (3, to_date('01-jun-2016',
    ❸ 'dd-mon-yyyy' ), 'xyz');
1 row created.

```

1 строка создана.

```

EODA@ORA12CR1> create table c_2016
2      (order#      number,
3        line#      number,
4        data       varchar2(30),
5        constraint cel_pk primary key(order#,line#),
6        constraint cel_fk_p foreign key(order#) references part_2016);

```

Table created.

Таблица создана.

```

EODA@ORA12CR1> insert into c_2016 values(3, 1, 'xyz');
1 row created.
1 строка создана.

```

Теперь мы можем произвести собственно обмен в одной транзакции. Обратите внимание на указание конструкции CASCADE:

```

EODA@ORA12CR1> alter table orders
2  exchange partition part_2016
3  with table part_2016
4  without validation
5  CASCADE
6  update global indexes;
Table altered.
Таблица изменена.

```

Вот и все. С помощью одного оператора DDL мы одновременно обменяли две таблицы, связанные ограничением внешнего ключа, и таблицы, секционированной по ссылкам. Любой получающий доступ к базе данных будет видеть секции родительской и дочерней таблиц, добавленные в виде одной единицы работы.

Аудит и сжатие пространства сегментов

Несколько лет тому назад еще не действовали ограничения правительства США, подобные HIPAA (Health Insurance Portability and Accountability Act — Акт о передаче и учете информации о страховании здоровья); <http://www.hhs.gov/ocr/hipaa>. Компании вроде Enron еще присутствовали на рынке, а другие государственные требования по соответствию системы бухгалтерского учета, изложенные позднее в законе Сарбейнса-Оксли (Sarbanes-Oxley), пока не были введены. Тогда аудит считался чем-то таким, что может быть сделано когда-то в будущем, да и то — не факт. Однако в наши дни аудит выходит на передний план, и многие администраторы баз данных обязаны в течение семи лет удерживать в оперативном режиме информацию аудита для своих финансовых, деловых и медицинских баз данных.

Информация аудита — это порция данных, которая добавляется, но никогда не извлекается во время нормальной работы. Она преимущественно служит судебным доказательством после свершившегося факта. Мы должны иметь ее, но с многих точек зрения это нечто такое, что находится на дисках и занимает место — очень и очень много места. Таким образом, ежемесячно, ежегодно или с какой-то другой периодичностью мы должны очищать ее либо архивировать. Аудит представляет собой такую вещь, которая, не будучи правильно спроектированной с самого начала, может попросту прикончить вас в конце. Через семь лет, отделяющих текущий момент от первой очистки или архивирования старых данных, будет поздно задумываться о том, как вы будете это делать. Если вы не заложили такую возможность в проект, то выведение из эксплуатации старых данных может стать очень болезненным.

Обратитесь к двум технологиям, которые не только делают задачу аудита приемлемой по сложности, но также значительно облегчают управление данными аудита и сокращают потребление дискового пространства. Этими технологиями являются секционирование и сжатие пространства сегментов, как обсуждалось в главе 10. Второй вариант может быть не настолько очевидным, потому что *базовое* сжатие пространства сегментов работает только с крупными пакетными операциями, такими как загрузка в прямом режиме (сжатие OLTP входит в состав опции расширенного сжатия (Advanced Compression Option) — оно не доступно во всех редакциях базы данных), а информация аудита обычно вставляется в строку при наступлении определенных событий. Трюк заключается в комбинировании секций скользящего окна со сжатием пространства сегментов.

Предположим, что мы решили секционировать информацию аудита по месяцам. На протяжении первого месяца мы просто вставляем данные в секционированную таблицу; эти вставки производятся в обычном режиме, а не в прямом, а потому данные не сжимаются. Перед самым концом месяца мы добавим к таблице новую секцию, куда будет помещаться информация аудита будущего месяца. Вскоре после начала следующего месяца мы выполним крупную пакетную операцию над информацией аудита предыдущего месяца, а именно — воспользуемся командой `ALTER TABLE` для перемещения секции прошлого месяца, что также даст эффект сжатия данных. Если продвинуться на шаг дальше, мы могли бы переместить эту секцию из табличного пространства с доступом по чтению и записи, в котором она должна была находиться, в табличное пространство, поддерживающее только чтение (и содержащее другие секции для информации аудита). Таким образом, мы можем сохранять резервную копию этого табличного пространства раз в месяц после перемещения в него новой секции, гарантируя наличие актуальной, чистой и читабельной копии, которую мы не будем трогать в течение данного месяца. Для хранения информации аудита можно предусмотреть следующие табличные пространства.

- Текущее оперативное табличное пространство, доступное для чтения и записи, которое подвергается резервному копированию подобно любому другому табличному пространству в системе. Информация аудита в нем не сжата и постоянно пополняется.
- Табличное пространство, доступное только для чтения, которое содержит секции аудита текущего года в сжатом формате. В начале каждого месяца мы делаем это табличное пространство доступным для чтения и записи, перемещаем в него со сжатием информацию аудита последнего месяца, делаем его снова доступным только для чтения и производим резервное копирование.
- Последовательность табличных пространств за прошлый год, позапрошлый год и т.д. Все они доступны только для чтения и вполне могут размещаться на медленных, дешевых носителях. В случае отказа носителя понадобится просто выполнить восстановление из резервной копии. Периодически необходимо проводить выборочную проверку резервных копий, чтобы удостовериться в их читабельности (ленты иногда разрушаются).

В таком стиле очистка выполняется легко (т.е. удалением секции). Также легко производится архивирование — достаточно переместить табличное пространство и восстановить его позднее. Мы сокращаем потребление пространства за счет внедрения сжатия. Мы уменьшаем размеры резервных копий, т.к. во многих системах единственным самым крупным набором данных является *информация аудита*. Если вам удастся исключить часть или все это из процесса ежедневного резервного копирования, то разница окажется довольно заметной.

Короче говоря, требования к информации аудита и секционирование следуют рука об руку независимо от лежащего в основе типа системы, будь то хранилище данных или система OLTP.

Совет. Подумайте о применении предлагаемого Oracle средства архива ретроспективных данных (Flashback Data Archive) для удовлетворения требований аудита. Когда средство Flashback Data Archive включено для таблицы, оно будет автоматически создавать внутреннюю секционированную таблицу для записи транзакционной информации.

Резюме

Секционирование исключительно удобно для масштабирования крупных объектов в базе данных. Это масштабирование наглядно проявляется в масштабировании производительности, доступности и администрирования. Все три аспекта чрезвычайно важны для разных людей. Администраторы баз данных заинтересованы в масштабировании администрирования. Владельцев системы интересует ее доступность, поскольку любой простой означает потерю денег, и все, что сокращает простой — или снижает их влияние — повышает окупаемость системы. Конечные пользователи заинтересованы в масштабировании производительности. В конце концов, никому не нравится иметь дело с медленной системой.

Мы также обратили внимание на тот факт, что в системе OLTP секции могут и не увеличивать производительность, особенно в случае их неподходящего использования. Секционирование может увеличить производительность определенного класса запросов, но такие запросы обычно в системе OLTP не применяются. Данный момент важно понять, т.к. многие люди ассоциируют секционирование с “бесплатным увеличением производительности”. Это вовсе не означает, что секционирование *не* должно применяться в системах OLTP — оно обеспечивает много других заметных преимуществ в таких средах — просто не ожидайте массового роста пропускной способности. Ожидайте сокращения времени простоя. Ожидайте той же самой приемлемой производительности (при правильном использовании секционирование *не* замедляет работу системы). Ожидайте облегчения управляемости, что может привести к повышению производительности за счет того факта, что некоторые операции по обслуживанию будут выполняться администраторами баз данных более часто, поскольку появится такая возможность.

Мы исследовали разнообразные схемы секционирования таблиц, предлагаемые Oracle — секционирование по диапазонам, хеш-секционирование, секционирование по списку, секционирование по интервалам, секционирование по ссылкам, секционирование по интервалам и по ссылкам, секционирование по виртуальному столбцу и составное секционирование — и обсудили, где их целесообразно применять. Мы уделили много времени рассмотрению секционированных индексов и определению отличий между префиксными и непрефиксными, локальными и глобальными индексами. Мы описали выполнение операций над секциями в сочетании с глобальными индексами внутри хранилищ данных, а также достижение компромисса между потреблением ресурсов и доступностью. Мы также взглянули на новые удобные средства обслуживания в Oracle 12c, такие как возможность выполнения операций обслуживания для множества секций одновременно, каскадное усечение и каскадный обмен. С каждым новым выпуском базы данных секционирование продолжает обновляться и совершенствоваться.

С течением времени я наблюдаю рост значимости секционирования для более широкой аудитории по мере увеличения размера и масштаба приложений баз данных. Сеть Интернет с присущими ей потребностями в базах данных, наряду с требованием законодательства о долгосрочном хранении информации аудита, ведет к появлению все более и более крупных наборов данных, и секционирование — это естественный инструмент, помогающий справиться с этой задачей.

Параллельное выполнение

Параллельное выполнение, средство редакции Enterprise (в редакции Standard оно не доступно), впервые появилось в версии Oracle 7.1.6 в 1994 году. Оно позволяет физически разбивать крупную последовательную задачу (любую операцию DML или DDL в общем случае) на множество фрагментов меньшего размера, которые могут быть обработаны одновременно. Параллельное выполнение в Oracle имитирует реальные процессы, которые мы наблюдаем постоянно. Например, вряд ли можно ожидать, что строительством дома будет заниматься всего один человек; вместо этого люди объединяются в бригаду, чтобы работать параллельно с целью быстрого возведения дома. Таким образом, определенные операции могут быть разделены на более мелкие задачи, которые будут выполняться одновременно; например, водопроводно-канализационные работы и монтаж электропроводки могут проводиться параллельно, чтобы сократить общее время, необходимое для завершения всего строительства.

Параллельное выполнение в Oracle следует во многом похожей логике. Нередко Oracle в состоянии разделить определенную большую работу на меньшие части и выполнить каждую часть параллельно с другими. Другими словами, если требуется полное сканирование крупной таблицы, то нет причин, по которым база данных Oracle не смогла бы запустить четыре параллельных сеанса, P001—P004, вместе выполняющие полное сканирование, при этом каждый сеанс читал бы отдельную часть таблицы. Если данные, просканированные сеансами P001—P004, нуждаются в сортировке, то этого можно добиться с помощью еще четырех параллельных сеансов, P005—P008, которые в конечном итоге отправят результаты в общий координирующий сеанс для запроса.

Параллельное выполнение — это инструмент, который при подходящем использовании может на порядки улучшить показатели времени отклика некоторых операций. Однако при непродуманном применении результаты обычно будут прямо противоположными. Цель настоящей главы заключается вовсе не в детальном объяснении реализации параллельных запросов в Oracle, бесчисленных комбинаций планов, которые могут быть результатом параллельных операций, и других аспектов подобного рода. Эти сведения можно найти в руководстве для администратора баз данных Oracle (*Oracle Database Administrator's Guide*), в руководстве по концепциям базы данных Oracle (*Oracle Database Concepts*), в руководстве по очень крупным ба-

зам данных и секционированию Oracle (*Oracle Database VLDB and Partitioning Guide*) и особенно в руководстве по хранилищам данных (*Oracle Database Data Warehousing Guide*). Цель главы — дать общее представление о том, для решения каких классов задач параллельное выполнение подходит, а для каких оно неприемлемо. В частности, после того, как мы посмотрим, когда использовать параллельное выполнение, мы раскроем следующие темы.

- **Параллельный запрос.** Это способность Oracle выполнять одиночный запрос с применением множества процессов или потоков операционной системы. База данных Oracle будет находить операции, которые она может проводить параллельно, такие как полные сканирования таблиц или масштабные сортировки, и создавать планы запросов, выполняющие их параллельно.
- **Параллельный DML (Parallel DML — PDML).** По своей природе средство PDML очень похоже на параллельный запрос, но относится к выполнению модификаций (INSERT, UPDATE, DELETE и MERGE) с использованием параллельной обработки. В этой главе мы рассмотрим PDML и обсудим некоторые присущие ему ограничения.
- **Параллельный DDL.** Параллельный DDL — это возможность Oracle выполнять крупные операции DDL параллельно. Например, операции перестройки индекса, создания нового индекса, загрузки данных через CREATE TABLE AS SELECT и реорганизации больших таблиц могут применять параллельную обработку. Я уверен, что это наилучшие точки для внедрения параллелизма в базе данных, поэтому основное обсуждение будет сосредоточено на данной теме.
- **Параллельная загрузка.** Внешние таблицы и средство SQL*Loader обладают возможностью загрузки данных параллельно. Эта тема кратко рассматривается здесь и в главе 15.
- **Процедурный параллелизм.** Это возможность параллельного запуска разработанного вами кода. В настоящей главе мы обсудим два имеющихся подхода. Первый подход предусматривает параллельное выполнение кода PL/SQL в манере, прозрачной для разработчиков (сами разработчики не пишут параллельный код, взамен Oracle прозрачным образом распараллеливает готовый код). При втором подходе, который я обычно называю “самодельным параллелизмом”, разработчики проектируют код специально для его параллельного выполнения. Мы рассмотрим два метода употребления такого самодельного параллелизма — в высокой мере ручную реализацию, действующую в Oracle 11g Release 1 и предшествующих версиях, и новый автоматизированный метод, доступный в Oracle 11g Release 2 и последующих версиях.

Существуют еще два типа параллельного выполнения, которые заслуживают упоминания, но они выходят за рамки материала этой книги. Они приводятся здесь ради полноты, а также для того, чтобы вы знали о них, когда будете иметь дело с операциями восстановления или репликации.

- **Параллельное восстановление.** Другой формой параллельного выполнения в Oracle является возможность параллельного восстановления. Параллельное восстановление может проводиться на уровне экземпляра и увеличивать скорость процедуры восстановления, которая должна быть выполнена после

аварийного отказа программного обеспечения, операционной системы или системы в целом (например, из-за неожиданного перебоя электропитания). Параллельное восстановление может быть также применено во время восстановления носителя (скажем, из резервной копии). Раскрытие здесь всех вопросов, связанных с восстановлением, в мои цели не входит, поэтому я напоминаю о существовании параллельное восстановление лишь мимоходом. За детальными сведениями по этой теме обращайтесь в руководство пользователя по резервному копированию и восстановлению Oracle (Oracle Backup and Recovery User's Guide).

- **Параллельное распространение.** Тип параллельного выполнения, используемый опцией Oracle Advanced Replication (Расширенная репликация Oracle), который обеспечивает асинхронное параллельное распространение транзакций. В этом режиме Oracle применяет параллельные процессы для увеличения полосы пропускания операций репликации. Дополнительные сведения о параллельном распространении ищите в руководстве по расширенной репликации баз данных Oracle (Oracle Database Advanced Replication).

После такого краткого введения в параллельное выполнение давайте начнем с исследования ситуаций, когда уместно использовать это средство.

Использование параллельного выполнения

Параллельное выполнение может быть фантастически полезным. Оно позволяет взять процесс, выполняющийся много часов или дней, и завершить его за считанные минуты. Разбиение крупной задачи на мелкие компоненты в некоторых случаях может значительно сократить время обработки. Тем не менее, есть одна основополагающая концепция, которую следует иметь в виду относительно параллельного выполнения. Она подытожена следующей цитатой из книги Джонатана Льюиса *Practical Oracle8i: Building Efficient Databases* (Addison-Wesley, 2001 год):

Параллельные запросы по существу не являются масштабируемыми.

Хотя этой цитате уже почти 15 лет, она справедлива и сегодня, если не сказать, что даже в большей степени, чем на то время. Параллельное выполнение — по своей сути немасштабируемое решение. Оно было разработано для того, чтобы позволить индивидуальному пользователю или отдельному оператору PL/SQL потреблять все ресурсы базы данных. Если есть средство, позволяющее кому-то задействовать все доступные ресурсы, которое затем предоставляется в распоряжение двум пользователям, то возникнут очевидные проблемы конкуренции. Когда количество параллельных пользователей в системе начинает превышать объем имеющихся ресурсов (памяти, центрального процессора и системы ввода-вывода), то возможность развертывания параллельных операций оказывается под вопросом. Например, на четырехпроцессорной машине при наличии в среднем 32 пользователей, одновременно выполняющих запросы, скорее всего, вы *не* захотите распараллеливать их операции. Позволив каждому пользователю выполнять только 2 запроса параллельно, вы получили бы 64 параллельных операции на машине, располагающей всего четырьмя процессорами. Если машина не была перегружена до применения параллельного выполнения, то теперь почти наверняка это произойдет.

Короче говоря, параллельное выполнение также может быть совершенно неподходящей идеей. Во многих случаях использование параллельной обработки будет приводить лишь к повышенному потреблению ресурсов, т.к. параллельное выполнение пытается захватить *все доступные ресурсы*. В системе, где ресурсы должны совместно применяться множеством параллельных транзакций, как это имеет место в системе OLTP, скорее всего, будет наблюдаться *увеличение* времени отклика. Система избегает использования определенных приемов, которые могут эффективно применяться в последовательном плане выполнения. Вместо этого она использует пути выполнения, такие как полное сканирование, в надежде на то, что запуск многочисленных частей крупной пакетной операции в параллельном режиме даст лучшие результаты, чем последовательный план. Параллельное выполнение, примененное неподходящим образом, может стать причиной возникновения проблем с производительностью, а не их решением.

Таким образом, перед использованием параллельного выполнения должны быть удовлетворены следующие два условия.

- Вы должны иметь дело с очень большой задачей, подобной полному сканированию 50 Гбайт данных.
- В вашем распоряжении должен быть достаточный объем доступных ресурсов. Перед параллельным полным сканированием 50 Гбайт данных необходимо удостовериться в наличии достаточного количества свободных процессоров для обслуживания параллельных процессов и достаточной пропускной способности ввода-вывода. Все 50 Гбайт должны быть разбросаны по более чем одному физическому диску, чтобы позволить множеству запросов чтения происходить одновременно, должно быть в наличии достаточное число каналов дискового ввода-вывода, чтобы извлекать данные из диска параллельно, и т.д.

В случае небольшой задачи, типичным примером которой служат запросы в системе OLTP, или недостаточном количестве *доступных* ресурсов, что, опять-таки, характерно для системы OLTP, где процессор и система ввода-вывода зачастую уже задействованы по максимуму, параллельное выполнение не подойдет. Итак, для лучшего понимания концепции ниже представлена аналогия.

Аналогия параллельной обработки

Для описания параллельной обработки и объяснения, почему необходимо иметь крупную задачу и достаточный объем свободных ресурсов в базе данных, я часто привожу аналогию следующего вида. Предположим, что есть две задачи, которые должны быть решены. Первая задача состоит в написании одностраничной сводки по новому товару. Вторая задача заключается в написании всеобъемлющего отчета из десяти разделов, каждый из которых относительно независим от других. Возьмем для примера настоящую книгу: эта глава, посвященная параллельному выполнению, довольно самостоятельна и не связана с главой о повторе и отмене — их не обязательно писать последовательно.

Как же подойти к решению каждой задачи? Какая из них, по вашему мнению, могла бы выиграть от параллельной обработки?

Одностраничная сводка

В этой аналогии подготовка одностраничной сводки является небольшой задачей. Ее можно выполнить самостоятельно или поручить одному исполнителю. Почему? Потому что объем работы по распараллеливанию этого процесса превысил бы объем работы по написанию документа самостоятельно. Вы должны оценить, что документ должен содержать 12 абзацев, выяснить, что каждый абзац не зависит от других, собрать команду на совещание, выбрать 12 исполнителей, объяснить им задачу и поручить каждому из них написание отдельного абзаца. Затем вам потребуется выступить координатором и собрать вместе написанные ими абзацы, расположить их в правильном порядке, проверить правильность и распечатать результат. Все это вероятно займет больше времени, чем когда вы *просто напишете весь документ сами, последовательно*. Накладные расходы по управлению большой группой людей в проекте такого масштаба намного перевесят любой выигрыш, который получится от параллельного написания 12 абзацев.

Точно такой же принцип применяется к параллельному выполнению в базе данных. В случае работы, требующей при последовательном выполнении нескольких секунд, введение параллельного выполнения вместе со связанными накладными расходами, скорее всего, приведет к тому, что задача потребует более длительного времени на завершение.

Отчет из десяти разделов

Теперь рассмотрим вторую задачу. Если вы хотите получить этот отчет из десяти разделов насколько возможно быстро, то самым медленным путем для достижения этой цели было бы поручение всей работы единственному исполнителю (поверьте, я знаю — взгляните на эту книгу; иногда мне хотелось, чтобы ее писали одновременно 15 моих клонов). Итак, вы должны организовать совещание, оценить процесс, выдать задания, взять на себя функцию координатора, собрать результаты, связать их в окончательный отчет и отправить его по назначению. Это позволит сделать работу если не за 1/10, то примерно за 1/8 часть времени, которое бы заняло написание отчета одним человеком. Я снова подчеркиваю: это возможно только при условии, что *вы располагаете достаточным количеством свободных ресурсов*. Если в вашем распоряжении имеется большой штат работников, которые в текущий момент ничем не заняты, то разделение работы на части совершенно оправдано.

Однако представьте, что персонал одновременно решает много задач и в текущий момент достаточно плотно загружен работой. В этом случае следует проявить осторожность с таким большим проектом. Вы должны быть уверены в том, что исполнители не окажутся перегруженными; вы вовсе не хотите, чтобы они работали до полного изнеможения. Вы не можете поручить выполнение большего объема работ, чем могут справиться ресурсы (персонал), иначе они попросту уволятся. Если персонал уже полностью занят, то добавление новых работ сдвинет календарные планы и задержит завершение всех проектов.

Параллельное выполнение в Oracle во многом похоже. При наличии задачи, решение которой занимает много минут, часов или дней, использование параллельного выполнения может в восемь раз увеличить скорость. Но если вы уже серьезно ограничены в ресурсах (сотрудники перегружены работой), то параллельного выполнения следует избегать, иначе система может еще более тормозить. И хотя сервер-

ные процессы Oracle не уволятся в знак протеста, они могут переполнять память и отказывать либо просто испытывать настолько длительные ожидания ввода-вывода или процессора, что будут выглядеть так, как будто бы никакой работы не происходит.

Если иметь это в виду и помнить, что никогда не стоит доводить аналогию до крайностей, то вы получаете разумное руководящее правило для определения целесообразности применения параллелизма. Работу, требующую нескольких секунд, параллельное выполнение вряд ли ускорит — более вероятно противоположное. Если ресурсов уже недостаточно (т.е. ресурсы полностью задействованы), то добавление параллельного выполнения, скорее всего, ухудшит картину, а не улучшит ее. Параллельное выполнение великолепно себя ведет в ситуации с действительно крупной задачей и при избытке ресурсов. В этой главе мы взглянем на некоторые способы эффективной эксплуатации параллельного выполнения.

Oracle Exadata

Машина баз данных Oracle Exadata Database Machine — это комбинация оборудования и программного обеспечения, предлагаемая Oracle Corporation, которая переносит параллельные операции на новый уровень. Продукт Oracle Exadata представляет собой решение с массовым параллелизмом для задач в крупных базах данных (с объемами в сотни и тысячи терабайтов). В нем сочетается оборудование — специализированная сеть хранения данных (storage area network — SAN) для базы данных и программное обеспечение — Oracle Enterprise Linux или Oracle Solaris и части программного обеспечения Oracle на уровне физических дисков. Эти части предназначены для разгрузки того, что обычно являлось типовыми функциями обработки сервера базы данных, и их выполнения на уровне самого хранилища:

- сканирование блоков;
- обработка блоков;
- шифрование/дешифрация блоков;
- фильтрация блоков (применение конструкции WHERE);
- выборка столбцов из блоков;
- обработка регулярных выражений.

Результатом является гигантское сокращение использования ресурсов на сервере базы данных, т.к. вместо полного набора данных он получает и обрабатывает только те строки и столбцы, которые необходимы (устройства хранения помимо прочего уже производят обработку конструкции WHERE и формирование списка столбцов оператора SELECT). При этом не только значительно (во много раз) ускоряется сканирование диска, но и увеличивается скорость обработки самих данных за счет массового параллелизма.

Таким образом, в основе Oracle Exadata лежит параллельная обработка, но в настоящей книге этот продукт подробно не рассматривается. В будущем я предвижу, что ему придется посвятить главу (или две). А пока сосредоточим внимание на собственных возможностях параллельной обработки Oracle без продукта Exadata.

Параллельный запрос

Параллельный запрос позволяет разделять одиночный SQL-оператор `SELECT` на множество запросов меньших размеров, каждый из которых выполняется одновременно с другими, после чего их результаты объединяются для получения окончательного ответа. Например, рассмотрим следующий запрос:

```
EODA@ORA12CR1> select count(status) from big_table;
```

В режиме параллельного выполнения этот запрос мог бы инициировать некоторое количество параллельных сеансов, разбить таблицу `BIG_TABLE` на небольшие неперекрывающиеся срезы, и затем поручить каждому сеансу чтение своей части таблицы и подсчет в ней строк. Координатор параллельного запроса получил бы все значения агрегированных счетчиков от индивидуальных параллельных сеансов и собрал их в одно целое, возвращая финальный ответ клиентскому приложению. Графическое представление всего процесса показано на рис. 14.1.

На заметку! В действительности между процессами и файлами нет однозначного соответствия, как изображено на рис. 14.1. На самом деле все данные таблицы `BIG_TABLE` могут находиться в одном файле, обрабатываемом четырьмя параллельными процессами. Данные могут также храниться в двух файлах, которые обрабатываются теми же четырьмя параллельными процессами, или в общем случае — в любом количестве файлов.

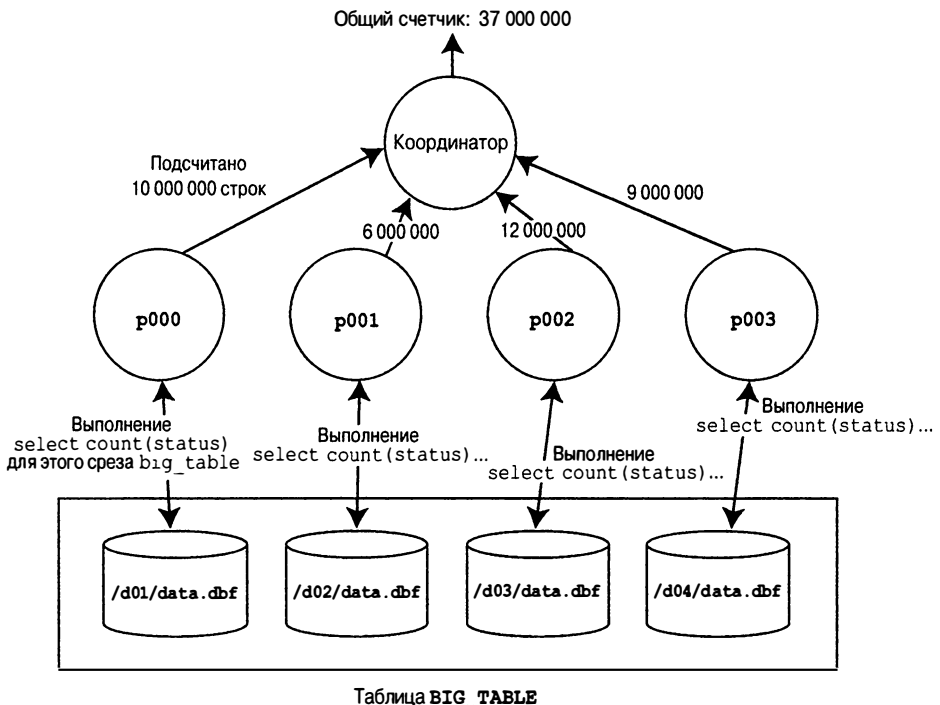


Рис. 14.1. Схема выполнения запроса `select count(status)`

Процессы p000, p001, p002 и p003 известны как *серверы параллельного выполнения*, иногда также называемые *подчиненными процессами параллельного запроса* (parallel query — PQ). Каждый из этих серверов параллельного выполнения является отдельным сеансом, подключенным так, как если бы он был процессом выделенного сервера. Каждый из них отвечает за сканирование неперекрывающейся области таблицы BIG_TABLE, формирование результирующего поднабора и отправку своего вывода координирующему серверу — серверному процессу исходного сеанса, — который будет объединять эти частичные результаты в окончательный ответ.

Это можно видеть в плане выполнения. Мы возьмем таблицу BIG_TABLE с 10 миллионами строк (сценарий создания BIG_TABLE описан в разделе “Настройка среды” в самом начале книги), включим параллельный запрос для этой таблицы и посмотрим на него в действии. Пример выполнялся на четырехпроцессорной машине в среде Oracle 12c Release 1 со стандартными значениями для всех параметров, касающихся параллелизма. То есть речь идет о стандартной установке, где были заданы только самые необходимые параметры, включая MEMORY_TARGET (установлен в 4 Гбайт), CONTROL_FILES, DB_BLOCK_SIZE (установлен в 8 Кбайт) и PGA_AGGREGATE_TARGET (установлен в 512 Мбайт). В начале можно ожидать следующего плана:

```
EODA@ORA12CR1> explain plan for select count(status) from big_table;
Explained.
Объяснено.
```

```
EODA@ORA12CR1> select * from table(dbms_xplan.display(null, null,
2      'TYPICAL -ROWS -BYTES -COST'));

```

| Id | Operation | Name | Time |
|----|-------------------|-----------|----------|
| 0 | SELECT STATEMENT | | 00:00:03 |
| 1 | SORT AGGREGATE | | |
| 2 | TABLE ACCESS FULL | BIG_TABLE | 00:00:03 |

На заметку! Разные выпуски Oracle имеют отличающиеся стандартные настройки для разнообразных средств, связанных с параллелизмом — иногда отличающиеся радикально. Не удивляйтесь, если при запуске некоторых примеров в старых выпусках будут получены совершенно другие результаты.

Это типичный *последовательный* план. Здесь нет никакого параллелизма, потому что возможность параллельного запроса не была включена, а по умолчанию она неактивна.

Включить возможность параллельного запроса можно многими способами, в том числе с применением подсказки прямо в запросе либо путем изменения таблицы, чтобы учитывались пути параллельного выполнения (именно этот вариант здесь используется).

Мы можем специально навязывать степень параллелизма для учета в путях выполнения относительно этой таблицы. Например, мы можем сообщить Oracle о том, что при создании планов выполнения для этой таблицы нас интересует степень параллелизма 4. Это делается с помощью такого кода:

```
EODA@ORA12CR1> alter table big_table parallel 4;
Table altered.
Таблица изменена.
```

Я предпочитаю просто указывать Oracle о том, что необходимо учитывать параллельное выполнение, но определять степень параллелизма на основе текущей рабочей нагрузки в системе и самого запроса. Это значит, что степени параллелизма позволено варьироваться во времени по мере роста и уменьшения нагрузки в системе. При наличии множества свободных ресурсов степень параллелизма будет возрастет, а в периоды ограниченного объема доступных ресурсов — снижаться. Вместо чрезмерной нагрузки машины при фиксированной степени параллелизма такой подход позволяет Oracle динамически увеличивать и уменьшать количество одновременно используемых ресурсов, требуемых запросом.

Мы просто включаем возможность параллельных запросов для данной таблицы посредством команды ALTER TABLE:

```
EODA@ORA12CR1> alter table big_table parallel;
Table altered.
Таблица изменена.
```

Теперь параллельные запросы будут учитываться для операций над данной таблицей. На этот раз план выполнения выглядит следующим образом:

```
EODA@ORA12CR1> explain plan for select count(status) from big_table;
Explained.
Объяснено.
EODA@ORA12CR1> select * from table(dbms_xplan.display(null, null,
2      'TYPICAL -ROWS -BYTES -COST'));

```

| Id | Operation | Name | Time | TQ | IN-OUT | PQ Distrib |
|----|---------------------|-----------|----------|-------|--------|------------|
| 0 | SELECT STATEMENT | | 00:00:01 | | | |
| 1 | SORT AGGREGATE | | | | | |
| 2 | PX COORDINATOR | | | | | |
| 3 | PX SEND QC (RANDOM) | :TQ10000 | | Q1,00 | P->S | QC (RAND) |
| 4 | SORT AGGREGATE | | | Q1,00 | PCWP | |
| 5 | PX BLOCK ITERATOR | | 00:00:01 | Q1,00 | PCWC | |
| 6 | TABLE ACCESS FULL | BIG_TABLE | 00:00:01 | Q1,00 | PCWP | |

Обратите внимание, что общее время для запроса, выполняемого параллельно, составило 00:00:01 в противоположность предыдущей оценке 00:00:03 для последовательного плана. Помните, что это *оценки*, а не обещания!

Если просмотреть этот план снизу вверх, начиная со строки с ID, равным 6, он отразит шаги, представленные на рис. 14.1. Полное сканирование таблицы будет разделено на множество сканирований меньших размеров (шаг 5). Каждое из них займется сбором своего значения COUNT(STATUS) (шаг 4). Эти промежуточные результаты будут переданы координатору параллельного запроса (шаги 2 и 3), который объединить их (шаг 1) и выведет ответ.

Стандартные серверы параллельного выполнения

При запуске экземпляра Oracle применяет значение параметра инициализации `PARALLEL_MIN_SERVERS` для определения количества серверов параллельного выполнения, которые должны запускаться автоматически. Эти процессы используются для обслуживания параллельно выполняемых операторов. В версии Oracle 11g стандартным значением `PARALLEL_MIN_SERVERS` было 0, т.е. по умолчанию во время старта экземпляра параллельные процессы не запускаются.

Начиная с версии Oracle 12c, минимальное значение `PARALLEL_MIN_SERVERS` вычисляется по формуле `CPU_COUNT × PARALLEL_THREADS_PER_CPU × 2`. В среде Linux/UNIX эти процессы можно просмотреть с применением команды `ps`:

```
$ ps -aef | grep '^oracle.*ora_p00._ORA12CR1'
oracle 18518 1 0 10:13 ? 00:00:00 ora_p000_ORA12CR1
oracle 18520 1 0 10:13 ? 00:00:00 ora_p001_ORA12CR1
oracle 18522 1 0 10:13 ? 00:00:00 ora_p002_ORA12CR1
oracle 18524 1 0 10:13 ? 00:00:00 ora_p003_ORA12CR1
oracle 18526 1 0 10:13 ? 00:00:00 ora_p004_ORA12CR1
...
```

В версиях, предшествующих Oracle 12c, если вы не изменили значение `PARALLEL_MIN_SERVERS` со стандартного 0, то при первом запуске экземпляра не увидите никаких серверных процессов параллельного выполнения. Такие процессы появятся после того, как вы запустите оператор, обрабатываемый параллельно.

Если вы достаточно любопытны, чтобы проследить параллельный запрос, можно воспользоваться для этого двумя сеансами. В сеансе, в котором будет запущен параллельный запрос, начнем с определения идентификатора `SID`:

```
EODA@ORA12CR1> select sid from v$mystat where rownum = 1;
SID
-----
258
```

В другом сеансе подготовим запрос к запуску (пока не запуская его, а только набрав):

```
EODA@ORA12CR1> select sid, qcsid, server#, degree
2 from v$px_session
3 where qcsid = 258
```

Теперь, вернувшись в первоначальный сеанс, идентификатор `SID` которого мы запрашивали, запустим параллельный запрос. В сеансе, где был подготовлен запрос, теперь можно запустить его и наблюдать примерно такой вывод:

```
4 /
```

| SID | QCSID | SERVER# | DEGREE |
|-----|-------|---------|--------|
| 26 | 258 | 1 | 8 |
| 102 | 258 | 2 | 8 |
| 177 | 258 | 3 | 8 |
| 267 | 258 | 4 | 8 |
| 23 | 258 | 5 | 8 |
| 94 | 258 | 6 | 8 |
| 169 | 258 | 7 | 8 |
| 12 | 258 | 8 | 8 |
| 258 | 258 | | |

9 rows selected.

Здесь видно, что наш сеанс параллельного запроса (SID=258) имеет идентификатор SID координатора запроса (query coordinator SID — QCSID) для девяти строк в этом динамическом представлении производительности. Наш сеанс теперь является *координирующим* или управляющим ресурсами этого параллельного запроса. Мы можем видеть, что каждый сеанс имеет собственный SID; фактически каждый из них — это отдельный сеанс Oracle и, как таковой, он отражается в V\$SESSION во время выполнения параллельного запроса:

```
EODA@ORA12CR1> select sid, username, program
2  from v$session
3  where sid in ( select sid
4                  from v$px_session
5                  where qcsid = 258 )
6  /
```

| SID | USERNAME | PROGRAM |
|-----|----------|-----------------------------|
| 12 | EODA | oracle@heera07 (P007) |
| 23 | EODA | oracle@heera07 (P004) |
| 26 | EODA | oracle@heera07 (P000) |
| 94 | EODA | oracle@heera07 (P001) |
| 102 | EODA | oracle@heera07 (P005) |
| 169 | EODA | oracle@heera07 (P006) |
| 177 | EODA | oracle@heera07 (P002) |
| 258 | EODA | sqlplus@heera07 (TNS V1-V3) |
| 267 | EODA | oracle@heera07 (P003) |

9 rows selected.

9 строк выбрано.

На заметку! Если в системе параллельное выполнение не происходит, то серверов параллельного выполнения в представлении V\$SESSION вы не увидите. Они присутствуют в представлении V\$PROCESS, но не будут иметь установившихся сеансов, пока не начнется их использование. Серверы параллельного выполнения будут подключены к базе данных, но без связанных с ними сеансов. Подробные сведения об отличиях между сеансом и соединением приведены в главе 5.

В сущности именно так работает параллельный запрос — и в действительности параллельное выполнение в целом. Он влечет за собой появление последовательности серверов параллельного выполнения, работающих совместно для получения промежуточных результатов, которые передаются либо другим серверам параллельного выполнения, либо координатору параллельного запроса.

В этом конкретном примере у нас есть таблица BIG_TABLE, разнесенная по четырем отдельным устройствам в одном табличном пространстве (табличном пространстве с четырьмя файлами данных). При внедрении параллельного выполнения обычно оптимально иметь данные, разнесенные по максимально возможному числу физических устройств. Достичь этого можно несколькими способами.

- Применение расщепления RAID между дисками.
- Использование ASM со встроенным расщеплением.
- Применение секционирования для физического разнесения BIG_TABLE по нескольким дискам.

- Использование множества файлов данных в одном табличном пространстве, что позволит Oracle размещать экстенды для сегмента `BIG_TABLE` во многих файлах.

В общем случае параллельное выполнение работает лучше всего, когда имеется доступ к как можно большему набору ресурсов (процессору, памяти и системе ввода-вывода). Тем не менее, это не значит, что нельзя получить выигрыш от параллельного запроса, если все данные находятся на одном диске, хотя вероятно он не будет настолько значительным, как в случае множества дисков. Причина возможного получения более короткого времени отклика даже в ситуации с единственным диском, связана с тем, что когда определенный сервер параллельного выполнения подсчитывает строки, он не читает их — и наоборот. Таким образом, два сервера параллельного выполнения могут быть в состоянии выполнить подсчет всех строк за меньшее время, чем при последовательном плане выполнения.

Подобным образом получить выигрыш от параллельного запроса можно даже на однопроцессорной машине. Сомнительно, что последовательный оператор `SELECT COUNT (*)` потребует 100% ресурсов процессора на однопроцессорной машине — часть времени будет тратиться на выполнение (и ожидание) физического ввода-вывода на диск. Параллельный запрос позволил бы полностью утилизировать ресурсы машины (в этом случае процессор и ввод-вывод), каковыми бы они ни были.

Последний момент возвращает нас к приведенной ранее цитате из книги *Practical Oracle8i: Building Efficient Databases*: параллельные запросы по существу не являются масштабируемыми. Если вы разрешите четырем сеансам одновременно выполнять запросы с двумя серверами параллельного выполнения на однопроцессорной машине, то вероятно обнаружите, что их времена отклика станут более длительными, чем при последовательной обработке. Чем больше процессов соперничают за ограниченный ресурс, тем дольше они будут обрабатывать запросы.

Помните, что параллельные запросы требуют соблюдения двух условий. Во-первых, решаемая задача должна быть крупной — например, длительно выполняющийся запрос, требующий для своего завершения минут, часов или дней, а не секунд или долей секунды. Это значит, что параллельный запрос будет неподходящим решением для типичной системы OLTP, в которой длительно выполняющиеся задачи отсутствуют. Включение параллельного выполнения на таких системах часто является крайне неудачным действием.

Во-вторых, в вашем распоряжении должно быть достаточно свободных ресурсов, таких как процессор, ввод-вывод и память. При нехватке любого из этих ресурсов параллельный запрос может привести к чрезмерной его утилизации, негативно повлияв на общую производительность и время выполнения.

В Oracle 11g Release 2 и последующих версиях была добавлена новая функциональность, призванная ограничить расход ресурсов: организация очереди параллельных операторов (Parallel Statement Queuing — PSQ). Когда используется PSQ, база данных будет ограничивать количество одновременно выполняющихся параллельных запросов и помещать любые поступающие впоследствии параллельные запросы в очередь выполнения. Когда ресурсы процессора исчерпаны (что измеряется количеством одновременно задействованных серверов параллельного выполнения), база данных будет препятствовать активизации новых запросов. Эти запросы не потеряют неудачу — просто их запуск будет отложен, т.е. они будут поставлены в очередь.

По мере освобождения ресурсов (применяемые серверы параллельного выполнения завершают свои задачи и переходят в режим простоя), база данных начнет выполнять запросы из очереди. В итоге одновременно может выполняться разумное количество параллельных запросов, не перегружающее систему, а последующие запросы будут терпеливо ожидать, пока до них дойдет дело. В целом все запросы получают свои ответы быстрее, но при этом участвует очередь.

Однажды параллельные запросы были признаны обязательными для многих хранилищ данных — просто потому, что в прошлом (скажем, в 1995 году) хранилища данных являлись редкими и обычно имели очень небольшую, специализированную пользовательскую базу. В наши дни хранилища данных встречаются буквально везде и поддерживают сообщества пользователей, размер которых не уступает аналогичным сообществам для многих транзакционных систем. Другими словами, в таких системах в любой заданный момент времени может не оказаться достаточно свободных ресурсов, чтобы включить параллельные запросы. Это вовсе не означает, что параллельное выполнение в таком случае бесполезно — просто оно может служить больше инструментом администратора базы данных, как будет показано в разделе “Параллельный DDL”, чем средством выполнения параллельных запросов.

Параллельный DML

Документация Oracle ограничивает область использования параллельного DML (PDML) только операторами INSERT, UPDATE, DELETE и MERGE (сюда не входит оператор SELECT, как в обычном DML). В случае PDML для обработки INSERT, UPDATE, DELETE или MERGE база данных Oracle может применять много серверов параллельного выполнения вместо одного последовательного процесса. В многопроцессорной машине с широкой пропускной способностью ввода-вывода потенциальный рост скорости массовых операций DML может оказаться огромным.

Тем не менее, вы *не* должны рассматривать PDML как средство ускорения приложений OLTP. Как утверждалось ранее, параллельные операции спроектированы для максимизации использования машины. Они задуманы так, что одиночный пользователь может полностью задействовать все диски, процессоры и память машины. В определенном хранилище данных (с большим объемом данных и несколькими пользователями) это может быть как раз то, чего желательно достичь. В системе OLTP (с многочисленными пользователями, которые выполняют короткие и быстрые транзакции) пользователям не следует предоставлять возможность полного захвата ресурсов машины.

Это звучит противоречиво: мы применяем параллельные запросы для масштабирования, так как же они могут быть немасштабируемыми? Данное утверждение особенно точно в отношении системы OLTP. Параллельный запрос не является тем, что может масштабироваться по мере увеличения количества одновременно работающих пользователей. Параллельный запрос призван позволить одиночному сеансу сгенерировать столько работы, сколько могли бы породить 100 одновременных сеансов. В системе OLTP мы определенно не хотим, чтобы один пользователь загрузил базу работой как целая сотня пользователей.

Средства PDML полезны в средах крупных хранилищ данных для облегчения пакетного обновления массивных объемов данных. Операция PDML выполняется во многом так же, как выполнялся бы распределенный запрос Oracle, при этом каждый

сервер параллельного выполнения действует подобно процессу в отдельном экземпляре базы данных. Каждый срез таблицы модифицируется отдельным потоком с собственной независимой транзакцией (и, следовательно, с собственным сегментом отмены, как можно надеяться). После того, как все они закончат работу, выполняется эквивалент быстрой двухфазной фиксации для сохранения результатов отдельных независимых транзакций. На рис. 14.2 показано параллельное обновление, использующее четыре сервера параллельного выполнения. Каждый сервер параллельного выполнения имеет собственную независимую транзакцию; эти транзакции либо все фиксируются координирующим сеансом PDML, либо не фиксируется ни одна из них.

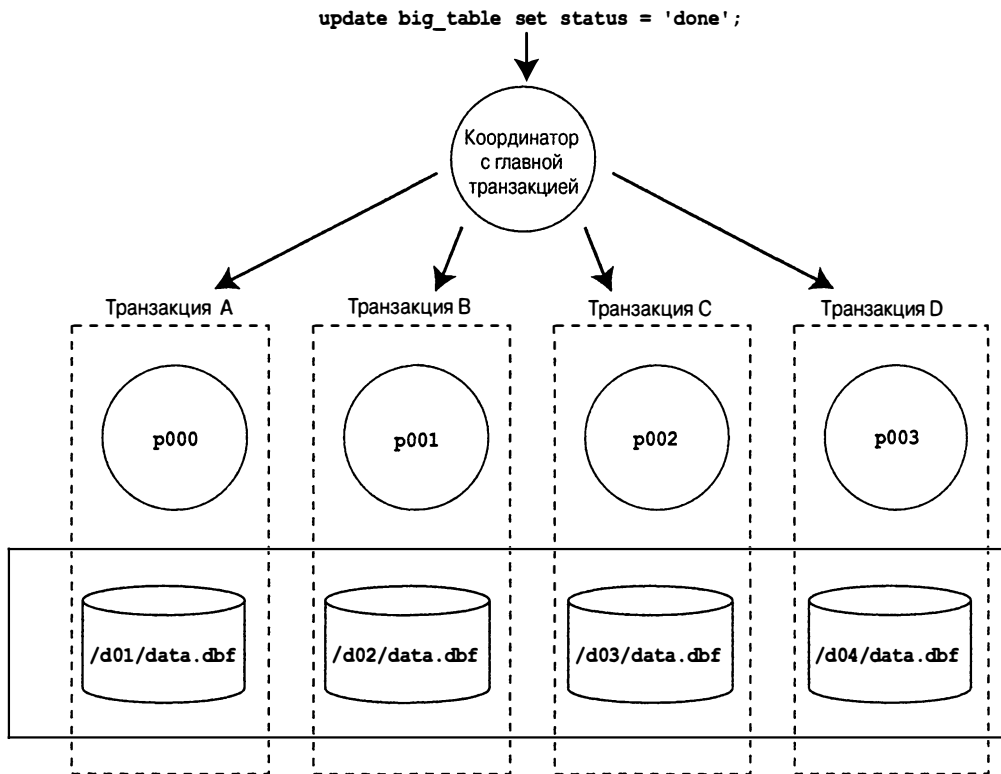


Рис. 14.2. Схема параллельного обновления (PDML)

На самом деле факт создания отдельных независимых транзакций для серверов параллельного выполнения можно отследить. Мы снова будем применять два сеанса. В сеансе с SID=258 мы явно включим параллельный DML. В этом отношении PDML отличается от параллельных запросов — если он явно не затребован, то вы его *не* получите:

```

EODA@ORA12CR1> alter session enable parallel dml;
Session altered.
Сеанс изменен.

```

Вот как удостовериться в том, что параллельный DML был включен для вашего сеанса:

```

EODA@ORA12CR1> select pdml_enabled from v$session where sid = sys_
context('userenv','sid');

PDM
---
YES

```

Того факта, что таблица является “параллельной”, недостаточно, как это было для параллельного запроса. Причина необходимости явного включения PDML в сеансе состоит в том, что с PDML связаны определенные ограничения, которые будут перечислены после рассмотрения примера.

В том же самом сеансе мы запускаем массовую операцию UPDATE, которая фактически выполнится параллельно, т.к. таблица объявлена “параллельной”:

```

EODA@ORA12CR1> update big_table set status = 'done';

```

В другом сеансе мы осуществляем соединение представления V\$SESSION с V\$TRANSACTION, чтобы отобразить активные сеансы для операции PDML, а также информацию об их независимых транзакциях:

```

EODA@ORA12CR1> select a.sid, a.program, b.start_time, b.used_ublk,
2      b.xidusn || '.' || b.xidslot || '.' || b.xidsqn trans_id
3   from v$session a, v$transaction b
4  where a.taddr = b.addr
5    and a.sid in ( select sid
6                    from v$px_session
7                    where qcsid = 258)
8  order by sid
9  /

```

| SID | PROGRAM | START_TIME | USED_UBLK | TRANS_ID |
|-----|-----------------------------|-------------------|-----------|-----------|
| 11 | oracle@heera07 (P00B) | 02/25/14 14:10:17 | 13985 | 26.32.15 |
| 12 | oracle@heera07 (P000) | 02/25/14 14:10:17 | 1 | 70.16.6 |
| 21 | oracle@heera07 (P00F) | 02/25/14 14:10:17 | 13559 | 20.18.37 |
| 23 | oracle@heera07 (P007) | 02/25/14 14:10:17 | 1 | 12.3.62 |
| 26 | oracle@heera07 (P004) | 02/25/14 14:10:17 | 1 | 33.4.11 |
| 95 | oracle@heera07 (P005) | 02/25/14 14:10:17 | 1 | 48.15.10 |
| 97 | oracle@heera07 (P00C) | 02/25/14 14:10:17 | 12676 | 9.5.1730 |
| 103 | oracle@heera07 (P008) | 02/25/14 14:10:17 | 14434 | 44.32.10 |
| 105 | oracle@heera07 (P001) | 02/25/14 14:10:17 | 1 | 64.0.9 |
| 169 | oracle@heera07 (P002) | 02/25/14 14:10:17 | 1 | 34.19.11 |
| 176 | oracle@heera07 (P00D) | 02/25/14 14:10:17 | 14621 | 4.22.1739 |
| 177 | oracle@heera07 (P006) | 02/25/14 14:10:17 | 1 | 74.14.6 |
| 191 | oracle@heera07 (P009) | 02/25/14 14:10:17 | 13070 | 54.11.10 |
| 258 | sqlplus@heera07 (TNS V1-V3) | 02/25/14 14:10:17 | 1 | 59.8.12 |
| 261 | oracle@heera07 (P00A) | 02/25/14 14:10:17 | 13521 | 7.13.1748 |
| 263 | oracle@heera07 (P00E) | 02/25/14 14:10:17 | 12186 | 14.23.76 |
| 267 | oracle@heera07 (P003) | 02/25/14 14:10:17 | 1 | 28.23.19 |

17 rows selected.
17 строк выбрано.

Как видите, здесь происходит больше, чем при простом параллельном запросе таблицы. Эту операцию выполняют 17 процессов, а не 9, как было раньше. Причина в том, что разработанный план включает шаг обновления таблицы и независимые шаги обновления элементов индекса.

Взглянем на план выполнения BASIC плюс PARALLEL, выведенный с помощью пакета DBMS_XPLAN:

```
EODA@ORA12CR1> explain plan for update big_table set status = 'done';
Explained.
Объяснено.
```

```
EODA@ORA12CR1> select * from table(dbms_xplan.display(null,null,
  ↳'BASIC +PARALLEL'));
```

Вот как он выглядит:

| Id | Operation | Name | TQ | IN-OUT | PQ Distrib |
|----|---------------------|-----------|-------|--------|------------|
| 0 | UPDATE STATEMENT | | | | |
| 1 | PX COORDINATOR | | | | |
| 2 | PX SEND QC (RANDOM) | :TQ10001 | Q1,01 | P->S | QC (RAND) |
| 3 | INDEX MAINTENANCE | BIG_TABLE | Q1,01 | PCWP | |
| 4 | PX RECEIVE | | Q1,01 | PCWP | |
| 5 | PX SEND RANGE | :TQ10000 | Q1,00 | P->P | RANGE |
| 6 | UPDATE | BIG_TABLE | Q1,00 | PCWP | |
| 7 | PX BLOCK ITERATOR | | Q1,00 | PCWC | |
| 8 | TABLE ACCESS FULL | BIG_TABLE | Q1,00 | PCWP | |

Как результат псевдораспределенной реализации PDML, этому средству присущи определенные ограничения.

- Во время операции PDML триггеры не поддерживаются. По моему мнению, это разумное ограничение, т.к. триггеры обычно добавляют к обновлению большой объем накладных расходов, а вы используете PDML для достижения высокой скорости — эти два средства не сочетаются друг с другом.
- Существуют определенные декларативные ограничения ссылочной целостности, которые не поддерживаются во время PDML, поскольку каждый срез таблицы модифицируется как отдельная транзакция в отдельном сеансе. Например, не поддерживается рефлексивная ссылочная целостность. Если бы она поддерживалась, возникли бы взаимоблокировки и другие проблемы, связанные с блокировкой.
- Получать доступ к таблице, модифицируемой PDML, нельзя до тех пор, пока не произойдет фиксация или откат.
- Вместе с PDML не поддерживается расширенная репликация (потому что реализация расширенной репликации основана на триггерах).
- При выполнении PDML отложенные ограничения (т.е. ограничения, функционирующие в отложенном режиме) не поддерживаются.
- В отношении таблицы, имеющей битовые индексы или столбцы LOB, операции PDML могут выполняться, только если таблица секционирована, и степень параллелизма будет ограничен количеством секций. В этом случае распараллеливать операции в пределах одной секции невозможно, т.к. каждая секция получает отдельный сервер параллельного выполнения. Следует отметить, что, начиная с версии Oracle 12c, можно запускать операции PDML на столбцах LOB вида SECUREFILE без секционирования.

- При выполнении PDML распределенные транзакции не поддерживаются.
- Кластеризованные таблицы вместе с PDML не поддерживаются.

Если вы нарушите любое из этих ограничений, произойдет одно из двух: либо оператор будет выполнен последовательно (безо всякого параллелизма), либо возникнет ошибка. Например, если вы уже запустили операцию PDML на таблице T и пытаетесь выполнить запрос к этой таблице до окончания транзакции, то получите ошибку ORA-12838: cannot read/modify an object after modifying it in parallel (ORA-12838: не удастся прочитать/модифицировать объект после его модификации в параллельном режиме).

В отношении предшествующего примера также полезно отметить, что если не включить параллельный DML для оператора UPDATE, то план выполнения будет выглядеть несколько по-другому:

| Id | Operation | Name | TQ | IN-OUT | PQ Distrib |
|----|---------------------|-----------|-------|--------|------------|
| 0 | UPDATE STATEMENT | | | | |
| 1 | UPDATE | BIG_TABLE | | | |
| 2 | PX COORDINATOR | | | | |
| 3 | PX SEND QC (RANDOM) | :TQ10000 | Q1,00 | P->S | QC (RAND) |
| 4 | PX BLOCK ITERATOR | | Q1,00 | PCWC | |
| 5 | TABLE ACCESS FULL | BIG_TABLE | Q1,00 | PCWP | |

Для неопытного взгляда может показаться, что оператор UPDATE выполняется параллельно, но на самом деле это не так. Вывод показывает, что оператор UPDATE является последовательным, но полное сканирование (чтение) таблицы производится параллельно. Таким образом, здесь задействован параллельный запрос, а не PDML.

Контроль параллельных операций

Вы можете быстро проконтролировать параллельные операции, которые произошли в сеансе, запросив словарь данных. Например, мы имеем следующую параллельную операцию DML:

```
EODA@ORA12CR1> alter session enable parallel dml;
EODA@ORA12CR1> update big_table set status='AGAIN';
```

Ниже показано, как проверить тип и количество параллельных действий:

```
EODA@ORA12CR1> select name, value from v$statname a, v$mystat b
2 where a.statistic# = b.statistic# and name like '%parallel%';
```

Вот часть вывода для этого сеанса:

| NAME | VALUE |
|--|-------|
| DBWR parallel query checkpoint buffers written | 0 |
| queries parallelized | 0 |
| DML statements parallelized | 1 |
| DDL statements parallelized | 0 |
| DFO trees parallelized | 1 |

Вывод подтверждает, что в сеансе был выполнен один параллельный оператор DML.

Параллельный DDL

Я уверен, что параллельный DDL является действительно ярким аспектом технологии параллельного выполнения Oracle. Как уже обсуждалось, параллельное выполнение обычно не подходит для систем OLTP. В действительности для многих хранилищ данных параллельные запросы становятся все менее и менее приемлемым вариантом. Раньше считалось, что хранилище данных строилось для небольшого, специализированного сообщества пользователей, иногда включающего всего одного или двух аналитиков. Однако на протяжении последнего десятилетия я наблюдал за их ростом до сообществ, насчитывающих сотни и тысячи пользователей. Представьте хранилище данных с пользовательским интерфейсом в виде веб-приложения: оно может быть доступно буквально тысячам пользователей одновременно.

Но что касается администраторов баз данных, выполняющих крупные пакетные операции, возможно, в рамках окна обслуживания, то здесь совсем другая история. Администратор базы данных остается единственным лицом, в распоряжении которого находится огромная машина с громадным объемом вычислительных ресурсов. Администратор базы данных должен делать только что-то одно: загрузить нужные данные, реорганизовать определенную таблицу или перестроить какой-нибудь индекс. Без параллельного выполнения он вряд ли сумеет в полной мере воспользоваться всеми возможностями оборудования, а с параллельным выполнением — сможет. Распараллеливание допускают перечисленные ниже команды DDL.

- **CREATE INDEX.** Множество серверов параллельного выполнения могут сканировать таблицу, сортировать данные и записывать отсортированные сегменты в индексную структуру.
- **CREATE TABLE AS SELECT.** Запрос, выполняющий SELECT, может функционировать в параллельном режиме, и сама загрузка таблицы также может осуществляться параллельно.
- **ALTER INDEX REBUILD.** Структуру индекса можно перестраивать параллельно.
- **ALTER TABLE MOVE.** Таблицу можно перемещать параллельно.
- **ALTER TABLE SPLIT|COALESCE PARTITION.** Индивидуальные секции таблицы могут расщепляться и перераспределяться параллельно.
- **ALTER INDEX SPLIT PARTITION.** Секция индекса может расщепляться параллельно.
- **CREATE/ALTER MATERIALIZED VIEW.** С помощью параллельных процессов можно создавать материализованное представление либо изменять стандартную степень параллелизма.

На заметку! Полный список операторов, которые поддерживают параллельные операции, приведен в руководстве по языку SQL базы данных Oracle (*Oracle Database SQL Language Reference*).

Первые четыре команды работают также и с отдельными секциями таблиц/индексов, т.е. можно выполнять перемещение индивидуальных секций таблицы в параллельном режиме.

Для меня параллельный DDL — это то, где преимущества параллельного выполнения Oracle наиболее заметны. Несомненно, он может применяться с параллельными запросами для ускорения определенных длительно выполняющихся операций, но с точки зрения обслуживания и администрирования параллельный DDL является той областью, где параллельные операции затрагивают администраторов баз данных и разработчиков в наибольшей мере. Если считать, что параллельные запросы спроектированы в большей степени для конечного пользователя, то параллельный DDL предназначен для администраторов баз данных и разработчиков.

Параллельный DDL и загрузка данных с использованием внешних таблиц

Одним из моих излюбленных средств, появившихся в Oracle9i, были внешние таблицы, которые особенно удобны в области загрузки данных. Загрузка данных и внешние таблицы подробно обсуждаются в следующей главе, а здесь мы кратко обратимся к этим темам, чтобы изучить эффекты, оказываемые параллельным DDL на установку размеров и усечение экстентов.

Внешние таблицы позволяют легко выполнять параллельную прямую загрузку, не особенно о ней задумываясь. В версии Oracle 7.1 была предоставлена возможность выполнения параллельной прямой загрузки, посредством которой множество сеансов могут производить запись непосредственно в файлы данных Oracle, полностью минуя буферный кеш, генерирование информации отмены для табличных данных и, возможно, даже генерирование информации повторения. Это достигалось через SQL*Loader. Администратор базы данных должен был написать сценарий для множества сеансов загрузки SQL*Loader, разбить входные файлы данных для ручной загрузки, определить степень параллелизма и координировать все процессы SQL*Loader. Короче говоря, это можно было сделать, но совсем нелегко.

Благодаря параллельному DDL или параллельному DML в сочетании с внешними таблицами мы получаем параллельную прямую загрузку, которая реализуется посредством простого оператора `CREATE TABLE AS SELECT` или `INSERT /*+ APPEND */`. Больше не требуется написание сценариев, разбиение файлов и координация *N* сценариев, которые должны быть выполнены. Другими словами, такая комбинация обеспечивает легкость применения без каких-либо потерь в производительности.

Давайте взглянем на простой пример. Вскоре вы увидите, как создается внешняя таблица. (Загрузка данных из внешних таблиц более подробно рассматривается в следующей главе.) В настоящий момент для загрузки данных из другой таблицы мы будем использовать реальную таблицу, что очень похоже на то, как многие поступают с промежуточными таблицами в хранилище данных. Этот прием кратко описан далее.

1. Создать входные файлы с применением некоторого инструмента извлечения, преобразования и загрузки (extract, transform, load — ETL).
2. Загрузить эти входные файлы в промежуточные таблицы.
3. Загрузить новую таблицу с использованием запросов в отношении промежуточных таблиц.

Мы будем применять ту же самую представленную ранее таблицу `BIG_TABLE`, которая готова для параллельной обработки и содержит 10 миллионов записей. Мы соединим ее со второй таблицей, `USER_INFO`, хранящей связанную с `OWNER` инфор-

мацию из представления словаря данных ALL_USERS. Цель заключается в денормализации этой информации и превращении ее в плоскую структуру.

Для начала мы создадим таблицу USER_INFO, разрешим в ней параллельные операции и затем соберем статистику:

```
EOGA@ORA12CR1> create table user_info as select * from all_users;
Table created.
Таблица создана.

EOGA@ORA12CR1> alter table user_info parallel;
Table altered.
Таблица изменена.

EOGA@ORA12CR1> exec dbms_stats.gather_table_stats( user, 'USER_INFO' );
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.
```

Теперь мы хотели бы выполнить параллельную прямую загрузку новой таблицы этой информацией. Запрос, который понадобится для этого, прост:

```
create table new_table parallel
as
select a.*, b.user_id, b.created user_created
  from big_table a, user_info b
 where a.owner = b.username;
```

План выполнения этого оператора CREATE TABLE AS SELECT в Oracle 12c выглядит следующим образом:

| ----- | | | | | | |
|-------|--------------------------------|-----------|----------|-------|--------|------------|
| Id | Operation | Name | Time | TQ | IN-OUT | PQ Distrib |
| ----- | | | | | | |
| 0 | CREATE TABLE STATEMENT | | 00:00:01 | | | |
| 1 | PX COORDINATOR | | | | | |
| 2 | PX SEND QC (RANDOM) | :TQ10000 | 00:00:01 | Q1,00 | P->S | QC (RAND) |
| 3 | LOAD AS SELECT | NEW_TABLE | | Q1,00 | PCWP | |
| 4 | OPTIMIZER STATISTICS GATHERING | | 00:00:01 | Q1,00 | PCWP | |
| 5 | HASH JOIN | | 00:00:01 | Q1,00 | PCWP | |
| 6 | TABLE ACCESS FULL | USER_INFO | 00:00:01 | Q1,00 | PCWP | |
| 7 | PX BLOCK ITERATOR | | 00:00:01 | Q1,00 | PCWC | |
| 8 | TABLE ACCESS FULL | BIG_TABLE | 00:00:01 | Q1,00 | PCWP | |
| ----- | | | | | | |

Если взглянуть на шаги, начиная с пятого, то это компонент запроса (SELECT). Сканирование BIG_TABLE и хеш-соединение с USER_INFO было выполнено параллельно, и каждый из промежуточных результатов загружен в часть таблицы (шаг 3, LOAD AS SELECT). После того, как каждый сервер параллельного выполнения завершил свою порцию соединения и загрузки, он отправил результаты координатору запроса. В этом случае результаты просто указывают на “успех” или “отказ”, т.к. работа уже сделана.

Вот и все; параллельная прямая загрузка осуществляется легко. Наиболее важный момент, который следует учитывать в этих операциях, связан с тем, как используется (или не используется) пространство. Особенно важен побочный эффект, который называется *усечением экстенгов* (extent trimming). Давайте посвятим некоторое время его исследованию.

Параллельный DDL и усечение экстентов

Параллельный DDL полагается на операции прямого пути. Другими словами, данные не передаются в буферный кеш для последующей записи. Вместо этого операция вроде `CREATE TABLE AS SELECT` создаст новые экстенты и произведет запись непосредственно в них, и данные следуют из запроса прямо на диск в эти вновь выделенные экстенты. Каждый сервер параллельного выполнения при обработке своей части операции `CREATE TABLE AS SELECT` будет осуществлять запись в собственный экстенст. Операция `INSERT /*+ APPEND */` (прямая вставка) производит запись “за пределами” HWM-маркера сегмента, и каждый сервер параллельного выполнения будет записывать в собственный набор экстентов, никогда не разделяя их с другими серверами параллельного выполнения.

Таким образом, если происходит параллельная операция `CREATE TABLE AS SELECT` и для создания таблицы применяются четыре сервера параллельного выполнения, то получится *минимум* четыре экстента, а может быть и больше. Но каждый из серверов параллельного выполнения будет выделять собственный экстенст, производить в него запись и, когда экстенст заполнится, выделять очередной новый экстенст. Сервер параллельного выполнения никогда не будет использовать экстенст, выделенный другим сервером параллельного выполнения.

Этот процесс изображен на рис. 14.3. Здесь оператор `CREATE TABLE NEW TABLE AS SELECT` обрабатывается четырьмя серверами параллельного выполнения. Каждый сервер параллельного выполнения представлен своим оттенком (белым, светло-серым, темно-серым и черным). Прямоугольники на барабане диска представляют экстенты, созданные в некотором файле данных этим оператором `CREATE TABLE`. Каждый экстенст представлен одним из упомянутых выше четырех цветов, указывая на то, что все данные в любом экстенсте были загружены только одним из четырех серверов параллельного выполнения. На рисунке видно, что сервер p003 создал и загрузил четыре экстента, сервер p000 — пять экстентов и т.д.

На первых порах подход выглядит правильным, но в среде хранилища данных он может привести к излишним расходам пространства после крупномасштабной загрузки. Предположим, что вы хотите загрузить 1 010 Мбайт (около 1 Гбайт) данных и применяете табличное пространство с экстентами 100 Мбайт. Для загрузки этих данных вы решили использовать 10 серверов параллельного выполнения. Каждый сервер должен начинать работу с выделения собственного экстенста размером 100 Мбайт (всего их будет 10) и его заполнения.

Поскольку каждому серверу требуется загрузить 101 Мбайт данных, он заполнит свой первый экстенст и затем выделит еще один экстенст размером 100 Мбайт, в котором будет занято только 1 Мбайт. Теперь имеется 20 экстентов (10 полных и 10 с занятым пространством 1 Мбайт), в которых пространство размером 990 Мбайт выделено, но не задействовано. Оно может быть занято в следующий раз при загрузке дополнительных данных, но в настоящий момент 990 Мбайт расходуются впустую. Именно здесь на помощь приходит усечение экстентов. База данных Oracle будет брать последние экстенты, принадлежащие серверам параллельного выполнения, и усекать их до минимально возможных размеров.

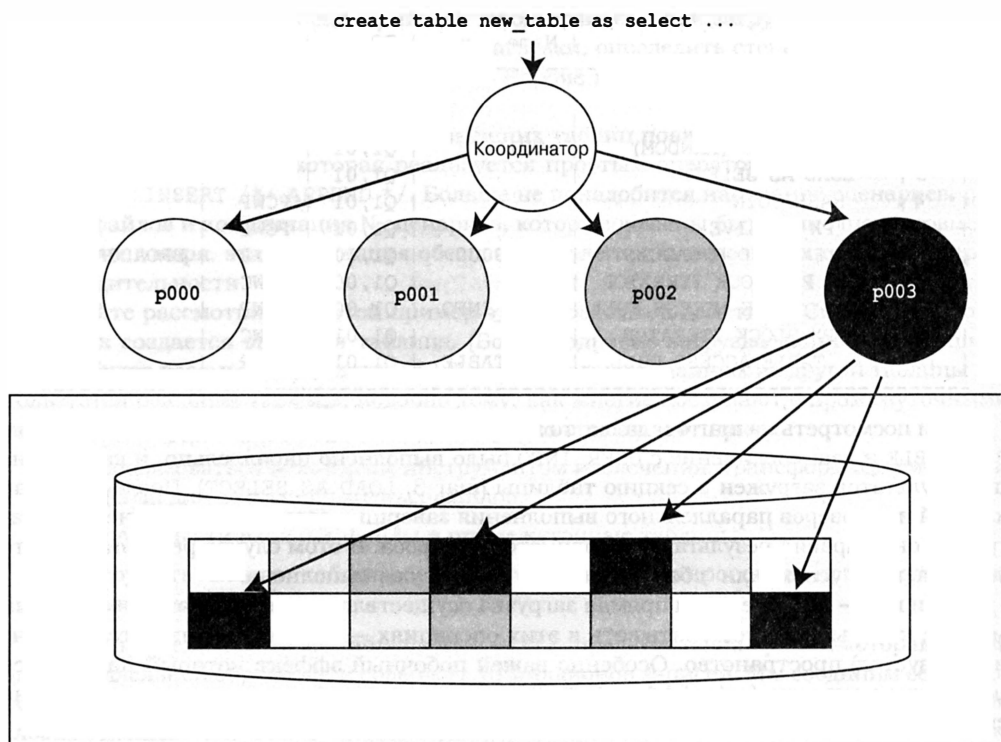


Рис. 14.3. Схема выделения экстенгов параллельным DDL

Усечение экстенгов и табличные пространства, управляемые словарем

Если вы применяете унаследованные табличные пространства, управляемых словарем, то Oracle сможет преобразовать экстенги размером 100 Мбайт, содержащие всего 1 Мбайт данных, в экстенги размером 1 Мбайт. К сожалению, это обычно приводит (в управляемых словарем табличных пространствах) к тому, что остаются свободными 10 несмежных экстенгов по 99 Мбайт, а поскольку схема выделения рассчитана на экстенги размером 100 Мбайт, эти 990 Мбайт пространства окажутся не особенно полезными! Следующее выделение 100 Мбайт, вероятно, *не* сможет использовать существующее пространство, т.к. будет доступно только 99 Мбайт свободного пространства, за которым следует 1 Мбайт занятого пространства, далее снова 99 Мбайт свободного пространства и т.д. Подход с управлением словарем в этой книге дополнительно рассматриваться не будет.

Усечение экстенгов и локально управляемые табличные пространства

Обратимся к локально управляемым табличным пространствам. Различают два их типа: `UNIFORM SIZE`, при котором каждый экстент в табличном пространстве имеет строго один и тот же размер, и `AUTOALLOCATE`, когда Oracle с помощью внутреннего алгоритма решает, насколько большим должен быть каждый экстент. Оба подхода изящно решают проблему вида *свободное пространство размером 99 Мбайт, за которым следует занятое пространство размером 1 Мбайт, затем снова свободное пространство 99 Мбайт*.

Однако каждый тип локально управляемого табличного пространства решает ее совершенно по-другому. Подход `UNIFORM SIZE` вообще исключает усечение экстен-тов из рассмотрения. Когда применяется табличное пространство `UNIFORM SIZE`, усечение экстен-тов не может выполняться. Все экстен-ты имеют единый размер — ни один из них не может быть меньше (или больше) этого единого размера. С другой стороны, табличное пространство `AUTOALLOCATE` поддерживает усечение экстен-тов, но интеллектуальным образом. В нем используется несколько специфических размеров экстен-тов и есть возможность применять пространство разного размера, т.е. алгоритм допускает использование со временем всего свободного места в табличном пространстве. В отличие от табличного пространства, управляемого словарем, где запрос экстен-та 100 Мбайт будет отклонен, когда удастся найти свободные экстен-ты только размером 99 Мбайт (близким, но недостаточным), локально управляемое табличное пространство `AUTOALLOCATE` может оказаться более гибким. Оно способно сокращать размер запрошенного пространства, чтобы попытаться задействовать все свободное место.

Теперь давайте посмотрим на отличия между двумя подходами в локально управляемом табличном пространстве. Для этого понадобится реалистичный пример. Мы подготовим внешнюю таблицу, которую можно будет применять в ситуации с прямой загрузкой, что приходится делать довольно часто. Даже если вы по-прежнему используете `SQL*Loader` для прямой загрузки данных, сведения в этом разделе полностью актуальны — нужно только иметь вручную написанный сценарий для выполнения действительной загрузки данных. Итак, чтобы исследовать усечение экстен-тов, необходимо настроить загрузку для рассматриваемого примера, затем произвести ее в разных условиях и исследовать полученные результаты.

Настройка локально управляемых табличных пространств

Для начала нам потребуется внешняя таблица. То и дело я обнаруживаю у себя унаследованный управляющий файл из `SQL*Loader`, который привык применять для загрузки данных — один из тех файлов, что выглядят следующим образом:

```
LOAD DATA
  INFILE '/tmp/big_table.dat'
  INTO TABLE big_table
  REPLACE
  FIELDS TERMINATED BY '|'
  ( id ,owner ,object_name ,subobject_name ,object_id
  ,data_object_id ,object_type ,created ,last_ddl_time
  ,timestamp ,status ,temporary ,generated ,secondary
  ,namespace ,edition_name
  )
```

Его можно *без труда* преобразовать в определение внешней таблицы с использованием самого инструмента `SQL*Loader`:

```
$ sqlldr eoda/foo big_table.ctl external_table=generate_only
SQL*Loader: Release 12.1.0.1.0 - Production on Mon Feb 17 14:39:21 2014
Copyright (c) 1982, 2013, Oracle and/or its affiliates. All rights reserved.
```

На заметку! Команда `SQLLDR` вместе с ее параметрами подробно рассматривается в главе 15.

Обратите внимание на параметр `EXTERNAL TABLE`, передаваемый `SQL*Loader`. В этом случае он указывает инструменту `SQL*Loader` о том, что нужно не загружать данные, а сгенерировать оператор `CREATE TABLE` в журнальном файле. Оператор `CREATE TABLE` имеет показанный ниже вид (это сокращенная форма; для экономии места повторяющиеся элементы были опущены):

```
CREATE TABLE "SYS_SQLLDR_X_EXT_BIG_TABLE"
(
  "ID" NUMBER,
  ...
  "EDITION_NAME" VARCHAR2(128)
)
ORGANIZATION external
(
  TYPE oracle_loader
  DEFAULT DIRECTORY my_dir
  ACCESS PARAMETERS
  (
    RECORDS DELIMITED BY NEWLINE CHARACTERSET US7ASCII
    BADFILE 'SYS_SQLLDR_XT_TMPDIR_00001': 'big_table.bad'
    LOGFILE 'SYS_SQLLDR_XT_TMPDIR_00001': 'big_table.log_xt'
    READSIZE 1048576
    FIELDS TERMINATED BY "|" LDRTRIM
    REJECT ROWS WITH ALL NULL FIELDS
    (
      "ID" CHAR(255)
        TERMINATED BY "|",
      ...
      "EDITION_NAME" CHAR(255)
        TERMINATED BY "|"
    )
  )
  location
  (
    'big_table.dat'
  )
) REJECT LIMIT UNLIMITED
```

Все, что понадобится здесь предпринять — это отредактировать желаемым образом имя внешней таблицы, возможно, изменить каталоги и т.д.:

```
EODA@ORA12CR1> create or replace directory my_dir as '/tmp/';
Directory created.
Каталог создан.
```

После этого остается только действительно создать внешнюю таблицу:

```
1 CREATE TABLE "BIG_TABLE_ET"
2 (
3   "ID" NUMBER,
...
18  "EDITION_NAME" VARCHAR2(128)
19 )
20 ORGANIZATION external
21 (
```

```

22 TYPE oracle_loader
23 DEFAULT DIRECTORY my_dir
24 ACCESS PARAMETERS
25 (
26   records delimited by newline
27   fields terminated by ','
28   missing field values are null
29 )
30 location
31 (
32   'big_table.dat'
33 )
34 )REJECT LIMIT UNLIMITED
35 /
Table created.

```

Далее необходимо разрешить параллельные операции для этой таблицы. В этом и заключается “магия” — именно данный шаг облегчит параллельную прямую загрузку:

```

EODA@ORA12CR1> alter table big_table_et PARALLEL;
Table altered.
Таблица изменена.

```

На заметку! Конструкция PARALLEL может также применяться в самом операторе CREATE TABLE. Ключевое слово PARALLEL можно было бы добавить сразу после конструкции REJECT LIMIT UNLIMITED. Оператор ALTER использовался только с целью привлечения внимания к тому факту, что внешняя таблица на самом деле допускает параллельные операции.

Усечение экстенгов в локально управляемых табличных пространствах UNIFORM SIZE и AUTOALLOCATE

Это и все, что нужно было сделать для настройки компонента загрузки. Теперь мы хотим сравнить особенности управления пространством в локально управляемых табличных пространствах (locally-managed tablespace — LMT) типа UNIFORM SIZE и AUTOALLOCATE. В рассматриваемом случае мы будем применять экстенги размером 100 Мбайт. Для начала мы создадим табличное пространство по имени LMT_UNIFORM, использующее экстенги унифицированного размера (UNIFORM SIZE):

```

EODA@ORA12CR1> create tablespace lmt_uniform
2   datafile '/u01/dbfile/ORA12CR1/lmt_uniform.dbf' size 1048640K reuse
3   autoextend on next 100m
4   extent management local
5   UNIFORM SIZE 100m;
Tablespace created.

```

Далее мы создадим табличное пространство по имени LMT_AUTO, в котором для определения размера экстенгов применяется подход AUTOALLOCATE:

```

EODA@ORA12CR1> create tablespace lmt_auto
2   datafile '/u01/dbfile/ORA12CR1/lmt_auto.dbf' size 1048640K reuse
3   autoextend on next 100m
4   extent management local
5   AUTOALLOCATE;
Tablespace created.

```

Каждое табличное пространство начинается с файла данных размером 1 Гбайт (плюс 64 Кбайт, которые используются локально управляемыми табличными пространствами для управления хранилищем; в случае размера блока 32 Кбайт дополнительно будет задействовано 128 Кбайт, а не 64 Кбайт). Мы позволим этим файлам данных автоматически расширяться на 100 Мбайт за раз. Мы собираемся загрузить следующий файл с 10 000 000 записей:

```
$ ls -lag big_table.dat
-rw-r----- 1 dba 1018586660 Feb 27 21:27 big_table.dat
```

Он был создан с применением сценария `big_table.sql`, который описан в разделе “Настройка среды” в самом начале книги, и затем выгружен с использованием сценария `flat.sql`, доступного по адресу <http://tkyte.blogspot.com/2009/10/httpasktomoraclecomtkyteflat.html>. Затем мы произведем параллельную прямую загрузку этого файла в каждое табличное пространство:

```
EODA@ORA12CR1> create table uniform_test
  2 parallel
  3 tablespace lmt_uniform
  4 as
  5 select * from big_table_et;
Table created.
Таблица создана.

EODA@ORA12CR1> create table autoallocate_test
  2 parallel
  3 tablespace lmt_auto
  4 as
  5 select * from big_table_et;
Table created.
Таблица создана.
```

В моей четырехпроцессорной системе эти операторы `CREATE TABLE` обрабатывались с помощью восьми серверов параллельного выполнения и одного координатора. Я проверил это путем запроса одного из динамических представлений производительности, относящихся к параллельному выполнению, `V$PX_SESSION`, пока эти операторы функционировали:

```
EODA@ORA12CR1> select sid, serial#, qcsid, qcserial#, degree from v$px_session;
```

| SID | SERIAL# | QCSID | QCserial# | DEGREE |
|-----|---------|-------|-----------|--------|
| 375 | 1275 | 243 | 1697 | 8 |
| 18 | 3145 | 243 | 1697 | 8 |
| 135 | 3405 | 243 | 1697 | 8 |
| 244 | 3065 | 243 | 1697 | 8 |
| 376 | 177 | 243 | 1697 | 8 |
| 5 | 875 | 243 | 1697 | 8 |
| 134 | 2829 | 243 | 1697 | 8 |
| 249 | 467 | 243 | 1697 | 8 |
| 243 | 1697 | 243 | | |

```
9 rows selected.
9 строк выбрано.
```

На заметку! При создании таблиц `UNIFORM_TEST` и `AUTOALLOCATE_TEST` мы просто указали `parallel`, поручив Oracle выбрать степень параллелизма самостоятельно. В этом случае я был единственным пользователем машины (со всеми доступными ресурсами), и база данных Oracle установила стандартную степень параллелизма в 8 на основе количества процессоров (4) и значения параметра `PARALLEL_THREADS_PER_CPU`, по умолчанию равного 2.

В столбцах `SID` и `SERIAL#` приведены идентификаторы сеансов параллельного выполнения, а столбцы `QCSID` и `QCSERIAL#` определяют идентификатор координатора запроса параллельного выполнения. Таким образом, при восьми запущенных сеансах параллельного выполнения мы хотим увидеть, как использовалось пространство. Быстрый запрос к `USER_SEGMENTS` дает хорошее представление об этом:

```
EODA@ORA12CR1> select segment_name, blocks, extents
  2 from user_segments
  3 where segment_name in ( 'UNIFORM_TEST', 'AUTOALLOCATE_TEST' );
```

| SEGMENT_NAME | BLOCKS | EXTENTS |
|-------------------|--------|---------|
| AUTOALLOCATE_TEST | 153696 | 131 |
| UNIFORM_TEST | 166400 | 13 |

Из-за применения размера блока 8 Кбайт разница составила около 100 Мбайт; с точки зрения отношения выделенное пространство в таблице `AUTOALLOCATE_TEST` соответствует примерно 92% выделенного пространства в таблице `UNIFORM_TEST`. Ниже приведены результаты по действительно используемому пространству:

```
EODA@ORA12CR1> exec show_space('UNIFORM_TEST' );
Unformatted Blocks ..... 0
FS1 Blocks (0-25) ..... 0
FS2 Blocks (25-50) ..... 0
FS3 Blocks (50-75) ..... 0
FS4 Blocks (75-100)..... 12,782
Full Blocks ..... 153,076
Total Blocks..... 166,400
Total Bytes..... 1,363,148,800
Total MBytes..... 1,300
Unused Blocks..... 0
Unused Bytes..... 0
Last Used Ext FileId..... 5
Last Used Ext BlockId..... 153,728
Last Used Block..... 12,800
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.
```

```
EODA@ORA12CR1> exec show_space('AUTOALLOCATE_TEST' );
Unformatted Blocks ..... 0
FS1 Blocks (0-25) ..... 0
FS2 Blocks (25-50) ..... 0
FS3 Blocks (50-75) ..... 0
FS4 Blocks (75-100)..... 0
Full Blocks ..... 153,076
Total Blocks..... 153,696
Total Bytes..... 1,259,077,632
```

```
Total MBytes..... 1,200
Unused Blocks..... 0
Unused Bytes..... 0
Last Used Ext FileId..... 6
Last Used Ext BlockId..... 147,584
Last Used Block..... 6,240
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.
```

На заметку! Процедура SHOW_SPACE описана в разделе “Настройка среды” в самом начале книги.

Обратите внимание на то, что в таблице AUTOALLOCATE_TEST занято около 1200 Мбайт пространства, тогда как в таблице UNIFORM_TEST — примерно 1300 Мбайт пространства. Все это объясняется тем, что усечение экстенгов не происходило. Если взглянуть на таблицу UNIFORM_TEST, то можно увидеть это ясно:

```
EOA@ORA12CR1> select segment_name, extent_id, blocks
  2 from user_extents where segment_name = 'UNIFORM_TEST';
```

| SEGMENT_NAME | EXTENT_ID | BLOCKS |
|--------------|-----------|--------|
| UNIFORM_TEST | 0 | 12800 |
| UNIFORM_TEST | 1 | 12800 |
| UNIFORM_TEST | 2 | 12800 |
| UNIFORM_TEST | 3 | 12800 |
| UNIFORM_TEST | 4 | 12800 |
| UNIFORM_TEST | 5 | 12800 |
| UNIFORM_TEST | 6 | 12800 |
| UNIFORM_TEST | 7 | 12800 |
| UNIFORM_TEST | 8 | 12800 |
| UNIFORM_TEST | 9 | 12800 |
| UNIFORM_TEST | 10 | 12800 |
| UNIFORM_TEST | 11 | 12800 |
| UNIFORM_TEST | 12 | 12800 |

13 rows selected.
13 строк выбрано.

Каждый экстенг имеет размер 100 Мбайт. Приводить полный список из 131 экстенгов, выделенных для таблицы AUTOALLOCATE_TEST, было бы напрасной тратой бумаги, так что давайте посмотрим на них в совокупности:

```
EOA@ORA12CR1> select segment_name, blocks, count(*)
  2 from user_extents
  3 where segment_name = 'AUTOALLOCATE_TEST'
  4 group by segment_name, blocks
  5 order by blocks;
```

| SEGMENT_NAME | BLOCKS | COUNT(*) |
|-------------------|--------|----------|
| AUTOALLOCATE_TEST | 1024 | 128 |
| AUTOALLOCATE_TEST | 6240 | 1 |
| AUTOALLOCATE_TEST | 8192 | 2 |

В целом это соответствует тому, каким образом локально управляемые табличные пространства AUTOALLOCATE распределяют пространство (результаты предшествующего запроса будут варьироваться в зависимости от объема данных и версии Oracle). Такие значения, как показанные для экстенстов блоков 1024 и 8192, являются нормальными; для таблиц типа AUTOALLOCATE мы их будем видеть постоянно. Тем не менее, остальные значения не относятся к нормальным; обычно мы их не наблюдаем. Это объясняется происходящим усечением экстенстов. Некоторые из серверов параллельного выполнения завершили свою часть загрузки — они заняли свой последний экстенст 64 Мбайт (8192 блока) и усекли его, образовав небольшой резерв. Какой-то другой сеанс параллельного выполнения при возникновении потребности в пространстве мог бы воспользоваться этим небольшим резервом. В свою очередь, когда эти другие сеансы параллельного выполнения завершат обработку собственных загрузок, они усекут свои последние экстенсты и также оставят небольшие резервы пространства.

Итак, какой же подход применять? Если ваша цель заключается в том, чтобы как можно чаще выполнять прямую загрузку параллельно, то я рекомендую избрать AUTOALLOCATE в качестве политики управления экстенстами. Параллельные операции прямого пути вроде этой не используют пространство до HWM-маркера объекта — пространство в списке свободных блоков. Таким образом, если в эти таблицы также не производятся какие-то обычные вставки, то выделение в стиле UNIFORM SIZE постоянно будет оставлять дополнительное свободное пространство, которое никогда не задействуется. Если размер экстенстов для локально управляемого табличного пространства UNIFORM SIZE не удастся установить в намного меньшее значение, то со временем вы увидите то, что я назвал бы чрезмерной потерей места. К тому же помните, что это пространство ассоциировано с сегментом и будет включено в полное сканирование таблицы.

Чтобы продемонстрировать сказанное, выполним еще одну параллельную прямую загрузку в существующие таблицы с применением тех же входных данных:

```

EODA@ORA12CR1> alter session enable parallel dml;
Session altered.
Сеанс изменен.

EODA@ORA12CR1> insert /*+ append */ into UNIFORM_TEST
  2 select * from big_table_et;
10000000 rows created.
10000000 строк создано.

EODA@ORA12CR1> insert /*+ append */ into AUTOALLOCATE_TEST
  2 select * from big_table_et;
10000000 rows created.
10000000 строк создано.

EODA@ORA12CR1> commit;
Commit complete.
Фиксация выполнена.

```

Если мы сравним утилизацию пространства в двух таблицах после этой операции, то заметим, что по мере загрузки все большего и большего объема данных в таблицу UNIFORM_TEST с помощью параллельных операций прямой загрузки использование пространства с течением времени становится все хуже:

```

EODA@ORA12CR1> exec show_space( 'UNIFORM_TEST' );
Unformatted Blocks ..... 0
FS1 Blocks (0-25) ..... 0
FS2 Blocks (25-50) ..... 0
FS3 Blocks (50-75) ..... 0
FS4 Blocks (75-100)..... 25,564
Full Blocks ..... 306,152
Total Blocks..... 332,800
Total Bytes..... 2,726,297,600
Total MBytes..... 2,600
Unused Blocks..... 0
Unused Bytes..... 0
Last Used Ext FileId..... 5
Last Used Ext BlockId..... 320,128
Last Used Block..... 12,800
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.

```

```

EODA@ORA12CR1> exec show_space( 'AUTOALLOCATE_TEST' );
Unformatted Blocks ..... 0
FS1 Blocks (0-25) ..... 0
FS2 Blocks (25-50) ..... 0
FS3 Blocks (50-75) ..... 0
FS4 Blocks (75-100)..... 0
Full Blocks ..... 306,152
Total Blocks..... 307,392
Total Bytes..... 2,518,155,264
Total MBytes..... 2,401
Unused Blocks..... 0
Unused Bytes..... 0
Last Used Ext FileId..... 6
Last Used Ext BlockId..... 301,312
Last Used Block..... 6,240
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.

```

Мы могли бы выбрать меньший унифицированный размер сегмента либо применить конструкцию `AUTOALLOCATE`. Конструкция `AUTOALLOCATE` также может со временем генерировать больше экстенгов, но утилизация свободного пространства будет намного лучше за счет происходящего усечения экстенгов.

На заметку! Ранее в этой главе я отмечал, что показатели, получаемые вами при выполнении предыдущих примеров реализации параллелизма, могут варьироваться. Полезно подчеркнуть снова: ваши результаты будут изменяться в зависимости от версии Oracle, используемой степени параллелизма и объема загружаемых данных. Показанный выше в разделе вывод был сгенерирован в среде Oracle 12c Release 1 со стандартной степенью параллелизма на четырехпроцессорной машине.

Процедурный параллелизм

Мы обсудим два типа процедурного параллелизма.

- Параллельные конвейерные функции, которые являются средством Oracle.
- Самодельный параллелизм (do-it-yourself — DIY), представляющий собой применение к созданным вами приложениям тех же самых приемов, которые Oracle применяет к параллельному полному сканированию таблиц. Параллелизм DIY — это в большей степени прием разработки, чем что-то встроенное напрямую в Oracle 11g Release 1 и предшествующие версии, но собственное средство базы данных в Oracle 11g Release 2 и последующих версиях.

Вы часто будете обнаруживать, что приложения (обычно пакетные процессы), спроектированные для выполнения в последовательном режиме, будут выглядеть подобно следующей процедуре:

```
Create procedure process_data
As
Begin
For x in ( select * from some_table )
    Выполнить сложную обработку для x
    Обновить какую-то другую таблицу или вставить запись куда-то в другое место
End loop
end
```

В этом случае средство параллельных запросов Oracle или PDML особо не поможет (на самом деле параллельное выполнение здесь SQL-кода базой данных Oracle, вероятно, вызвало бы потребление большего числа ресурсов и заняло бы более длительное время). Если Oracle попытается выполнить простой оператор `SELECT * FROM SOME_TABLE` параллельно, то заметного роста скорости наблюдаться не будет. Если Oracle выполнит параллельно оператор `UPDATE` или `INSERT` после сложной обработки, это также не даст положительного эффекта (в конце концов, это однострочные операторы `UPDATE/INSERT`).

Существует одно очевидное решение, которое можно было бы предпринять здесь: использование массивов для оператора `UPDATE/INSERT` после сложной обработки. Однако это не даст сокращения времени выполнения на 50 и более процентов, к чему вы обычно стремитесь. Не поймите меня неправильно, вы определенно захотите реализовать здесь обработку массива модификаций, но это не сделает процесс быстрее в два, три, четыре или более раз.

Теперь предположим, что данный процесс запускается ночью на машине с четырьмя процессорами и является единственным действием, которое происходит в это время. Вы видите, что в системе частично занят только один процессор, а дисковая система вообще особо не эксплуатируется. Более того, этот процесс занимает часы, причем с каждым днем, по мере добавления новых данных, длится немного дольше. Вам нужно существенно сократить время выполнения этой задачи, ускорив ее в четыре или в восемь раз, так что постепенного ускорения на несколько процентов будет недостаточно. Что можно предпринять?

Здесь применимы два подхода. Один заключается в реализации параллельной конвейерной функции, посредством которой Oracle самостоятельно определит нужные степени параллелизма (предполагая, что вы их выберете, т.к. это рекомендуется).

База данных Oracle будет создавать сеансы, координировать их и запускать. Это очень похоже на предыдущий пример с параллельным DDL или DML, где за счет использования операторов CREATE TABLE AS SELECT или INSERT /*+ APPEND */ база данных Oracle полностью автоматизировала прямую загрузку. Другой подход предусматривает применение самодельного (DIY) параллелизма. В последующих разделах рассматриваются оба подхода.

Параллельные конвейерные функции

Возьмем упомянутый ранее чрезвычайно последовательный процесс PROCESS_DATA и заставим Oracle выполнить его в параллельном режиме. Чтобы достичь этого, мы должны превратить его в подпрограмму. Вместо выбора строк из одной таблицы, их обработки и вставки в другую таблицу мы будем вставлять в другую таблицу результаты извлечения и обработки ряда строк. Мы удалим оператор INSERT из нижней части цикла и заменим его кодом с конструкцией PIPE ROW. Конструкция PIPE ROW позволит нашей подпрограмме PL/SQL генерировать табличные данные в качестве своего вывода, так что появится возможность производить выборку из процесса PL/SQL. Подпрограмма PL/SQL, используемая для процедурной обработки данных, по существу становится *таблицей*, а извлекаемые и обрабатываемые строки — выводом. Вы уже неоднократно видели это на протяжении настоящей книги, когда выдавался оператор следующего вида:

```
Select * from table(dbms_xplan.display);
```

Это подпрограмма PL/SQL, которая читает таблицу PLAN_TABLE, реструктурирует вывод вплоть до добавления строк и выводит эти данные с применением конструкции PIPE ROW для отправки их клиенту. Мы собираемся сделать здесь в сущности то же самое, но позволим выполнять обработку параллельно.

В примере участвуют две таблицы: T1 и T2. Из них T1 — таблица, которую мы читали ранее в строке select * from some_table процедуры PROCESS_DATA, а T2 — таблица, в которую необходимо перенести эту информацию. Предположим, что это определенный вид процесса ETL, который мы запускаем для получения данных о транзакциях за текущий день и преобразования их в отчетную информацию на завтра. Ниже показаны используемые таблицы:

```
EODA@ORA12CR1> create table t1
2 as
3 select object_id id, object_name text
4   from all_objects;
Table created.
Таблица создана.

EODA@ORA12CR1> begin
2   dbms_stats.set_table_stats
3   ( user, 'T1', numrows=>10000000,numblks=>100000 );
4 end;
5 /
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.

EODA@ORA12CR1> create table t2
2 as
```

```

3 select t1.*, 0 session_id
4   from t1
5  where 1=0;
Table created.
Таблица создана.

```

Мы применяем пакет DBMS_STATS для того, чтобы заставить оптимизатор считать, что входная таблица содержит 10 000 000 записей и потребляет 100 000 блоков. То есть мы хотим эмулировать крупную таблицу. Вторая таблица, T2, представляет собой копию структуры первой таблицы, но с дополнительным столбцом SESSION_ID. Этот столбец будет полезен, чтобы увидеть параллелизм в действии.

Далее нам необходимо настроить типы объектов, которые будут возвращать конвейерная функция. Тип объекта — это структурное определение вывода процедуры. В рассматриваемом случае он просто подобен таблице T2:

```

EODA@ORA12CR1> CREATE OR REPLACE TYPE t2_type
2 AS OBJECT (
3   id          number,
4   text        varchar2(30),
5   session_id  number
6 )
7 /
Type created.
Тип создан.

EODA@ORA12CR1> create or replace type t2_tab_type
2 as table of t2_type
3 /
Type created.
Тип создан.

```

Обратимся к конвейерной функции, которая является просто переписанной исходной процедурой PROCESS_DATA. Теперь это функция, производящая строки. Она принимает в REF CURSOR входные данные, подлежащие обработке. Функция возвращает только что созданный тип T2_TAB_TYPE. Это конвейерная функция с опцией PARALLEL_ENABLE. С помощью конструкции partition мы сообщаем Oracle о необходимости секционирования или разделения данных любым способом, который работает лучше других; мы не делаем никаких предположений относительно упорядочения данных.

Вы можете также использовать хеш-секционирование или секционирование по диапазонам на определенном столбце в REF CURSOR. Это предполагает применение строго типизированного REF CURSOR, так что компилятор знает, какие столбцы доступны. Хеш-секционирование будет просто отправлять равное количество строк каждому серверу параллельного выполнения для обработки, основываясь на хеш-значении указанного столбца. Секционирование по диапазонам будет отправлять неперекрывающиеся диапазоны данных каждому серверу параллельного выполнения на основе ключа секционирования. Например, если используется секционирование по диапазонам на столбце ID, то каждый сервер параллельного выполнения может получить диапазоны 1...1000, 1001...20000, 20001...30000 и т.д. (значения ID в этих диапазонах).

Здесь мы хотим всего лишь разделить данные. Как именно они будут разделены — для нашей обработки не важно, поэтому определение выглядит так:

```
EODA@ORA12CR1> create or replace
 2 function parallel_pipelined( l_cursor in sys_refcursor )
 3 return t2_tab_type
 4 pipelined
 5 parallel_enable ( partition l_cursor by any )
```

Нас интересует возможность выяснения, какие строки обрабатывались тем или иным сервером параллельного выполнения, поэтому объявим локальную переменную L_SESSION_ID и инициализируем ее из V\$MYSTAT:

```
6
7 is
8     l_session_id number;
9     l_rec         t1%rowtype;
10 begin
11     select sid into l_session_id
12         from v$mystat
13         where rownum =1;
```

Теперь можно приступать к обработке данных. Мы просто извлекаем строку (или строки, т.к. определенно могли бы применять BULK COLLECT для обработки как массива данных REF CURSOR), подвергаем ее сложной обработке и передаем результат в конвейер. Когда данные в REF CURSOR будут исчерпаны, мы закрываем курсор и возвращаем управление:

```
14     loop
15         fetch l_cursor into l_rec;
16         exit when l_cursor%notfound;
17         -- complex process here
18         pipe row(t2_type(l_rec.id,l_rec.text,l_session_id));
19     end loop;
20     close l_cursor;
21     return;
22 end;
23 /
```

Function created.

Функция создана.

Вот и все. Мы готовы обрабатывать данные параллельно, позволив Oracle самостоятельно вывести самую подходящую степень параллелизма на основе доступных ресурсов:

```
EODA@ORA12CR1> alter session enable parallel dml;
Session altered.
Сеанс изменен.

EODA@ORA12CR1> insert /*+ append */
 2 into t2(id,text,session_id)
 3 select *
 4 from table(parallel_pipelined
 5             (CURSOR(select /*+ parallel(t1) */ *
 6                       from t1 )
 7             ))
```

```

8 /
17914 rows created.
17914 строк создано.
EODA@ORA12CR1> commit;
Commit complete.
Фиксация выполнена.

```

Просто чтобы понаблюдать за происходящим, мы можем запросить вновь вставленные данные, сгруппировать их по SESSION_ID и увидеть, сколько использовалось серверов параллельного выполнения, и сколько строк обработал каждый из них:

```

EODA@ORA12CR1> select session_id, count(*)
2   from t2
3   group by session_id;

```

| SESSION_ID | COUNT(*) |
|------------|----------|
| 198 | 2166 |
| 11 | 2569 |
| 13 | 2493 |
| 185 | 1865 |
| 95 | 2613 |
| 17 | 2377 |
| 256 | 2331 |
| 103 | 1500 |

```

8 rows selected.
8 строк выбрано.

```

Очевидно, что для компонента SELECT этой параллельной операции были задействованы 8 серверов параллельного выполнения, каждый из которых в среднем обработал около 2000 записей.

Как видите, база данных Oracle распараллелила наш процесс, но он претерпел существенную переработку. Это длинный путь, который должен быть пройден от исходной реализации. Таким образом, хотя Oracle может выполнить нашу процедуру параллельно, процедур, закодированных для параллельного выполнения, вполне может не оказаться. Если довольно большой объем работы по переработке процедуры нежелателен, вас может заинтересовать реализация самодельного (DIY) параллелизма.

Самодельный параллелизм

Пусть имеется тот же самый процесс, что и в предыдущем разделе: простая последовательная процедура. Мы не можем позволить себе переписывать ее реализацию, но хотели бы выполнять ее параллельном режиме. Что можно предпринять?

В Oracle 11g Release 2 и последующих версиях мы располагаем новым способом реализации параллелизма — посредством встроенного пакета DBMS_PARALLEL_EXECUTE. С его помощью можно выполнять оператор SQL или PL/SQL параллельно, взяв подлежащие обработке данные и разделив их на множество потоков меньшего размера. Этот подход можно было бы реализовать вручную, как будет показано в следующем разделе “Самодельный параллелизм старой школы”. Однако преимущество нового пакета в том, что он избавляет от большинства утомительной рутинной работы, которую в противном случае пришлось бы взять на себя.

Давайте начнем с допущения, что имеется последовательная процедура SERIAL, которую желательно было бы выполнять параллельно в отношении крупной таблицы. Мы хотели бы сделать это с как можно меньшими затратами труда — другими словами, изменить минимум кода и написать очень мало нового кода. Обратимся к пакету DBMS_PARALLEL_EXECUTE. Мы не будем здесь раскрывать все возможные применения этого пакета (он полностью документирован в справочнике по пакетам и типам PL/SQL (*Oracle Database PL/SQL Packages and Types Reference*)), но используем достаточную его часть для реализации только что описанного процесса.

Предполагая, что мы начинаем с пустой таблицы T2, модифицируем последовательный процесс, чтобы он выглядел так, как показано ниже (дополнения в исходному простому последовательному процессу выделены полужирным):

```
EODA@ORA12CR1> create or replace
 2 procedure serial( p_lo_rowid in rowid, p_hi_rowid in rowid )
 3 is
 4 begin
 5     for x in ( select object_id id, object_name text
 6               from big_table
 7               where rowid between p_lo_rowid
 8                           and p_hi_rowid )
 9     loop
10         -- сложная обработка данных
11         insert into t2 (id, text, session_id )
12         values ( x.id, x.text, sys_context( 'userenv', 'sessionid' ) );
13     end loop;
14 end;
15 /
Procedure created.
```

Это все: нужно просто добавить входные значения ROWID и предикат. Код по большому счету не изменился. Я использую процедуру SYS_CONTEXT для получения идентификатора SESSIONID, чтобы можно было отслеживать, сколько работы было выполнено каждым потоком и каждым параллельным сеансом.

Перед запуском процесса понадобится разделить таблицу на небольшие части. Это можно сделать по некоторому числовому диапазону (удобно для таблиц, применяющих последовательность для заполнения первичного ключа), по любому произвольному SQL-коду, который вы захотите подготовить, или по диапазонам ROWID. Мы будем использовать диапазон ROWID. Я считаю это наиболее эффективным способом, т.к. он создает неперекрывающиеся диапазоны таблицы (свободные от конкуренции) и не требует запрашивания таблицы для определения диапазонов; применяется только словарь данных. Итак, выполним следующие API-вызовы:

```
EODA@ORA12CR1> begin
 2     dbms_parallel_execute.create_task('PROCESS BIG TABLE');
 3     dbms_parallel_execute.create_chunks_by_rowid
 4     ( task_name => 'PROCESS BIG TABLE',
 5       table_owner => user,
 6       table_name => 'BIG_TABLE',
 7       by_row      => false,
 8       chunk_size => 10000 );
 9 end;
10 /
PL/SQL procedure successfully completed.
```

Мы начали с создания именованной задачи, в данном случае 'PROCESS BIG TABLE'. Это просто уникальное имя, которое служит для ссылки на наш большой процесс. Затем мы вызвали процедуру CREATE_CHUNKS_BY_ROWID. Эта процедура разбивает таблицу на диапазоны по ROWID в манере, подобной тому, как мы только что делали. Мы сообщаем процедуре о том, что необходимо прочитать информацию о таблице BIG_TABLE текущего пользователя и разбить ее на фрагменты размером не более 10 000 блоков (CHUNK_SIZE). Параметр BY_ROW был установлен в false, подразумевая в этом случае то, что CHUNK_SIZE указывает не количество строк для создания диапазонов ROWID, а количество блоков для их создания.

Просмотреть число фрагментов и информацию о каждом из них можно сразу после выполнения этого блока кода, запросив новое представление DBA_PARALLEL_EXECUTE_CHUNKS:

```
EODA@ORA12CR1> select *
2   from (
3   select chunk_id, status, start_rowid, end_rowid
4   from dba_parallel_execute_chunks
5   where task_name = 'PROCESS BIG TABLE'
6   order by chunk_id
7   )
8   where rownum <= 5
9   /
```

| CHUNK_ID | STATUS | START_ROWID | END_ROWID |
|----------|------------|--------------------|---------------------|
| 1 | UNASSIGNED | AAAEyWAAEAAAAQoAAA | AAAEyWAAEAAAAQvCcP |
| 2 | UNASSIGNED | AAAEyWAAEAAAAQwAAA | AAAEyWAAEAAAAQ3CcP |
| 3 | UNASSIGNED | AAAEyWAAEAAAAQ4AAA | AAAEyWAAEAAAAQ/CcP |
| 4 | UNASSIGNED | AAAEyWAAEAAAAARAAA | AAAEyWAAEAAAAARHCcP |
| 5 | UNASSIGNED | AAAEyWAAEAAAAARIAA | AAAEyWAAEAAAAARPCcP |

Запрос в этом примере отображает первые пять строк в представлении; в моем случае всего было 218 итоговых строк для интересующей таблицы, каждая из которых представляла неперекрывающуюся часть таблицы для обработки. Это вовсе не означает, что таблица будет обрабатываться в 218 параллельных потоках, просто мы имеем суммарно 218 фрагментов, подлежащих обработке. Теперь мы готовы запустить задачу посредством следующего API-вызова:

```
EODA@ORA12CR1> begin
2   dbms_parallel_execute.run_task
3   ( task_name      => 'PROCESS BIG TABLE',
4     sql_stmt       => 'begin serial( :start_id, :end_id ); end;',
5     language_flag  => DBMS_SQL.NATIVE,
6     parallel_level => 4 );
7 end;
8 /
```

PL/SQL procedure successfully completed.

Процедура PL/SQL успешно завершена.

Здесь мы запрашиваем запуск задачи 'PROCESS BIG TABLE', указывающей на наши фрагменты. Оператор SQL, который необходимо выполнить, выглядит как 'begin serial(:start_id, :end_id); end;' — простой вызов нашей хранимой процедуры с диапазоном ROWID, предназначенным для обработки.

Параметр `PARALLEL_LEVEL` получил значение 4, т.е. при выполнении оператора будут задействованы четыре параллельных потока/процесса. Хотя имеются 218 фрагментов, мы будем обрабатывать по четыре фрагмента за раз. Для параллельного выполнения этих потоков пакет `DBMS_PARALLEL_EXECUTE` внутренне использует пакет `DBMS_SCHEDULER`.

Как только задача начинает функционировать, она создает четыре задания; каждому заданию сообщается о том, что нужно обработать фрагменты, идентифицированные значением ключа 'PROCESS BIG TABLE', и запустить для каждого фрагмента хранимую процедуру `SERIAL`. Таким образом, эти четыре сеанса запускаются, и каждый из них читает фрагмент из представления `DBA_PARALLEL_EXECUTE_CHUNKS`, обрабатывает его и обновляет столбец `STATUS`. Если фрагмент обработан успешно, то строка помечается как `PROCESSED`. Если по какой-то причине произошел отказ или данный фрагмент не может быть обработан, то строка помечается как `PROCESSED_WITH_ERROR`, а другие столбцы будут содержать подробное сообщение об ошибке, описывающее причину ее возникновения. В любом случае сеанс затем извлекает следующий фрагмент и обрабатывает его и т.д. Со временем эти четыре задания обработают все фрагменты, и вся задача будет завершена.

В случае если обработка любого из фрагментов терпит неудачу, можно устранить причину ошибки и возобновить выполнение задачи. Это приведет к повторной обработке отказавших фрагментов. Когда все они будут обработаны успешно, работа сделана.

```
EODA@ORA12CR1> select *
2   from (
3   select chunk_id, status, start_rowid, end_rowid
4   from dba_parallel_execute_chunks
5   where task_name = 'PROCESS BIG TABLE'
6   order by chunk_id
7   )
8   where rownum <= 5
9   /
```

| CHUNK_ID | STATUS | START_ROWID | END_ROWID |
|----------|-----------|--------------------|---------------------|
| 1 | PROCESSED | AAAEyWAAEAAAAQoAAA | AAAEyWAAEAAAAQvCcP |
| 2 | PROCESSED | AAAEyWAAEAAAAQwAAA | AAAEyWAAEAAAAQ3CcP |
| 3 | PROCESSED | AAAEyWAAEAAAAQ4AAA | AAAEyWAAEAAAAQ/CcP |
| 4 | PROCESSED | AAAEyWAAEAAAAARAAA | AAAEyWAAEAAAAARHCcP |
| 5 | PROCESSED | AAAEyWAAEAAAAARIAA | AAAEyWAAEAAAAARPCcP |

Задачу можно либо сохранить для истории, либо удалить. Ниже показано, как удалить задачу:

```
EODA@ORA12CR1> begin
2   dbms_parallel_execute.drop_task('PROCESS BIG TABLE' );
3 end;
4   /
```

PL/SQL procedure successfully completed.

Процедура PL/SQL успешно завершена.

Если мы посмотрим собственную таблицу приложения, то увидим, что задача выполнялась с применением четырех параллельных потоков, каждый из которых обработал примерно одинаковое количество строк:

```

EODA@ORA12CR1> select session_id, count(*)
2   from t2
3   group by session_id
4   order by session_id;

```

```

SESSION_ID  COUNT(*)
-----
22603       2521812
22604       2485273
22605       2529386
22606       2463529

```

4 rows selected.

4 строки выбрано.

Если вы пока еще не работаете в среде Oracle 11g Release 2 или последующей версии, то можете реализовать такой же вид параллелизма с использованием более трудоемкого подхода, описанного в следующем разделе. Однако новый пакет предлагает развитый API-интерфейс (который мы только слегка затронули здесь), включающий гораздо больше функциональности, чем можно обеспечить при ручной реализации.

Самодельный параллелизм старой школы

В версиях, предшествующих Oracle 11g Release 2, также можно реализовать подход к распараллеливанию, похожий на тот, что был описан в предыдущем разделе. Здесь нет развитого API-интерфейса, который бы оказывал поддержку, поэтому много рутинной работы придется проделать самостоятельно, но все же это возможно. Я много раз применял подход на основе диапазонов ROWID для разбиения таблицы на ряд диапазонов, которые не перекрываются (хотя полностью охватывают таблицу).

На заметку! Если вы работаете в среде Oracle 11g Release 2 или последующей версии, обратитесь к предыдущему разделу за примером применения пакета DBMS_PARALLEL_EXECUTE. Если этот пакет доступен, то вы должны использовать его вместо описанного здесь ручного подхода.

Этот трудоемкий подход концептуально очень похож на то, как Oracle выполняет параллельный запрос. Скажем, при обработке полного сканирования таблицы Oracle с помощью определенного метода разбивает таблицу на множество таблиц меньших размеров, каждая из которых обрабатывается сервером параллельного выполнения. Мы собираемся сделать то же самое с применением диапазонов ROWID. В ранних выпусках сама реализация параллелизма Oracle в действительности использовала диапазоны ROWID.

Мы будем применять таблицу BIG_TABLE с одним миллионом строк, т.к. описываемый прием работает лучше всего на больших таблицах с множеством экстенгов, а метод, используемый для создания диапазонов ROWID, зависит от границ экстенгов. Чем больше имеется экстенгов, тем лучше распределены данные. Итак, после создания таблицы BIG_TABLE, содержащей один миллион строк, мы создадим таблицу T2 следующим образом:


```

EODA@ORA12CR1> create table t2
2 as
3 select object_id id, object_name text, 0 session_id
4   from big_table
5  where 1=0;
Table created.
Таблица создана.

```

Для параллельной обработки нашей процедуры мы будем применять поддержку очередей заданий, встроенную в базу данных. Мы запланируем запуск некоторого числа заданий. Каждое задание — это наша процедура, слегка измененная для обработки только строк в заданном диапазоне ROWID.

На заметку! В Oracle 10g и последующих версиях для простоты можно также использовать планировщик. Здесь мы будем применять очереди заданий, чтобы сделать пример совместимым с версией Oracle9i.

Для эффективной поддержки очередей заданий будет использоваться таблица параметров, обеспечивающая передачу исходных значений нашим заданиям:

```

EODA@ORA12CR1> create table job_parms
2 ( job      number primary key,
3   lo_rid   rowid,
4   hi_rid   rowid
5 )
6 /
Table created.
Таблица создана.

```

Это позволит просто передать идентификатор задания процедуре SERIAL, так что она сможет запросить данную таблицу и получить диапазон ROWID для дальнейшей обработки. Теперь модифицируем нашу процедуру (вновь добавленный код выделен полужирным):

```

EODA@ORA12CR1> create or replace
2 procedure serial( p_job in number )
3 is
4   l_rec      job_parms%rowtype;
5 begin
6   select * into l_rec
7     from job_parms
8    where job = p_job;
9
10  for x in ( select object_id id, object_name text
11            from big_table
12            where rowid between l_rec.lo_rid
13                        and l_rec.hi_rid )
14  loop
15    -- сложная обработка данных
16    insert into t2 (id, text, session_id )
17      values ( x.id, x.text, p_job );
18  end loop;
19

```

```

20      delete from job_parms where job = p_job;
21      commit;
22 end;
23 /
Procedure created.
Процедура создана.

```

Как видите, изменения незначительны. Большая часть добавленного кода просто принимает входные данные и диапазон ROWID для обработки. Единственное изменение в логике связано с модификацией предиката в строках 12 и 13.

Теперь давайте запланируем запуск задания. Мы построим довольно сложный запрос, который применяет для разделения таблицы аналитические средства. Самый внутренний запрос в строках 19–26 разбивает данные на восемь групп. Первый вызов sum в строке 22 вычисляет промежуточный итог суммы блоков, а второй вызов sum в строке 23 — общее количество блоков. Если выполнить целочисленное деление промежуточного итога на желаемый размер фрагмента (здесь это общий размер, деленный на 8), мы сможем создать группы файлов/блоков, которые охватывают почти одинаковые объемы данных. Запрос в строках 8–28 находит максимальный и минимальный номера файлов и номера блоков по GRP и возвращает отдельные записи. Он строит входные данные, которые затем будут отправлены пакету DBMS_ROWID для создания идентификаторов ROWID, необходимых Oracle. Мы возьмем этот вывод и с помощью пакета DBMS_JOB отправим задание для обработки диапазона ROWID:

```

EODA@ORA12CR1> declare
2      l_job number;
3  begin
4  for x in (
5  select dbms_rowid.rowid_create
         ( 1, data_object_id, lo_fno, lo_block, 0 ) min_rid,
6         dbms_rowid.rowid_create
         ( 1, data_object_id, hi_fno, hi_block, 10000 ) max_rid
7  from (
8  select distinct grp,
9         first_value(relative_fno)
         over (partition by grp order by relative_fno, block_id
10             rows between unbounded preceding and unbounded following) lo_fno,
11         first_value(block_id)
         over (partition by grp order by relative_fno, block_id
12             rows between unbounded preceding and unbounded following) lo_block,
13         last_value(relative_fno)
         over (partition by grp order by relative_fno, block_id
14             rows between unbounded preceding and unbounded following) hi_fno,
15         last_value(block_id+blocks-1)
         over (partition by grp order by relative_fno, block_id
16             rows between unbounded preceding and unbounded following) hi_block,
17         sum(blocks) over (partition by grp) sum_blocks
18  from (
19  select relative_fno,
20         block_id,
21         blocks,
22         trunc( (sum(blocks) over (order by relative_fno, block_id)-0.01) /
23               (sum(blocks) over ()/8) ) grp

```

```

24  from dba_extents
25  where segment_name = upper('BIG_TABLE')
26  and owner = user order by block_id
27  )
28  ),
29  (select data_object_id
    from user_objects where object_name = upper('BIG_TABLE') )
30 )
31 loop
32     dbms_job.submit( l_job, 'serial(JOB);' );
33     insert into job_parms(job, lo_rid, hi_rid)
34     values ( l_job, x.min_rid, x.max_rid );
35 end loop;
36 end;
37 /

```

PL/SQL procedure successfully completed.

Процедура PL/SQL успешно завершена.

Этот блок PL/SQL способен запланировать запуск вплоть до восьми заданий (или меньшего их числа, если таблица не может быть разбита на восемь частей из-за недостаточного количества экстенстов либо размера). Просмотреть количество запланированных к запуску заданий и их входные данные можно следующим образом:

```
EODA@ORA12CR1> select * from job_parms;
```

| | JOB | LO_RID | HI_RID |
|---|--------------------|---------------------|--------|
| 1 | AAAEzwAAEAAABKAAAA | AAAEzwAAEAAABl/CcQ | |
| 2 | AAAEzwAAEAAACyAAAA | AAAEzwAAEAAADR/CcQ | |
| 3 | AAAEzwAAEAAAauAAAA | AAAEzwAAEAAABJ/CcQ | |
| 4 | AAAEzwAAEAAACCAAAA | AAAEzwAAEAAACR/CcQ | |
| 5 | AAAEzwAAEAAADSAAAA | AAAEzwAAEAAABGUnCcQ | |
| 6 | AAAEzwAAEAAAQoAAA | AAAEzwAAEAAAAt/CcQ | |
| 7 | AAAEzwAAEAAABmAAAA | AAAEzwAAEAAACB/CcQ | |
| 8 | AAAEzwAAEAAACSAAAA | AAAEzwAAEAAACx/CcQ | |

8 rows selected.

8 строк выбрано.

```
EODA@ORA12CR1> commit;
```

Commit complete.

Фиксация выполнена.

Произведенная фиксация освобождает задания для обработки. В файле параметров для JOB_QUEUE_PROCESSES указано значение 1000, так что все восемь заданий запустятся и быстро завершатся. После их завершения результаты выглядят так:

```

EODA@ORA12CR1> select session_id, count(*)
2  from t2
3  group by session_id;

```

| SESSION_ID | COUNT(*) |
|------------|----------|
| 1 | 127651 |
| 6 | 124431 |
| 2 | 147606 |
| 5 | 124590 |

```

4      72961
8      147544
3      127621
7      127596

```

8 rows selected.

8 строк выбрано.

Однако предположим, что мы не хотим использовать обработку по ROWID — возможно, запрос не настолько прост, как `SELECT * FROM T`, и содержит соединения и другие конструкции, которые делают применение ROWID непрактичным. Вместо этого можно использовать первичный ключ какой-то таблицы. Например, пусть необходимо разбить ту же самую таблицу `BIG_TABLE` на десять частей, подлежащих параллельной обработке, по первичному ключу. Задачу легко решить с помощью встроенной аналитической функции `NTILE`. Процесс достаточно прост:

```

EODA@ORA12CR1> select nt, min(id), max(id), count(*)
2   from (
3   select id, ntile(10) over (order by id) nt
4   from big_table
5   )
6   group by nt;

```

| NT | MIN(ID) | MAX(ID) | COUNT (*) |
|----|---------|---------|-----------|
| 1 | 1 | 100000 | 100000 |
| 6 | 500001 | 600000 | 100000 |
| 2 | 100001 | 200000 | 100000 |
| 5 | 400001 | 500000 | 100000 |
| 4 | 300001 | 400000 | 100000 |
| 8 | 700001 | 800000 | 100000 |
| 3 | 200001 | 300000 | 100000 |
| 7 | 600001 | 700000 | 100000 |
| 9 | 800001 | 900000 | 100000 |
| 10 | 900001 | 1000000 | 100000 |

10 rows selected.

10 строк выбрано.

Теперь в вашем распоряжении десять неперекрывающихся диапазонов первичного ключа (все имеют одинаковые размеры), которые можно применять для реализации описанного выше приема с пакетом `DBMS_JOB`, чтобы распараллелить процесс.

Резюме

В этой главе мы исследовали концепцию параллельного выполнения в Oracle. Сначала была представлена аналогия, упрощающая понимание того, когда и где следует использовать параллельное выполнение, а именно — при наличии длительно функционирующих операторов или процедур и массы доступных ресурсов.

Затем мы посмотрели, как Oracle может задействовать параллелизм. Мы начали с параллельных запросов и способа разбиения базой данных Oracle крупных последовательных операций вроде полного сканирования на меньшие части, которые могут выполняться одновременно. Затем мы перешли к параллельному DML (PDML) и раскрыли довольно обширный список присущих ему ограничений.

После этого мы взглянули на наилучшую точку приложения параллельных операций: параллельный DDL. Параллельный DDL — это инструмент, предназначенный для администраторов баз данных и также разработчиков, который позволяет быстро выполнять операции по обслуживанию, обычно производимые в период минимальной нагрузки, когда ресурсы свободны. Мы рассмотрели процедурный параллелизм и два приема распараллеливания процедур: применяемого самой базой данных Oracle и самодельного.

При проектировании процесса с нуля можно принимать во внимание возможность его распараллеливания базой данных Oracle, чтобы будущее добавление или сокращение ресурсов позволяло легко варьировать степень параллелизма. Однако если есть готовый код, который необходимо быстро скорректировать, чтобы сделать его параллельным, можно прибегнуть к самодельному (DIY) параллелизму, реализуемому с помощью двух приемов — ручного и автоматического. В каждом из них используются диапазоны ROWID или диапазоны первичного ключа, а также пакет DBMS_JOB или DBMS_SCHEDULER для параллельного запуска заданий в фоновом режиме.

глава 15

Загрузка и выгрузка данных

В этой главе речь пойдет о загрузке и выгрузке данных, т.е. о том, как помещать и извлекать данные из базы данных Oracle. Основное внимание будет уделено следующим инструментам пакетной загрузки данных.

- **Внешние таблицы.** Это средство из Oracle9i и последующих версий, которое разрешает доступ к файлам операционной системы, как если бы они были таблицами базы данных, а в Oracle 10g и последующих версиях оно даже позволяет создавать файлы операционной системы как фрагменты таблиц.
- **SQL*Loader (SQLLDR).** Это исторически установившийся загрузчик данных Oracle, который по-прежнему является распространенным методом для загрузки данных.

В области выгрузки данных мы рассмотрим перечисленные ниже приемы.

- **Выгрузка Data Pump.** Data Pump — это патентованный Oracle двоичный формат, доступный через инструмент Data Pump и внешние таблицы.
- **Выгрузка в плоские файлы.** Инструменты для выгрузки в плоские файлы будут специально разработанными реализациями, но они предлагают результаты, которые являются переносимыми в системы других типов (даже электронные таблицы).

Внешние таблицы

Внешние таблицы были впервые представлены в версии Oracle9i Release 1. Попросту говоря, они позволяют трактовать файл операционной системы, как если бы это была доступная только для чтения таблица базы данных. Они не задумывались как замена “подлинной” таблицы или в целях использования вместо нее, а предназначены для применения в качестве инструмента, облегчающего загрузку, а в Oracle 10g и последующих версиях — также и выгрузку данных.

Когда средство внешних таблиц только появилось, я часто ссылался на него как на “замену SQLLDR”. Такое утверждение все еще остается справедливым — в *большинстве* ситуаций. Но, учитывая сказанное, может возникнуть вопрос: почему тогда в этой главе приведен материал, посвященный SQLLDR? Причина в том, что инструмент SQLLDR существует на протяжении долгого времени, и за годы было

накоплено огромное количество унаследованных управляющих файлов. Инструмент SQLLDR все еще широко используется; его знают и применяют многие люди. Мы пока находимся на этапе перехода от использования SQLLDR к работе с внешними таблицами, поэтому SQLLDR по-прежнему является очень важным средством.

Однако многие администраторы баз данных не осознают, что их знание управляющих файлов SQLLDR может быть без труда перенесено в область применения внешних файлов. Во время проработки примеров в этой части главы вы обнаружите, что внешние таблицы заключают в себе порядочную часть синтаксиса SQLLDR и многие из его приемов.

С учетом сказанного инструменту SQLLDR должно отдаваться предпочтение *перед* внешними таблицами в следующих ситуациях.

- Требуется загружать данные по сети — другими словами, когда входной файл не находится на самом сервере базы данных. Одно из ограничений внешних таблиц заключается в том, что входной файл должен быть доступен на сервере базы данных.
- Множество пользователей должны параллельно взаимодействовать с той же самой внешней таблицей, обрабатывая разные входные файлы.

Памятуя об этих исключениях, в целом я настоятельно рекомендую использовать внешние таблицы из-за их расширенных возможностей. SQLLDR — довольно простой инструмент, который генерирует оператор INSERT и загружает данные. Его способность работы с SQL ограничивается вызовом SQL-функций на построчной основе. Внешние таблицы открывают доступ ко всему набору функциональности SQL для загрузки данных. Ниже перечислены некоторые ключевые функциональные возможности, благодаря наличию которых внешние таблицы превосходят SQLLDR.

- **Возможность применения сложных условий WHERE для выборочной загрузки данных.** SQLLDR поддерживает конструкцию WHEN для выбора строк, подлежащих загрузке, но вы ограничены использованием только выражений AND и выражений с равенствами — не поддерживаются диапазоны (больше, меньше), выражения OR, конструкции IS NULL и т.д.
- **Способность выполнять слияние (MERGE) данных.** Можно взять файл операционной системы, заполненный данными, и с его помощью обновить записи базы данных.
- **Возможность выполнения эффективного поиска по кодам.** В рамках процесса загрузки можно соединять внешнюю таблицу с другими таблицами базы данных.
- **Способность загрузки отсортированных данных** за счет включения конструкции ORDER BY в оператор CREATE TABLE или INSERT.
- **Упрощение многотабличных вставок с применением оператора INSERT.** Начиная с версии Oracle9i, оператор INSERT может производить вставку в одну или более таблиц, используя сложные условия WHEN. Хотя SQLLDR позволяет выполнять загрузку сразу во множество таблиц, формирование необходимого синтаксиса может оказаться довольно сложным.
- **Возможность указания одной и более команд операционной системы,** которые должны быть выполнены как первый шаг (предварительная обработка) при выборе данных из внешней таблицы.

- **Пологая кривая обучения для новых разработчиков.** SQLLDR — это еще один инструмент для изучения в дополнение к языку программирования, инструментам разработки, языку SQL и т.д. Если разработчик знает SQL, то он может немедленно применить свои навыки в области пакетной загрузки данных без потребности в изучении нового инструмента (SQLLDR).

Таким образом, имея все это в виду, давайте посмотрим, как использовать внешние таблицы.

Настройка внешних таблиц

Существуют два простых метода, чтобы приступить к работе с внешними таблицами:

- выполнение SQLLDR с параметром `EXTERNAL_TABLE`;
- начиная с версии Oracle 12c, запуск SQLLDR в скоростном режиме (express mode).

Интересно отметить, что оба метода задействуют инструмент командной строки SQLLDR. Они обсуждаются в последующих разделах.

Выполнение SQLLDR с параметром `EXTERNAL_TABLE`

Один из самых легких способов начать работу с внешними таблицами предусматривает применение существующего унаследованного управляющего файла, чтобы получить определение внешней таблицы. Для первой простой демонстрации мы воспользуемся управляющим файлом SQLLDR следующего вида (подробные сведения об управляющих файлах SQLLDR будут приведены позже в разделе “Инструмент SQLLDR”):

```
LOAD DATA
INFILE *
INTO TABLE DEPT
FIELDS TERMINATED BY ','
(DEPTNO, DNAME, LOC )
BEGINDATA
10,Sales, Virginia
20,Accounting, Virginia
30,Consulting, Virginia
40,Finance, Virginia
```

Вначале создадим таблицу DEPT:

```
EODA@ORA12CR1> create table dept as select * from scott.dept;
Table created.
Таблица создана.
```

Следующая команда SQLLDR сгенерирует оператор CREATE TABLE для внешней таблицы:

```
$ sqlldr eoda demo1.ctl external_table=generate_only
Password:
SQL*Loader: Release 12.1.0.1.0 - Production on Fri Mar 7 16:28:38 2014
Copyright (c) 1982, 2013, Oracle and/or its affiliates. All rights reserved.
Path used:      External Table
```


Параметр `EXTERNAL_TABLE` может принимать одно из трех значений.

- `NOT_USED`. Стандартное значение с очевидным смыслом — параметр не применяется.
- `EXECUTE`. Означает, что `SQLLDR` не будет генерировать SQL-оператор `INSERT` и выполнять его. Вместо этого он создаст внешнюю таблицу и использует однократный пакетный SQL-оператор для загрузки в нее данных.
- `GENERATE_ONLY`. Заставляет `SQLLDR` не загружать какие-либо данные, а только генерировать операторы DDL и DML, которые он бы выполнял, и помещать их в созданный журнальный файл.

Внимание! Параметр `DIRECT=TRUE` переопределяет `EXTERNAL_TABLE=GENERATE_ONLY`. Если вы укажете `DIRECT=TRUE` в Oracle 10g и предшествующих версиях, то данные загрузятся, но внешняя таблица создаваться не будет. В Oracle 11g Release 1 и последующих версиях вместо этого вы получите сообщение об ошибке `SQL*Loader-144: Conflicting load methods: direct=true/external_table=generate_only specified` (`SQL*Loader-144: конфликтующие методы загрузки: указано direct=true/external_table=generate_only`). Дело вовсе не в том, что вы решите так поступить; просто запомните, что упомянутые два параметра несовместимы.

Когда применяется значение `GENERATE_ONLY`, в файле `demol.log` можно увидеть следующие данные:

```
CREATE DIRECTORY statements needed for files
Операторы CREATE DIRECTORY, необходимые для файлов
-----
CREATE DIRECTORY SYS_SQLLDR_XT_TMPDIR_00000 AS '/home/tkyte'
```

Оператор `CREATE DIRECTORY` в журнальном файле вполне может отсутствовать. Во время генерации сценария внешней таблицы `SQLLDR` подключается к базе данных и запрашивает словарь данных, чтобы выяснить, существует ли подходящий каталог. В этом случае подходящий каталог не обнаружен, поэтому `SQLLDR` сгенерировал оператор `CREATE DIRECTORY`. Затем он генерирует оператор `CREATE TABLE` для создания внешней таблицы:

```
CREATE TABLE statement for external table:
Оператор CREATE TABLE для внешней таблицы:
-----
CREATE TABLE "SYS_SQLLDR_X_EXT_DEPT"
(
  "DEPTNO" NUMBER(2),
  "DNAME" VARCHAR2(14),
  "LOC" VARCHAR2(13)
)
```

Далее `SQLLDR` входит в базу данных — именно так он узнает точные типы данных, которые должны использоваться в этом определении внешней таблицы (например, что `DEPTNO` имеет тип `NUMBER(2)`). Загрузчик извлекает сведения о типах прямо из словаря данных. Следующим идет начало определения внешней таблицы:

```
ORGANIZATION external
(
  TYPE oracle_loader
  DEFAULT DIRECTORY SYS_SQLLDR_XT_TMPDIR_00000
```

Конструкция `ORGANIZATION EXTERNAL` сообщает Oracle о том, что это не “обычная” таблица. Эта конструкция уже встречалась ранее в главе 10 при рассмотрении индекс-таблиц. В настоящее время существуют три типа организации, указываемые в `ORGANIZATION`: `HEAP` для обычной таблицы, `INDEX` для индекс-таблицы и `EXTENAL` для внешней таблицы. После `ORGANIZATION EXTERNAL` начинается представление для Oracle дополнительных сведений о внешней таблице. Тип `ORACLE_LOADER` — один из двух поддерживаемых типов (в Oracle9i это *единственный* поддерживаемый тип). Другим типом является `ORACLE_DATAPUMP` — патентованный формат Data Pump, применяемый в Oracle 10g и последующих версиях. Мы взглянем на этот тип в разделе, посвященном выгрузке данных, далее в главе; он представляет собой формат, который может использоваться как при загрузке, так и при выгрузке данных. Внешняя таблица может применяться для создания файла формата Data Pump и для его последующего чтения.

Далее мы встречаем раздел `ACCESS PARAMETERS` (параметры доступа) внешней таблицы. В нем приводится описание для базы данных того, как она должна обрабатывать входной файл. Глядя на указанный раздел, вы должны заметить сходство с управляющим файлом `SQLldr`; это не случайно. По большей части `SQLldr` и внешние таблицы используют очень похожий синтаксис:

```
ACCESS PARAMETERS
(
  RECORDS DELIMITED BY NEWLINE CHARACTERSET US7ASCII
  BADFILE 'SYS_SQLldr_XT_TMPDIR_00000': 'demo1.bad'
  LOGFILE 'demo1.log_xt'
  READSIZE 1048576
  SKIP 6
  FIELDS TERMINATED BY "," LDRTRIM
  REJECT ROWS WITH ALL NULL FIELDS
  (
    "DEPTNO" CHAR(255)
      TERMINATED BY ",",
    "DNAME" CHAR(255)
      TERMINATED BY ",",
    "LOC" CHAR(255)
      TERMINATED BY ","
  )
)
```

Эти параметры доступа показывают, каким образом настроить внешнюю таблицу, чтобы она обрабатывала файлы почти идентично тому, как их обрабатывал бы загрузчик `SQLldr`.

- **RECORDS.** По умолчанию записи завершаются символами новой строки, как в `SQLldr`.
- **BADFILE.** Файл некорректных записей (файл, куда помещаются записи, которые не удалось обработать), находящийся в только что созданном каталоге.
- **LOGFILE.** Журнальный файл, который эквивалентен журнальному файлу `SQLldr`, созданный в текущем рабочем каталоге.

- **READSIZE.** Стандартный буфер, применяемый Oracle для чтения входного файла данных. В этом случае он составляет 1 Мбайт. Память под него берется из области PGA в режиме выделенного сервера и области SGA в режиме разделяемого сервера. Этот буфер предназначен для буферизации данных входного файла в сеансе (память PGA и SGA обсуждалась в главе 4). Помните об этом факте, если используете разделяемые серверы: память выделяется из области SGA.
- **SKIP 6.** Количество записей, которые должны быть пропущены во входном файле. Может возникнуть вопрос: почему пропускаются 6 записей? В этом примере применяется `INFILE *`, а `SKIP 6` позволяет пропустить начало управляющего файла и перейти непосредственно к встроенным данным. Если бы мы не использовали конструкцию `INFILE *`, то не было бы и `SKIP`.
- **FIELDS TERMINATED BY.** Это точно такая же конструкция, которая применялась в самом управляющем файле. Однако в случае внешней таблицы добавляется `LDRTRIM`, что означает `LoaDeR TRIM`. Опция `LDRTRIM` указывает режим усечения, эмулирующий способ, которым `SQLLDR` по умолчанию усекает данные. К другим опциям относятся `LTRIM`, `LTRIM` и `RTRIM` (для усечения пробельных символов слева/справа), а также `NOTRIM` для сохранения всех ведущих/хвостовых пробельных символов.
- **REJECT ROWS WITH ALL NULL FIELDS.** Это заставляет внешнюю таблицу заносить в файл некорректных записей любые полностью пустые строки и не загружать их.
- **Сами определения столбцов.** Это метаданные, которые описывают ожидаемые входные значения. Все они являются символьными строками в загружаемом файле данных, могут иметь длину до 255 символов (стандартный размер `SQLLDR`), завершаются запятыми и необязательно заключены в кавычки.

На заметку! Полный список опций, доступных при использовании внешних таблиц, можно найти в руководстве по утилитам базы данных Oracle (*Oracle Database Utilities*), которое содержит раздел, полностью посвященный внешним таблицам. В руководстве по языку SQL базы данных Oracle (*Oracle Database SQL Language Reference*) приведен базовый синтаксис, но без подробностей, касающихся раздела `ACCESS PARAMETERS`.

Итак, мы добрались до раздела `LOCATION` (местоположение) определения внешней таблицы:

```
location
(
  'demo1.ctl'
)
) REJECT LIMIT UNLIMITED
```

Здесь Oracle сообщается имя файла для загрузки, `demo1.ctl` в этом случае, т.к. в исходном управляющем файле применялась конструкция `INFILE *`. Следующим оператором в управляющем файле является стандартный `INSERT`, который может использоваться для загрузки в таблицу данных из самой внешней таблицы:

`INSERT statements used to load internal tables:`

Операторы INSERT, используемые для загрузки внутренних таблиц:

```

INSERT /*+ append */ INTO DEPT
(
  DEPTNO,
  DNAME,
  LOC
)
SELECT
  "DEPTNO",
  "DNAME",
  "LOC"
FROM "SYS_SQLLDR_X_EXT_DEPT"

```

Показанный оператор должен выполнять логический эквивалент прямой загрузки, если это возможно (предполагается, что подсказка APPEND может быть удовлетворена; наличие триггеров или ограничений внешнего ключа может предотвратить запуск операции прямого пути).

Наконец, в журнальном файле мы увидим операторы, которые могут применяться для удаления созданных SQLLDR объектов после завершения загрузки:

statements to cleanup objects created by previous statements:
операторы для очистки объектов, созданных предыдущими операторами:

```

-----
DROP TABLE "SYS_SQLLDR_X_EXT_DEPT"
DROP DIRECTORY SYS_SQLLDR_XT_TMPDIR_00000

```

Вот и все. Если мы возьмем этот журнальный файл и добавим в нужном месте /, чтобы сделать его допустимым сценарием SQL*Plus, то должны быть готовы двигаться дальше — или нет, в зависимости от имеющихся разрешений. Например, предполагая, что схема, в которую произведен вход, имеет привилегию CREATE ANY DIRECTORY или доступ по чтению и записи к существующему каталогу, можно было бы наблюдать следующий вывод:

```

EODA@ORA12CR1> INSERT /*+ append */ INTO DEPT
2  (
3   DEPTNO,
4   DNAME,
5   LOC
6  )
7  SELECT
8   "DEPTNO",
9   "DNAME",
10  "LOC"
11 FROM "SYS_SQLLDR_X_EXT_DEPT"
12 /
INSERT /*+ append */ INTO DEPT
*

```

ERROR at line 1:

ORA-29913: error in executing ODCIEXTTABLEOPEN callout

ORA-29400: data cartridge error

KUP-04063: unable to open log file demo1.log_xt

OS error Permission denied

ORA-06512: at "SYS.ORACLE_LOADER", line 19

ORA-06512: at line 1

ОШИБКА в строке 1:

ORA-29913: ошибка при выполнении внешнего вызова ODCIEXTTABLEOPEN
 ORA-29400: ошибка картриджа данных
 KUP-04063: не удается открыть журнальный файл demo1.log_xt
 Ошибка операционной системы Отказ в доступе
 ORA-06512: в SYS.ORACLE_LOADER, строка 19
 ORA-06512: в строке 1

Поначалу это не выглядит правильным. Я вошел в систему как TKYTE, моим каталогом является /home/tkyte и я его владелец, так что я определенно могу осуществлять в него запись (в конце концов, я же создал в нем журнальный файл SQLLDR). Что случилось? Дело в том, что код внешней таблицы выполняется серверным программным обеспечением Oracle в выделенном или разделяемом сервере. Процесс, пытающийся читать файл входных данных, принадлежит владельцу программного обеспечения Oracle, а не моей учетной записи. Процесс, пытающийся создать журнальный файл, также принадлежит владельцу программного обеспечения Oracle, а не моей учетной записи. Очевидно, что Oracle не располагает привилегией, необходимой для записи в мой каталог, поэтому попытка доступа внешней таблицы потерпела неудачу. Это важный момент. Чтобы читать таблицу, учетная запись, от имени которой функционирует база данных (владелец программного обеспечения Oracle), должна обладать следующими возможностями.

- **Читать файл, на который мы указываем.** В среде UNIX/Linux это означает, что владелец программного обеспечения Oracle должен иметь разрешения чтения и выполнения для всех каталогов пути, ведущем к файлу. В среде Windows владелец программного обеспечения Oracle должен иметь возможность читать этот файл.
- **Записывать в каталоги, куда будет сохраняться журнальный файл** (или полностью пропустить генерацию журнального файла, но это обычно не рекомендуется). В действительности, если журнальный файл уже существует, то владелец программного обеспечения Oracle должен иметь возможность записывать в этот существующий файл.
- **Записывать в любой из указанных файлов некорректных записей** точно так же, как в журнальный файл.

Давайте вернемся к примеру. Следующая команда предоставляет Oracle возможность записи в мой каталог:

```
EODA@ORA12CR1> host chmod a+rw .
```

Внимание! На самом деле эта команда предоставляет возможность записи в каталог абсолютно всем! Она приведена здесь только в целях демонстрации; обычно для этого мы будем использовать специальный каталог, возможно, принадлежащий владельцу программного обеспечения Oracle.

Затем снова выполним предыдущий оператор INSERT:

```
EODA@ORA12CR1> list
1 INSERT /*+ append */ INTO DEPT
2 (
3   DEPTNO,
4   DNAME,
```

```

5  LOC
6  )
7  SELECT
8  "DEPTNO",
9  "DNAME",
10 "LOC"
11* FROM "SYS_SQLLDR_X_EXT_DEPT"
EODA@ORA12CR1> /
4 rows created.
4 строки создано.

EODA@ORA12CR1> host ls -l demo1.log_xt
-rw-r----- 1 oracle dba 687 Mar  8 14:35 demo1.log_xt

```

Как видите, на этот раз удалось обратиться к файлу, успешно загрузить четыре строки и создать журнальный файл, владельцем которого фактически является oracle, а не применяемая учетная запись операционной системы.

Запуск SQLLDR в скоростном режиме

Начиная с версии Oracle 12c, скоростной режим SQLLDR позволяет быстро загружать данные из файла CSV в таблицу. Если используемая вами схема имеет привилегию CREATE ANY DIRECTORY, то скоростной режим попытается применить для загрузки данных внешнюю таблицу. В противном случае он будет использовать SQLLDR в прямом режиме.

Проиллюстрируем сказанное с помощью простого примера. Пусть имеется таблица, созданная следующим образом:

```

EODA@ORA12CR1> create table dept
2  ( deptno  number(2) constraint dept_pk primary key,
3  dname    varchar2(14),
4  loc      varchar2(13)
5  )
6  /

```

Кроме того, в файле CSV по имени dept.dat находятся показанные ниже данные:

```

10,Sales,Virginia
20,Accounting,Virginia
30,Consulting,Virginia
40,Finance,Virginia

```

Для начала мы проделаем этот пример в ситуации, когда пользователь *не* имеет привилегии CREATE ANY DIRECTORY. Скоростной режим SQLLDR иницируется посредством ключевого слова TABLE:

```

$ sqlldr eoda table=dept
Password:
SQL*Loader: Release 12.1.0.1.0 - Production on Sat Mar 8 14:42:02 2014
Copyright (c) 1982, 2013, Oracle and/or its affiliates. All rights reserved.
Express Mode Load, Table: DEPT
Path used:      External Table, DEGREE_OF_PARALLELISM=AUTO
SQL*Loader-816: error creating temporary directory object SYS_SQLLDR_XT_
TMPDIR_00000 for file dept.dat
ORA-01031: insufficient privileges
SQL*Loader-579: switching to direct path for the load

```

```

SQL*Loader-583: ignoring trim setting with direct path, using value of LDRTRIM
SQL*Loader-584: ignoring DEGREE_OF_PARALLELISM setting with direct path,
using value of NONE
Express Mode Load, Table: DEPT
Path used:      Direct
Load completed - logical record count 4.
Table DEPT:
  4 Rows successfully loaded.
Check the log file:
  dept.log
for more information about the load.
SQL*Loader-816: ошибка при создании объекта временного каталога
SYS_SQLLDR_XT_TMPDIR_00000 для файла dept.dat
ORA-01031: недостаточно привилегий
SQL*Loader-579: переключение в прямой режим для загрузки
SQL*Loader-583: игнорирование установки усечения в прямом режиме,
используется значение LDRTRIM
SQL*Loader-584: игнорирование установки DEGREE_OF_PARALLELISM в прямом
режиме, используется значение NONE
Загрузка в скоростном режиме, таблица: DEPT
Используемый путь:      прямой
Загрузка завершена - количество логических записей 4.
Таблица DEPT:
  4 строки успешно загружены.
Обратитесь в журнальный файл:
  dept.log
за дополнительной информацией о загрузке.

```

Предыдущий вывод сообщает о том, что схема не располагает привилегией на создание объекта каталога; по этой причине SQLLDR применяет метод загрузки данных в прямом режиме. Запрос к таблице подтверждает, что данные были успешно загружены:

```

EODA@ORA12CR1> select * from dept;

```

| DEPTNO | DNAME | LOC |
|--------|------------|----------|
| 10 | Sales | Virginia |
| 20 | Accounting | Virginia |
| 30 | Consulting | Virginia |
| 40 | Finance | Virginia |

Теперь выдадим пользователю привилегию CREATE ANY DIRECTORY:

```

SYS@ORA12CR1> grant create any directory to eoda;
Grant succeeded.
Выдано успешно.

```

Чтобы выполнить пример снова, понадобится удалить записи из таблицы DEPT:

```

EODA@ORA12CR1> truncate table dept;

```

Теперь запустим SQLLDR из командной строки в скоростном режиме:

```

$ sqlldr eoda table=dept
Password:
SQL*Loader: Release 12.1.0.1.0 - Production on Sat Mar 8 14:45:53 2014

```

Copyright (c) 1982, 2013, Oracle and/or its affiliates. All rights reserved.

Express Mode Load, Table: DEPT

Path used: External Table, DEGREE_OF_PARALLELISM=AUTO

Table DEPT:

4 Rows successfully loaded.

Check the log files:

dept.log

dept_%p.log_xt

for more information about the load.

Загрузка в скоростном режиме, таблица: DEPT

Используемый путь: внешняя таблица, DEGREE_OF_PARALLELISM=AUTO

Таблица DEPT:

4 строки успешно загружены.

Обратитесь в журнальные файлы:

dept.log

dept_%p.log_xt

за дополнительной информацией о загрузке.

Поскольку у скоростного режима была возможность создать объект каталога, оператор INSERT обращается к внешней таблице, использующей файл dept.dat в качестве своего источника данных, для загрузки данных в таблицу DEPT. После завершения загрузки внешняя таблица удаляется.

Весь код, требуемый для выполнения предшествующих шагов, генерируется и заносится в файл dept.log. Этот файл автоматически заполняется при выполнении SQLLDR в скоростном режиме. Если вы хотите, чтобы инструмент SQLLDR только сгенерировал журнальный файл, не запуская его содержимое, то укажите параметр EXTERNAL_TABLE=GENERATE_ONLY, например:

```
$ sqlldr eoda table=dept external_table=generate_only
```

Внимательно просмотрев файл dept.log, вы найдете код, который был сгенерирован. Первым встречается управляющий файл SQLLDR (управляющие файлы SQLLDR подробно описаны в разделе “Инструмент SQLLDR” далее в этой главе):

```
OPTIONS (EXTERNAL_TABLE=EXECUTE, TRIM=LRTRIM)
LOAD DATA
INFILE 'dept'
APPEND
INTO TABLE DEPT
FIELDS TERMINATED BY ","
(
  DEPTNO,
  DNAME,
  LOC
)
```

Далее находится код SQL, который будет создавать внешнюю таблицу:

```
CREATE TABLE "SYS_SQLLDR_X_EXT_DEPT"
(
  "DEPTNO" NUMBER(2),
  "DNAME" VARCHAR2(14),
  "LOC" VARCHAR2(13)
)
```



```

ORGANIZATION external
(
  TYPE oracle_loader
  DEFAULT DIRECTORY SYS_SQLLDR_XT_TMPDIR_00000
  ACCESS PARAMETERS
  (
    RECORDS DELIMITED BY NEWLINE CHARACTERSET US7ASCII
    BADFILE 'SYS_SQLLDR_XT_TMPDIR_00000': 'dept_%p.bad'
    LOGFILE 'dept_%p.log_xt'
    READSIZE 1048576
    FIELDS TERMINATED BY "," LRTRIM
    REJECT ROWS WITH ALL NULL FIELDS
    (
      "DEPTNO" CHAR(255),
      "DNAME" CHAR(255),
      "LOC" CHAR(255)
    )
  )
  location
  (
    'dept.dat'
  )
) REJECT LIMIT UNLIMITED

```

За ним следует оператор INSERT в прямом режиме, который может применяться для загрузки данных из внешнего файла в обычную таблицу базы данных (DEPT в этом примере):

```

INSERT /*+ append parallel(auto) */ INTO DEPT
(
  DEPTNO,
  DNAME,
  LOC
)
SELECT
  "DEPTNO",
  "DNAME",
  "LOC"
FROM "SYS_SQLLDR_X_EXT_DEPT"

```

Наконец, временная таблица и объект каталога удаляются:

```

DROP TABLE "SYS_SQLLDR_X_EXT_DEPT"
DROP DIRECTORY SYS_SQLLDR_XT_TMPDIR_00000

```

Идея состоит в том, что когда вы запускаете SQLLDR в скоростном режиме с параметром EXTERNAL_TABLE=GENERATE_ONLY, вы можете использовать содержимое этого журнального файла для выполнения находящихся в нем операторов SQL вручную в среде SQL (при желании).

Совет. За подробными сведениями по всем параметрам, доступным в скоростном режиме SQLLDR, обращайтесь в руководство по утилитам базы данных Oracle (*Oracle Database Utilities*).

Обработка ошибок

Ошибок не бывает только в идеальном мире. Данные во входном файле были бы идеальными, и все они загружались бы корректно. В реальности подобное почти никогда не встречается. Каким же образом отслеживать ошибки в процессе загрузки?

Наиболее распространенный метод предусматривает применение параметра BADFILE. В файл некорректных записей Oracle заносит все записи, которые не удалось обработать. Например, если бы управляющий файл содержал запись с DEPTNO, равным 'ABC', то ее обработка потерпела бы неудачу и запись попала бы в файл некорректных записей, т.к. значение 'ABC' не может быть преобразовано в число. Мы продемонстрируем это в следующем примере.

Для начала поместим в конец файла `demol.ctl` такую строку (добавив строку данных, которая не может быть загружена):

```
ABC,XYZ,Hello
```

Затем с помощью показанной ниже команды удостоверимся в том, что файл `demol.bad` пока еще не существует:

```
EODA@ORA12CR1> host ls -l demol.bad
ls: demol.bad: No such file or directory
```

Отобразим содержимое внешней таблицы:

```
EODA@ORA12CR1> select * from SYS_SQLLDR_X_EXT_DEPT;
```

| DEPTNO | DNAME | LOC |
|--------|------------|----------|
| 10 | Sales | Virginia |
| 20 | Accounting | Virginia |
| 30 | Consulting | Virginia |
| 40 | Finance | Virginia |

Теперь файл некорректных записей существует и можно просмотреть его содержимое:

```
EODA@ORA12CR1> host ls -l demol.bad
-rw-r----- 1 oracle dba 14 Mar  9 10:38 demol.bad

EODA@ORA12CR1> host cat demol.bad
ABC,XYZ,Hello
```

Но как программно проинспектировать эти некорректные записи, а также журнал, который был сгенерирован? К счастью, это легко сделать с использованием еще одной внешней таблицы. Предположим, что мы настроили эту внешнюю таблицу следующим образом:

```
EODA@ORA12CR1> create table et_bad
2  ( text1 varchar2(4000) ,
3    text2 varchar2(4000) ,
4    text3 varchar2(4000)
5  )
6  organization external
7  (type oracle_loader
8    default directory SYS_SQLLDR_XT_TMPDIR_00000
9    access parameters
10  (
```

```

11 records delimited by newline
12 fields
13 missing field values are null
14 ( text1 position(1:4000),
15   text2 position(4001:8000),
16   text3 position(8001:12000)
17 )
18 )
19 location ('demo1.bad')
20 )
21 /

```

Table created.

Таблица создана.

Это просто таблица, которая может читать любой файл без возникновения отказа из-за ошибки, связанной с типом данных, при условии, что строки в этом файле содержат менее 12 000 символов. Если они длиннее 12 000 символов, можно просто добавить больше текстовых столбцов, чтобы уместить их.

Просмотреть отклоненные записи можно с помощью такого запроса:

```
EODA@ORA12CR1> select * from et_bad;
```

```
TEXT1          TEXT2          TEXT3
```

```
-----
ABC,XYZ,Hello
```

Функция COUNT (*) сообщает, сколько записей было отклонено. Другая внешняя таблица, созданная на основе журнального файла, который ассоциирован с исходной внешней таблицей, могла бы сообщать, по какой причине запись была отклонена. Тем не менее, чтобы сделать этот процесс повторяемым, необходимо продвигаться на шаг дальше. Дело в том, что файл некорректных записей не опустошается, даже если во время работы с внешней таблицей ошибки не возникали. Таким образом, если существует какой-то файл некорректных записей, содержащий данные, а внешняя таблица не сгенерировала ошибок, мы будем введены в заблуждение, думая, что ошибки были.

Для решения этой проблемы я использовал четыре подхода.

- Сбросить файл некорректных записей с помощью пакета UTL_FILE, т.е. очистить его, просто открыв для записи и сразу закрыв.
- Переименовать существующие файлы некорректных записей посредством пакета UTL_FILE, сохранив их содержимое и позволив создать новый файл.
- Включить идентификатор процесса (PID) в имена файлов некорректных записей и журнальных файлов. Это будет демонстрироваться в разделе “Проблемы многопользовательского доступа” далее в главе.
- Применить для решения проблем команды операционной системы (вроде переименования файла, его удаления и т.д.).

Поступая в подобной манере, мы будем в состоянии выяснить, были сгенерированы записи в файле некорректных записей только что или же они относятся к какой-то более старой версии этого файла и не имеют значения.

Проецирование указанных столбцов

Функция `COUNT(*)`, использованная ранее в этом разделе, напомнила мне о средстве, доступном в Oracle 10g и последующих версиях — возможности оптимизировать доступ к внешним таблицам за счет обращения только к тем полям во внешнем файле, которые упомянуты в запросе. То есть, если во внешней таблице определено 100 числовых полей, но нужно выбрать только одно из них, то вы можете указать Oracle о том, что преобразовывать остальные 99 строк в числа не нужно. Это звучит замечательно, но может привести к тому, что каждый запрос будет возвращать разное количество строк. Предположим, что внешняя таблица имеет 100 строк данных. Все данные для столбца C1 “допустимы” и преобразуются в числа. Данные для столбца C2 не являются “допустимыми” и в числа не преобразуются. Если вы выберете из этой внешней таблицы данные для столбца C1, то получите 100 строк, а если данные для столбца C2, то 0 строк.

Вы обязаны явно включить упомянутую выше оптимизацию, и должны обдумать, безопасно ли ее применять (только вы достаточно хорошо знаете свое приложение и характер его работы, чтобы суметь ответить на этот вопрос). Используя предшествующий пример с добавлением строки с некорректными данными, при запросе внешней таблицы можно ожидать следующего вывода:

```
EODA@ORA12CR1> select dname from SYS_SQLLDR_X_EXT_DEPT;
DNAME
-----
Sales
Accounting
Consulting
Finance

EODA@ORA12CR1> select deptno from SYS_SQLLDR_X_EXT_DEPT;
DEPTNO
-----
10
20
30
40
```

Нам известно, что некорректные записи были занесены в `BADFILE`. Но если мы просто выполним следующий оператор `ALTER` в отношении внешней таблицы и сообщим Oracle о том, что нужно проецировать (обрабатывать) только указанные столбцы, то получим из каждого запроса разное количество строк:

```
EODA@ORA12CR1> alter table SYS_SQLLDR_X_EXT_DEPT
2 project column referenced
3 /

Table altered.
Таблица изменена.

EODA@ORA12CR1> select dname from SYS_SQLLDR_X_EXT_DEPT;
DNAME
-----
Sales
Accounting
Consulting
Finance
XYZ
```

```
EODA@ORA12CR1> select deptno from SYS_SQLLDR_X_EXT_DEPT;
```

```
DEPTNO
-----
      10
      20
      30
      40
```

Столбец DNAME был допустимым для каждой отдельно взятой записи во входном файле, но столбец DEPTNO таковым не был. Если мы не извлекаем столбец DEPTNO, он не приведет к неудаче при обработке записи — результирующий набор по существу изменился.

Использование внешней таблицы для загрузки разных файлов

Часто возникает потребность в применении внешней таблицы для загрузки данных из файлов, которым в течение определенного периода времени назначались разные имена. То есть на текущей неделе мы должны загрузить file1.dat, на следующей неделе — file2.dat и т.д. До сих пор мы загружали из файла с фиксированным именем, demo1.ct1. Что если после этого нам необходимо выполнить загрузку из второго файла, demo2.ct1?

К счастью, обеспечить это довольно легко. С помощью команды ALTER TABLE можно переопределить настройку расположения внешней таблицы:

```
EODA@ORA12CR1> alter table SYS_SQLLDR_X_EXT_DEPT location( 'demo2.ct1' );
Table altered.
Таблица изменена.
```

После этого следующий запрос к внешней таблице получит доступ к файлу demo2.ct1.

Проблемы многопользовательского доступа

Ранее были описаны три ситуации, когда внешние таблицы могут оказаться не настолько удобными, как SQLLDR. Одна из них была связана со специфической проблемой многопользовательского доступа. Только что было показано, как изменить местоположение внешней таблицы — заставить ее читать данные из файла 2 вместо файла 1 и т.д. Проблема возникает, если несколько пользователей пытаются параллельно работать с внешней таблицей, направляя ее в своих сеансах на разные файлы.

Сделать это невозможно. В любой момент времени внешняя таблица будет указывать на единственный файл (или набор файлов). Если один пользователь войдет в базу данных и изменит таблицу, чтобы она указывала на файл 1, а другой пользователь выполнит аналогичное действие примерно в то же время, после чего они оба обратятся к этой таблице, то оба будут обрабатывать тот же самый файл.

Обычно такая проблема возникать не должна — внешние таблицы не являются заменой таблиц базы данных; это средства для загрузки данных, и потому их не следует использовать на повседневной основе как часть приложения. Внешние таблицы в основном выступают в качестве инструмента, который предназначен для администратора базы данных или разработчика и применяется с целью загрузки информации — либо единовременно, либо периодически как в случае хранилища данных.

Если администратор базы данных имеет два файла, предназначенные для загрузки в базу данных с использованием той же самой внешней таблицы, он *не* обязан делать это последовательно — т.е. нацеливать внешнюю таблицу на файл 1 и обработать его, а затем на файл 2 и обработать его. Вместо этого администратор должен просто направить внешнюю таблицу на *оба файла* и позволить базе данных обработать их оба:

```
EODA@ORA12CR1> alter table SYS_SQLLDR_X_EXT_DEPT
  2 location( 'file1.dat', 'file2.dat')
  3 /
Table altered.
Таблица изменена.
```

Если требуется *параллельная обработка*, то для этого база данных располагает встроенными средствами, как было показано в предыдущей главе.

Таким образом, одна проблема многопользовательского доступа все-таки может возникать, когда оба сеанса пытаются изменить местоположение примерно в одно и то же время (при наличии у них привилегии на выполнение команды ALTER TABLE для таблицы). Однако это просто возможность, о которой следует знать, а не что-то такое, с чем придется сталкиваться очень часто.

Другая проблема многопользовательского доступа связана с именами файла некорректных записей и журнального файла. Что, если вы имеете множество сеансов, которые параллельно просматривают одну и ту же внешнюю таблицу или применяют параллельную обработку (что в некотором смысле является многопользовательской ситуацией)? Было бы неплохо разделять упомянутые файлы на основе сеансов — и к счастью это возможно. В имена этих файлов можно включить следующие специальные строки:

- %p — идентификатор процесса (PID);
- %a — идентификатор агента серверов параллельного выполнения; серверам параллельного выполнения назначаются номера 001, 002, 003 и т.д.

В таком случае каждый сеанс будет генерировать собственный файл некорректных записей и журнальный файл. Например, предположим, что вы используете приведенный ниже синтаксис для BADFILE в показанном ранее операторе CREATE TABLE:

```
RECORDS DELIMITED BY NEWLINE CHARACTERSET US7ASCII
BADFILE 'SYS_SQLLDR_XT_TMPDIR_00000': 'demo1_%p.bad'
LOGFILE 'demo1.log_xt'
```

Если при загрузке встретятся записи, которые не удалось обработать, вы можете обнаружить файл с именем следующего вида:

```
$ ls *.bad
demo1_7108.bad
```

Однако в долговременной перспективе вы все равно можете столкнуться с проблемами. В большинстве операционных систем идентификаторы PID задействуются повторно. Таким образом, рассмотренные приемы работы с ошибками по-прежнему будут актуальными — вам придется очищать файл некорректных записей или переименовывать его, если обнаружится, что он является источником проблем.

Предварительная обработка

Предварительная обработка (preprocessing) — это средство внешних таблиц, которое позволяет выполнять одну и более команд операционной системы в качестве первого шага при выборке из внешней таблицы. Данное средство было добавлено в версию Oracle 11g Release 2, но затем его перенесли также и в версию 10.2.0.5 (так что оно присутствует в Oracle 10g Release 2 и последующих версиях). Предварительная обработка инициируется путем указания конструкции PREPROCESSOR. Входным значением этой конструкции может быть команда операционной системы или сценарий оболочки.

Иногда польза от предварительной обработки не является интуитивно понятной. Рассмотрим перечисленные ниже сценарии применения.

- Динамическое отображение вывода команды операционной системы (или комбинации команд) за счет выборки из внешней таблицы.
- Поиск файлов и фильтрация столбцов и/или строк перед отображением данных.
- Обработка и модификации содержимого файла перед возвращением данных.

Далее будут приведены примеры этих сценариев.

Мониторинг файловой системы посредством SQL

В базах данных некоторых моих заказчиков администраторы баз данных используют файлы данных с автоматическим расширением, но со многими файлами данных, разделяющими ту же самую файловую систему, например:

```
tablespace A, datafiles /u01/oradata/ts_A_file01.dbf autoextend unlimited
tablespace B, datafiles /u01/oradata/ts_B_file01.dbf autoextend unlimited
```

Запрошенное расширение заключается в том, что все файлы данных должны быть способны к увеличению минимум на 20% от их текущего размера; так что если, скажем, файл `ts_A_file01.dbf` имеет текущий размер 100 Гбайт, а файл `ts_B_file01.dbf` — 200 Гбайт, то мы должны обеспечить наличие в файловой системе /u01/oradata, по меньшей мере, 20 Гбайт + 40 Гбайт = 60 Гбайт свободного пространства.

Возникает вопрос: как это можно отследить в единственном запросе внутри базы данных? В настоящий момент у нас имеется сложный сценарий, который собирает информацию о свободном пространстве из команды `df` в текстовый файл, открывает курсор, подсчитывает текущее выделенное пространство с помощью представления `DBA_DATA_FILES` и читает данные команды `df`, находящиеся в текстовом файле, через внешнюю таблицу.

Все перечисленное может быть выполнено в одном SQL-запросе. Чтобы достичь этого, понадобится возможность запрашивать утилиту `df` интерактивно — без сложных манипуляций вроде запуска сценария или перенаправления вывода. Мы начнем с того, что сделаем так, чтобы вывод утилиты `df` можно было запрашивать, как если бы он являлся таблицей. Для этого будет применяться конструкция PREPROCESSOR.

Создадим каталог, куда можно поместить небольшой сценарий оболочки, который будет генерировать вывод `df`:

```
EODA@ORA12CR1> create or replace directory exec_dir as '/orahome/oracle/bin';
Directory created.
```

Каталог создан.

Для безопасного выполнения этой операции нам необходимо выдать привилегию EXECUTE для объекта каталога, содержащего программу, которую мы хотим запускать. Это позволит точно управлять тем, какую программу будет выполнять база данных Oracle, и избегать непреднамеренного запуска программ типа троянских копей. От имени привилегированной учетной записи выдадим привилегию EXECUTE пользователю EODA:

```
SYS@ORA12CR1> grant execute on directory exec_dir to eoda;
```

Далее мы создадим в этом каталоге сценарий оболочки по имени run_df.bsh, который будет содержать следующий код:

```
#!/bin/bash
/bin/df -Pl
```

Сделаем данный сценарий оболочки исполняемым:

```
$ chmod +x run_df.bsh
```

Вывод этого сценария будет выглядеть примерно так:

```
EODA@ORA12CR1> !./run_df.bsh
```

| Filesystem | 512-blocks | Used | Available | Capacity | Mounted on |
|----------------------|------------|-----------|-----------|----------|------------|
| rpool/ROOT/solaris-1 | 205406208 | 4576073 | 92886960 | 5% | / |
| ... | | | | | |
| orapool1/ora01 | 629145600 | 382371882 | 246773718 | 61% | /ora01 |
| orapool2/ora02 | 629145600 | 429901328 | 199244272 | 69% | /ora02 |
| orapool1/ora03 | 629145600 | 415189806 | 213955794 | 66% | /ora03 |
| orapool2/ora04 | 629145600 | 343152972 | 285992628 | 55% | /ora04 |

Обратите внимание, что в сценарии run_df.bsh для запуска утилиты df мы используем явные путевые имена; мы не полагаемся на среду, в частности, на переменную среды path. Это очень важно: при написании сценариев для внешних таблиц — и вообще при написании сценариев — всегда лучше указывать явные пути, чтобы запускать именно те программы, которые действительно должны быть запущены. Вы не имеете никакого реального контроля над средой, в которой будет выполняться сценарий, поэтому расчет на то, что среда должным образом настроена, приведет к катастрофическому отказу.

Итак, располагая сценарием и каталогом, мы готовы к созданию внешней таблицы. В следующем коде видно, что для этого необходимо заставить внешнюю таблицу пропустить первую запись и произвести разбор всех последующих строк, трактуя пробельные символы в качестве разделителей. Как показано ниже, это то, с чем внешняя таблица справляется легко:

```
EODA@ORA12CR1> create table df
```

```
2 (
3   fsname   varchar2(100),
4   blocks   number,
5   used     number,
6   avail    number,
7   capacity varchar2(10),
8   mount    varchar2(100)
9 )
10 organization external
11 (
12   type oracle_loader
```



```

13  default directory exec_dir
14  access parameters
15  (
16    records delimited
17    by newline
18    preprocessor
19    exec_dir:'run_df.bsh'
20    skip 1
21    fields terminated by
22    whitespace ldrtrim
23  )
24  location
25  (
26    exec_dir:'run_df.bsh'
27  )
28  )
29  /

```

Table created.

Таблица создана.

После создания внешней таблицы df вывод утилиты df можно легко просматривать в запросе:

```
EODA@ORA12CR1> select * from df;
```

| FSNAME | BLOCKS | USED | AVAIL | CAPACITY | MOUNT |
|----------------|-----------|-----------|-----------|----------|--------|
| orapool1/ora01 | 629145600 | 382371882 | 246773718 | 61% | /ora01 |
| orapool2/ora02 | 629145600 | 429901326 | 199244274 | 69% | /ora02 |
| orapool1/ora03 | 629145600 | 415189808 | 213955792 | 66% | /ora03 |
| orapool2/ora04 | 629145600 | 343152974 | 285992626 | 55% | /ora04 |

Совет. Вы сочтете этот подход удобным также при работе с утилитами ps, ls, du и т.д. Теперь все утилиты UNIX/Linux могут трактоваться как “таблицы”!

Имея данные, доступные во внешней таблице df, можно приступить к работе над запросом. Понадобится лишь выполнить соединение df с DBA_DATA_FILES, условие которого сопоставляет самую длинную из числа возможных точку монтирования с каждым именем файла. В следующем фрагменте кода показано решение в виде единственного запроса. Ниже обсуждается то, что происходит в некоторых строках.

```

EODA@ORA12CR1> with fs_data
2  as
3  (select /*+ materialize */ *
4    from df
5  )
6  select mount,
7    file_name,
8    bytes/1024/1024 mbytes,
9    tot_bytes/1024/1024 tot_mbytes,
10   avail_bytes/1024/1024 avail_mbytes,
11   case
12     when 0.2 * tot_bytes < avail_bytes

```

```

13      then 'OK'
14      else 'Short on disk space'
15      end status
16  from (
17  select file_name, mount, avail_bytes, bytes,
18         sum(bytes) over
19           (partition by mount) tot_bytes
20  from (
21  select a.file_name,
22         b.mount,
23         b.avail*1024 avail_bytes, a.bytes,
24         row_number() over
25           (partition by a.file_name
26            order by length(b.mount) DESC) rn
27  from dba_data_files a,
28       fs_data b
29  where a.file_name
30        like b.mount || '%'
31        )
32  where rn = 1
33        )
34  order by mount, file_name
35  /

```

| MOUNT | FILE_NAME | MBYTES | TOT_MBYTES | AVAIL_MBYTES | STATUS |
|--------|---|--------|------------|--------------|--------|
| /ora01 | /ora01/dbfile/ORA12CR1/cia_data_01.dbf | 1024 | 93486 | 240989.959 | OK |
| /ora01 | /ora01/dbfile/ORA12CR1/config_tbsp_1_01.dbf | 1500 | 93486 | 240989.959 | OK |
| ... | | | | | |
| /ora04 | /ora04/dbfile/ORA12CR1/dim_data08.dbf | 30720 | 136202 | 279289.674 | OK |
| /ora04 | /ora04/dbfile/ORA12CR1/dim_data_02.dbf | 30720 | 136202 | 279289.674 | OK |

48 rows selected.
48 строк выбрано.

В строках 3 и 4 производится запрос внешней таблицы `df`. Здесь намеренно применяется подсказка `materialize`, чтобы вынудить оптимизатор загружать данные `df` в эквивалент глобальной временной таблицы, поскольку запрос, как правило, будет читать и перечитывать внешнюю таблицу снова и снова, а во время выполнения запроса данные в таблице `df` могут измениться. Это предоставляет логический эквивалент согласованного чтения данных `df`. К тому же, если план выполнения запроса не предусматривал повторного чтения внешней таблицы, будет получено сообщение об ошибке:

```
KUP-04108 unable to reread file string
KUP-04108 не удастся повторно прочитать строку файла
```

Вот как эта ошибка объясняется в документации.

- **Причина.** Выполняемый запрос требует многократного чтения источника данных для внешней таблицы. Тем не менее, источник данных является последовательным устройством, которое не может быть прочитано повторно. Примерами источников данных такого типа служат магнитная лента и конвейер.

- **Действие.** Обойти эту проблему можно несколькими способами. Один из них предусматривает переписывание запроса, чтобы ссылка на внешнюю таблицу производилась только однажды. Второй способ заключается в перемещении источника данных на устройство, поддерживающее повторного чтение, такое как дисковый файл. Третий способ состоит в загрузке данных для внешней таблицы во временную таблицу и изменении запроса так, чтобы он работал с этой временной таблицей.

В строках 27–30 производится соединение представления `DBA_DATA_FILES` и данных `df` с использованием конструкции `WHERE`, внутри которой применяется операция `LIKE`. В результате произойдет соединение каждого файла в `DBA_DATA_FILES` с каждой возможной точкой монтирования в выводе `df`. Однако цель заключается в том, чтобы найти “самую длинную” совпадающую точку монтирования, поэтому в строках 24–26 каждой строке данных назначается `ROW_NUMBER`. Значения `ROW_NUMBER` будут последовательно присваиваться всем дублированным строкам в `DBA_DATA_FILES`, так что если `FILE_NAME` соответствует более чем одному `MOUNT`, то каждому вхождению `FILE_NAME` назначается уникальное увеличивающееся значение `ROW_NUMBER`. Значение `ROW_NUMBER` будет присваиваться после сортировки данных по длине `MOUNT`, от большей к меньшей.

Имея такие данные, к ним можно применить конструкцию `WHERE`, чтобы сохранить только первую запись для каждого значения `FILE_NAME` — предикат выглядит как `WHERE rn = 1` и находится в строке 32. В то же самое время в строках 18–19 добавляется еще один столбец — `TOT_MBYTES`. Он позволит проверять достижение порога в 20%.

На последнем шаге выполняется форматирование и вывод данных. Выводятся только значения интересующих столбцов. Оператор `CASE` в строках 11–15 предназначен для проверки того, что 20% общего количества байтов пространства, выделенного на заданной точке монтирования, не превысило оставшегося количества доступных байтов свободного пространства.

Итак, вы видели, каким образом использовать внешние таблицы для запрашивания вывода команд операционной системы, подобных `df`, `ps`, `find` и `ls`. Кроме того, внешние таблицы можно применять для запрашивания результатов выполнения любых утилит, которые осуществляют запись в стандартный вывод, включая `gunzip`, `sed` и т.д.

Чтение и фильтрация сжатых файлов в дереве каталогов

При работе с крупными загрузками данных вполне обычно использовать сжатые файлы. Когда при обработке загрузки данных приходится иметь дело со сжатыми файлами, как правило, предусматривается начальный шаг для их распаковки и еще один шаг для их загрузки в промежуточные таблицы. Посредством внешних таблиц процесс можно модернизировать, указав внешней таблице о необходимости распаковки данных по мере их чтения из сжатого файла.

Рассмотрим пример. Предположим, что у вас есть множество сжатых с помощью `gzip` файлов, которые вы хотите прочитать и обработать. Более того, пусть эти файлы расположены в разных каталогах и вам необходимо выполнить поиск по разнообразным уровням каталогов, найти сжатые файлы и сделать данные доступными за счет их выбора из внешней таблицы.

Для начала создадим три каталога:

```
$ mkdir /tmp/base
$ mkdir /tmp/base/base2a
$ mkdir /tmp/base/base2b
```

Затем создадим три тестовых файла, сожмем их и поместим каждый в отдельный каталог:

```
$ echo 'base col1,base col2' | gzip > /tmp/base/filebase.csv.gz
$ echo 'base2a col1,base2a col2' | gzip > /tmp/base/base2a/filebase2a.csv.gz
$ echo 'base2b col1,base2b col2' | gzip > /tmp/base/base2b/filebase2b.csv.gz
```

Далее мы создаем два объекта каталогов — один указывает на каталог, который будет содержать сценарий оболочки, а другой указывает на базовый каталог, который служит отправной точкой для поиска файлов, подлежащих обработке:

```
EODA@ORA12CR1> create or replace directory exec_dir as '/orahome/oracle/bin';
Directory created.
Каталог создан.

EODA@ORA12CR1> create or replace directory data_dir as '/tmp';
Directory created.
Каталог создан.
```

После этого мы создадим в каталоге /oracle/home/bin сценарий оболочки имени search_dir.bsh. Внутри сценария находится следующий код:

```
#!/bin/bash
/usr/bin/find $* -name "*.gz" -exec /bin/zcat {} \; | /usr/bin/cut -f1 -d,
```

Сценарий производит поиск, начиная с каталога, который ему передается; он будет искать в базовом каталоге и всех его подкаталогах файлы с расширением .gz. Для каждого найденного файла с помощью zcat производится просмотр распакованных данных. Наконец, посредством cut выводится первый столбец данных.

Следующая команда делает сценарий оболочки исполняемым:

```
$ chmod +x search_dir.bsh
```

Теперь осталось лишь создать внешнюю таблицу, которая применяет search_dir.bsh для отображения данных в сжатых файлах:

```
EODA@ORA12CR1> create table csv
2 ( col1 varchar2(20)
3 )
4 organization external
5 (
6 type oracle_loader
7 default directory data_dir
8 access parameters
9 (
10 records delimited by newline
11 preprocessor exec_dir:'search_dir.bsh'
12 fields terminated by ',' ldrtrim
13 )
14 location
15 (
16 data_dir:'base'
```

```

17      )
18  )
19  /

```

Table created.

Таблица создана.

В выводе видно, что отображается только первый столбец из сжатых файлов; причина в том, что в коде сценария `search_dir.bsh` используется утилита `cut`:

```

EODA@ORA12CR1> select * from csv;
COL1
-----
base2a col1
base col1
base2b col1

```

Обратите внимание на возможность динамического изменения каталога, где необходимо начинать поиск файлов:

```

EODA@ORA12CR1> create or replace directory data_dir as '/tmp/base';
EODA@ORA12CR1> alter table csv location( 'base2a' );

```

Выборка из таблицы возвращает только одну запись:

```

COL1
-----
base2a col1

```

Кроме того, мы можем без труда модифицировать код сценария оболочки и заставить его фильтровать данные на основе другого критерия, такого как результат поиска строки внутри файла CSV (например, `base2`). Создадим для этого сценарий `search_dir2.bsh`:

```

#!/bin/bash
/usr/bin/find $* -name "*.gz" -print0 | /usr/bin/xargs -0 -I {}
/usr/bin/zgrep "base2" {}

```

Сделаем его исполняемым:

```
$ chmod +x search_dir2.bsh
```

Ниже приведено определение новой внешней таблицы, которая позволяет помещать в вывод два столбца:

```

EODA@ORA12CR1> create table csv2
2  ( col1 varchar2(20)
3    ,col2 varchar2(20)
4  )
5  organization external
6  (
7    type oracle_loader
8    default directory data_dir
9    access parameters
10   (
11     records delimited by newline
12     preprocessor exec_dir:'search_dir2.bsh'
13     fields terminated by ',' ldrtrim
14   )

```

```

15     location
16     (
17         data_dir:'base'
18     )
19 )
20 /

```

Table created.

Таблица создана.

Произведем выборку из этой внешней таблицы:

```
EODA@ORA12CR1> select * from csv2;
```

В выводе видно, что таблица возвратила только две строки — те, которые содержат внутри себя строку base2:

| COL1 | COL2 |
|-------------|-------------|
| base2a col1 | base2a col2 |
| base2b col1 | base2b col2 |

Поиск самых крупных файлов

При решении проблем, связанных с дисковым пространством, иногда удобно отображать *N* самых крупных файлов в дереве каталогов. Для этого можно применять предварительно обрабатываемую внешнюю таблицу. Создадим два объекта каталогов:

```

EODA@ORA12CR1> create or replace directory exec_dir as '/orahome/oracle/bin';
EODA@ORA12CR1> create or replace directory data_dir as '/';

```

Затем создадим сценарий оболочки (в /orahome/oracle/bin) по имени flf.bsh со следующим кодом внутри:

```

#!/bin/bash
/usr/bin/find $1 -ls|/bin/sort -nrk7|/usr/bin/head -10|/bin/awk '{print $11,$7}'

```

Сделаем flf.bsh исполняемым:

```
$ chmod +x flf.bsh
```

Создадим внешнюю таблицу с конструкцией PREPROCESSOR:

```

create table flf (fname varchar2(200), bytes number)
organization external (
    type oracle_loader
    default directory exec_dir
    access parameters
    ( records delimited by newline
      preprocessor exec_dir:'flf.bsh'
      fields terminated by whitespace ldrtrim)
    location (data_dir:'u01'));

```

Произведем выборку из внешней таблицы, чтобы получить 10 самых крупных файлов в каталоге /u01 (и его подкаталогах):

```
EODA@ORA12CR1> select * from flf;
```

| FNAME | BYTES |
|--|------------|
| /u01/dbfile/ORA12CR1/temp01.dbf | 1.0737E+10 |
| ... | |
| /u01/app/oracle/unloadir/big_table.dat | 2786618287 |

А теперь предположим, что вы хотите изменить каталог, в котором начинается поиск, с /u01 на /orahome/oracle. Введите следующие команды:

```
EODA@ORA12CR1> create or replace directory data_dir as '/orahome';
EODA@ORA12CR1> alter table flf location(data_dir:'oracle');
```

При выборке из внешней таблицы путь каталога, в котором начинается поиск самых крупных файлов в дереве каталогов, изменился:

```
EODA@ORA12CR1> select * from flf;
FNAME                                                    BYTES
-----
/orahome/oracle/orainst/12.1.0.2/database1.zip          1652417511
...
/orahome/oracle/orainst/12.1.0.2/database2.zip          1212882524
```

Это не является типичным сценарием использования внешней таблицы, а скорее иллюстрацией того, что становится возможным благодаря предварительной обработке.

Усечение символов в файле

Как-то мне приходилось сотрудничать с администратором базы данных, который временами получал от обучающихся менеджеров почтовые сообщения с вложенными электронными таблицами, сопровождаемые вопросом о том, не мог ли бы он загрузить такие таблицы в производственную базу данных. В имеющейся среде прямой доступ через сеть к производственной базе данных отсутствовал, поэтому администратор базы данных не имел возможности загрузки по сети посредством инструмента вроде SQLLDR. В такой ситуации шаги для загрузки данных выглядят следующим образом.

1. Сохранить электронную таблицу на ноутбуке Windows в виде файла CSV.
2. Скопировать этот файл CSV на защищенный сервер, сконфигурированный специально для передачи файлов на производственный сервер, после чего скопировать этот файл на сам производственный сервер.
3. Применить утилиту операционной системы для удаления скрытых символов DOS, встроенных в файл CSV.
4. От имени пользователя Oracle создать внешнюю таблицу на основе этого файла CSV.
5. Использовать SQL для вставки в производственную таблицу данных, выбранных из внешней таблицы.

Я обращаю ваше внимание на шаг 3, поскольку именно его позволяет устранить предварительная обработка (что касается остальных шагов, то администратор базы данных в перспективе найдет для них другие решения).

Предположим, что загруженный файл имеет имя load.csv, находится в каталоге /tmp и содержит такие данные:

```
emergency data|load now^M
more data|must load data^M
```

Символы ^M обозначают возврат каретки в среде Windows и перед загрузкой данных должны быть удалены. Чтобы достичь этого, мы подготовим пару каталогов и

затем применим внешнюю таблицу с конструкцией PREPROCESSOR, которая вызовет сценарий оболочки для удаления специальных символов ^M перед извлечением данных:

```
EODA@ORA12CR1> create or replace directory data_dir as '/tmp';
Directory created.
Каталог создан.

EODA@ORA12CR1> create or replace directory exec_dir as '/orahome/oracle/bin';
Directory created.
Каталог создан.
```

Создадим сценарий оболочки по имени run_sed.bsh, который использует утилиту sed для удаления символов ^M из файла. В этом примере сценарий run_sed.bsh помещается в каталог /orahome/oracle/bin:

```
#!/bin/bash
/bin/sed -e 's/^M//g' $*
```

Совет. При передаче файлов из среды Windows/DOS в UNIX/Linux также рассмотрите возможность применения утилиты dos2unix для удаления нежелательных символов.

Символ ^M в сценарии run_sed.bsh вводится нажатием комбинации <CTRL+V> и затем <CTRL+M> (или вместо <CTRL+M> можете нажать здесь клавишу <Enter>); нельзя просто ввести ^ и M. Это должен быть специальный символ ^M.

Сделаем этот сценарий исполняемым:

```
$ chmod +x run_sed.bsh
```

А вот определение внешней таблицы:

```
EODA@ORA12CR1> create table csv3
2   ( col1 varchar2(20)
3     ,col2 varchar2(20)
4   )
5   organization external
6   (
7     type oracle_loader
8     default directory data_dir
9     access parameters
10    (
11      records delimited by newline
12      preprocessor exec_dir:'run_sed.bsh'
13      fields terminated by '|' ldrtrim
14    )
15    location
16    (
17      data_dir:'load.csv'
18    )
19  )
20  /
Table created.
Таблица создана.
```


Выберем данные из внешней таблицы:

```
EODA@ORA12CR1> select * from csv3;
```

| COL1 | COL2 |
|----------------|----------------|
| emergency data | load now |
| more data | must load data |

Как узнать, что посторонние символы ^M были удалены из всех строк? Проверьте длину COL2:

```
EODA@ORA12CR1> select length(col2) from csv3;
```

| LENGTH (COL2) |
|---------------|
| 8 |
| 14 |

Если скрытые символы ^M не были удалены, то длина столбца COL2 окажется минимум на один байт больше размера символьных данных внутри столбца (не говоря уже о том, что возникали бы некоторые сюрпризы при выполнении поисков и соединений).

Заключительные соображения по поводу предварительной обработки

Предшествующие примеры продемонстрировали всю мощь конструкции PREPROCESSOR для внешней таблицы. Безусловно, получить те же самые результаты можно и без предварительной обработки, но тогда у вас будет больше шагов, больше кода для сопровождения и больше мест, где может произойти отказ.

Специфичные примеры были подобраны так, чтобы показать гибкость предварительной обработки. В первом примере обрабатываемым файлом был сценарий оболочки; какие-либо файлы данных отсутствовали. Данными, которые возвращались из внешней таблицы, служил вывод команды операционной системы.

Во втором примере в предварительной обработке участвовал базовый каталог. Он предоставлял отправную точку для поиска в дереве каталогов сжатых файлов. Затем сжатые файлы распаковывались на лету и с помощью таких команд, как `find`, `cut` и `zgrep`, производилась фильтрация столбцов и/или строк.

В третьем примере использовался сценарий, который с помощью утилиты `sed` удалял нежелательные символы из файла CSV, после чего данные отображались. Это распространенная потребность, возникающая при переносе файлов между платформами DOS и UNIX/Linux.

Упомянутые три примера представили характерные способы применения конструкции PREPROCESSOR, что станет основой для упрощения реализации существующих у вас требований к загрузке данных.

Заключительные соображения по поводу внешних таблиц

В этом разделе мы исследовали внешние таблицы. Они являются средством, доступным в Oracle9i и последующих версиях, которое в большинстве случаев может заменить инструмент SQLLDR. Был показан самый быстрый способ обращения с внешними таблицами: использование SQLLDR для преобразования управляющих файлов, оставшихся из прошлой практики. Мы продемонстрировали ряд приемов для обнаружения и обработки ошибок с помощью файлов некорректных записей,

и также уделили внимание проблемам многопользовательского доступа, касающихся внешних таблиц. Наконец, мы рассмотрели приемы предварительной обработки, позволяющие выполнять команды операционной системы в качестве первого шага при выборке данных из внешней таблицы.

Теперь мы готовы перейти к следующему разделу этой главы, который посвящен выгрузке данных из базы.

Выгрузка Data Pump

В версии Oracle9i появились внешние таблицы как метод чтения внешних данных для их помещения в базу данных. В Oracle 10g была введена возможность движения в противоположном направлении — применение оператора CREATE TABLE для создания внешних данных, т.е. выгрузки информации из базы данных. Начиная с версии Oracle 10g, эти данные извлекаются в патентованном двоичном формате, известном как *формат Data Pump*, который используется в инструментах EXPDP и IMPDP, предоставленных Oracle для перемещения данных между базами. Применять выгрузку во внешнюю таблицу на самом деле очень легко — так же просто, как оператор CREATE TABLE AS SELECT. Для начала понадобится объект каталога:

```
EODA@ORA12CR1> create or replace directory tmp as '/tmp';
Directory created.
Каталог создан.
```

Теперь можно выгрузить данные в этот каталог с использованием простого оператора SELECT, например:

```
EODA@ORA12CR1> create table all_objects_unload
2 organization external
3 ( type oracle_datapump
4   default directory TMP
5   location( 'allobjects.dat' )
6 )
7 as
8 select
9 *
10 from all_objects
11 /
Table created.
Таблица создана.
```

Я преднамеренно выбрал представление ALL_OBJECTS, т.к. это довольно сложное представление с многочисленными соединениями и предикатами. Это показывает, что такой прием выгрузки Data Pump можно применять для извлечения произвольных данных из базы. Вы можете добавить предикаты или что-нибудь другое для извлечения интересующего среза данных.

На заметку! Приведенный пример демонстрирует возможность извлечения произвольных данных из базы. Да, я повторяюсь. С точки зрения безопасности это облегчает любому, кто имеет доступ к информации, извлекать ее откуда угодно. Необходимо контролировать группу людей, которые в состоянии создавать объекты каталогов и производить в них запись, а также имеют доступ к физическому серверу для получения выгруженных данных.

Завершающим шагом может быть копирование файла `allobjects.dat` на другой сервер, скажем, на машину разработчика с целью последующего тестирования и извлечения кода DDL для воссоздания этой таблицы заново:

```
EODA@ORA12CR1> select dbms_metadata.get_ddl( 'TABLE', 'ALL_OBJECTS_UNLOAD' )
from dual;

DBMS_METADATA.GET_DDL('TABLE','ALL_OBJECTS_UNLOAD')
-----
CREATE TABLE "EODA"."ALL_OBJECTS_UNLOAD"
(
  "OWNER" VARCHAR2(128),
  "OBJECT_NAME" VARCHAR2(128),
  "SUBOBJECT_NAME" VARCHAR2(128),
  "OBJECT_ID" NUMBER,
  "DATA_OBJECT_ID" NUMBER,
  "OBJECT_TYPE" VARCHAR2(23),
  "CREATED" DATE,
  "LAST_DDL_TIME" DATE,
  "TIMESTAMP" VARCHAR2(19),
  "STATUS" VARCHAR2(7),
  "TEMPORARY" VARCHAR2(1),
  "GENERATED" VARCHAR2(1),
  "SECONDARY" VARCHAR2(1),
  "NAMESPACE" NUMBER,
  "EDITION_NAME" VARCHAR2(128),
  "SHARING" VARCHAR2(13),
  "EDITIONABLE" VARCHAR2(1),
  "ORACLE_MAINTAINED" VARCHAR2(1)
)
ORGANIZATION EXTERNAL
( TYPE ORACLE_DATAPUMP
  DEFAULT DIRECTORY "TMP"
  LOCATION
  ( 'allobjects.dat'
  )
)
```

В дальнейшем этот фрагмент кода легко загрузить в другую базу данных, как показано ниже, и все — данные загружены:

```
EODA@ORA12CR1> insert /*+ append */ into some_table select * from all_
objects_unload;
```

На заметку! Начиная с версии Oracle 11g, вы можете создавать сжатые и зашифрованные файлы дампа. Такие возможности требуют наличия редакции Enterprise базы данных Oracle, а также опции расширенного сжатия (Advanced Compression Option) и опции расширенной безопасности (Advanced Security Option).

Инструмент SQLLDR

Инструмент SQLLDR представляет собой высокоскоростной пакетный загрузчик данных Oracle. Это исключительно полезный инструмент для помещения данных в базу Oracle из широкого разнообразия форматов плоских файлов. SQLLDR может

использоваться для загрузки громадных объемов данных за невероятно короткий период времени. Он работает в двух режимах:

- обычный путь — для загрузки данных SQLLDR задействует SQL-операторы INSERT;
- прямой путь — SQLLDR не применяет SQL, а форматирует блоки базы данных напрямую.

Загрузка в прямом режиме позволяет читать данные из плоского файла и записывать их непосредственно в форматированные блоки базы данных, минуя механизм SQL, генерирование информации отмены и дополнительно генерирование информации повторения. Параллельная загрузка в прямом режиме является одним из самых быстрых способов наполнения базы данных.

Здесь мы не будем раскрывать каждый отдельный аспект SQLLDR. Вы найдете подробные сведения в руководстве по утилитам базы данных Oracle (*Oracle Database Utilities*), где инструменту SQLLDR в Oracle 12c Release 1 посвящены семь глав. Тот факт, что для SQLLDR выделено целых семь глав, тогда как описание любой другой утилиты, например, DBVERIFY, DBNEWID и LogMiner, занимает одну главу или даже меньше, весьма примечателен. Полный синтаксис и все параметры можно найти в руководстве *Oracle Database Utilities*, а настоящая глава призвана дать ответы на вопросы вида “Как сделать...?”, которые в руководстве не рассматриваются.

Следует отметить, что интерфейс уровня вызовов Oracle (Oracle Call Interface — OCI) позволяет написать собственный загрузчик в прямом режиме на языке C. Это удобно, когда операция, которую нужно выполнить, в SQLLDR неосуществима или требуется тесная интеграция с разрабатываемым приложением. SQLLDR — инструмент командной строки (т.е. отдельная программа). Это не API-интерфейс и не что-то такое, что может быть вызвано, например, из PL/SQL.

Запуск SQLLDR в командной строке без параметров приводит к выводу следующей справки:

```
$ sqlldr
SQL*Loader: Release 12.1.0.1.0 - Production on Sun Mar 9 11:57:29 2014
Copyright (c) 1982, 2013, Oracle and/or its affiliates. All rights reserved.
Использование: SQLLDR ключевое_слово=значение [,ключевое_
слово=значение,...]
Допустимые ключевые слова:
  userid -- имя пользователя/пароль ORACLE
  control -- имя управляющего файла
  log -- имя журнального файла
  bad -- имя файла некорректных записей
  data -- имя файла данных
  discard -- имя файла отвергнутых записей
  discardmax -- разрешенное количество отвергнутых записей (по умолчанию - all)
  skip -- количество пропускаемых логических записей (по умолчанию - 0)
  load -- количество загружаемых логических записей (по умолчанию - all)
  errors -- разрешенное количество ошибок (по умолчанию - 50)
  rows -- количество строк в массиве привязки обычного режима
        или между сохранениями данных прямого режима (по умолчанию:
        обычный режим - 64, прямой режим - all)
```

```

bindsize -- размер массива привязки обычного режима в байтах
           (по умолчанию - 256000)
silent -- подавить выдачу сообщений во время выполнения (header
         (заголовки), feedback (обратная связь), errors (ошибки),
         discards (отброшенные данные), partitions (секции))
direct -- использовать прямой режим (по умолчанию - FALSE)
parfile -- файл параметров: имя файла, содержащего спецификации параметров
parallel -- выполнять параллельную загрузку (по умолчанию - FALSE)
         file -- файл для выделения экстенгов
skip_unusable_indexes -- запретить/разрешить неиспользуемые индексы или
                       индексные секции (по умолчанию - FALSE)
skip_index_maintenance -- не обслуживать индексы, помечать неисправные
                       индексы как неиспользуемые (по умолчанию - FALSE)
commit_discontinued -- фиксировать загруженные строки, когда загрузка
                       прекращена (по умолчанию - FALSE)
readsize -- размер буфера чтения (по умолчанию - 1048576)
external_table -- использовать для загрузки внешнюю таблицу; допустимые
                 значения - NOT_USED, GENERATE_ONLY, EXECUTE (по
                 умолчанию - NOT_USED)
columnarrayrows -- количество строк для массива привязки прямого режима
                 (по умолчанию - 5000)
streamsize -- размер буфера потока в прямом режиме в байтах (по
             умолчанию - 256000)
multithreading -- использовать многопоточность в прямом режиме
resumable -- разрешить или запретить возможность возобновления
            текущего сеанса (по умолчанию - FALSE)
resumable_name -- текстовая строка для помощи в идентификации оператора
                RESUMABLE
resumable_timeout -- время ожидания (в секундах) для оператора RESUMABLE
                  (по умолчанию - 7200)
date_cache -- размер (в записях) кеша преобразования дат
            (по умолчанию - 1000)
no_index_errors -- прерывать загрузку при любой ошибке индекса
                 (по умолчанию - FALSE)
...

```

Формальные определения всех параметров можно найти в главе 8 руководства *Oracle Database Utilities*. Применение нескольких параметров будет продемонстрировано ниже.

Для использования SQLLDR необходим *управляющий файл*. Управляющий файл SQLLDR просто содержит информацию, описывающую входные данные — их компоновку, типы и т.д. — и сведения о целевой таблице (таблицах).

На заметку! Не путайте управляющий файл SQLLDR с управляющим файлом базы данных. Вспомните из главы 3, что управляющий файл базы данных — это небольшой двоичный файл, который хранит каталог файлов, требуемый Oracle, наряду с другой информацией, такой как данные контрольных точек, имя базы данных и т.д.

Управляющий файл может даже содержать данные для загрузки. В следующем примере мы построим простой управляющий файл в пошаговой манере с объяснением всех команд. (Обратите внимание, что числа в скобках слева от кода *не* являются частью управляющего файла; они предназначены только для ссылки на строки.)

В дальнейшем мы будем называть этот файл `demol.ct1`.

```
(1) LOAD DATA
(2) INFILE *
(3) INTO TABLE DEPT
(4) FIELDS TERMINATED BY ' ,'
(5) (DEPTNO, DNAME, LOC )
(6) BEGINDATA
(7) 10,Sales,Virginia
(8) 20,Accounting,Virginia
(9) 30,Consulting,Virginia
(10) 40,Finance,Virginia
```

- **LOAD DATA (1).** Сообщает SQL*Loader, что необходимо делать (в этом случае — загрузить данные). Еще одним действием, которое может выполнять SQL*Loader, является `CONTINUE_LOAD` — продолжить загрузку. Вариант `CONTINUE_LOAD` должен применяться только в случае продолжения многотабличной загрузки в прямом режиме.
- **INFILE * (2).** Сообщает SQL*Loader, что данные, подлежащие загрузке, содержатся внутри самого управляющего файла, как показано выше в строках 6–10. В качестве альтернативы можно указать имя другого файла, который хранит данные. При желании оператор `INFILE` можно переопределить, используя параметр командной строки. Помните о том, что параметры командной строки переопределяют настройки управляющего файла.
- **INTO TABLE DEPT (3).** Сообщает SQL*Loader, в какую таблицу должны загружаться данные (в этом случае — `DEPT`).
- **FIELDS TERMINATED BY ' ,' (4).** Сообщает SQL*Loader, что данные будут находиться в форме значений, разделенных запятыми. Существуют десятки способов описания входных данных для SQL*Loader; это лишь один из наиболее распространенных методов.
- **(DEPTNO, DNAME, LOC) (5).** Сообщает SQL*Loader, какие столбцы загружаются, их порядок во входных данных, а также их типы данных. Типы данных имеют отношение к входному потоку, а не к базе данных. В этом случае принимается стандартная установка `CHAR(255)`, которой вполне достаточно.
- **BEGINDATA (6).** Сообщает SQL*Loader, что описание входных данных завершено, и в последующих строках — с 7 по 10 — находятся действительные данные, предназначенные для загрузки в таблицу `DEPT`.

Это управляющий файл в одном из его наиболее простых и распространенных форматов, который предназначен для загрузки разделенных запятыми данных в таблицу. В главе мы взглянем и на более сложные примеры, но приведенного выше содержимого файла вполне достаточно для первого знакомства. Для применения этого управляющего файла, которому назначено имя `demol.ct1`, понадобится создать пустую таблицу `DEPT`:

```
EODA@ORA12CR1> create table dept
2 ( deptno      number(2) constraint dept_pk primary key,
3   dname       varchar2(14),
4   loc         varchar2(13)
```

```
5 )
6 /
```

Table created.

Таблица создана.

Затем нужно запустить следующую команду:

```
$ sqlldr userid=eoda control=demol.ctl
Password:
SQL*Loader: Release 12.1.0.1.0 - Production on Sun Mar 9 12:03:26 2014
Copyright (c) 1982, 2013, Oracle and/or its affiliates. All rights reserved.
Path used:      Conventional
Commit point reached - logical record count 4
Table DEPT:
  4 Rows successfully loaded.
Check the log file:
  demol.log
Используемый путь:      обычный
Достигнута точка фиксации - количество логических записей 4
Таблица DEPT:
  4 строки успешно загружены.
Обратитесь в журнальный файл:
  demol.log
```

Если таблица не пуста, будет выдано сообщение об ошибке:

```
SQL*Loader-601: For INSERT option, table must be empty. Error on table DEPT
SQLLDR-601: для опции INSERT таблица должна быть пуста. Ошибка в таблице DEPT
```

Причина в том, что в управляющем файле мы разрешили принимать почти всем параметрам стандартные значения, а опцией загрузки по умолчанию является INSERT (в противоположность APPEND, TRUNCATE или REPLACE). В случае INSERT загрузчик SQLLDR предполагает, что таблица пуста. Для *добавления* записей в таблицу DEPT потребовалось бы указать опцию APPEND, а для замены данных в таблице DEPT — опцию REPLACE или TRUNCATE. Для удаления записей опция REPLACE предусматривает использование традиционного оператора DELETE, поэтому если загружаемая таблица содержит много записей, выполнение может проходить довольно медленно. Опция TRUNCATE предполагает применение оператора TRUNCATE и обычно работает быстрее, т.к. физически удалять каждую запись не приходится.

Каждая загрузка будет генерировать журнальный файл. Его содержимое для нашей простой загрузки выглядит следующим образом:

```
SQL*Loader: Release 12.1.0.1.0 - Production on Sun Mar 9 12:03:26 2014
Copyright (c) 1982, 2013, Oracle and/or its affiliates. All rights reserved.

Control File:  demol.ctl
Data File:     demol.ctl
Bad File:      demol.bad
Discard File:  none specified

(Allow all discards)

Number to load: ALL
Number to skip: 0
Errors allowed: 50
```

Bind array: 64 rows, maximum of 256000 bytes
 Continuation: none specified
 Path used: Conventional

Table DEPT, loaded from every logical record.
 Insert option in effect for this table: INSERT

| Column Name | Position | Len | Term | Encl | Datatype |
|-------------|----------|-----|------|------|-----------|
| DEPTNO | FIRST | * | , | | CHARACTER |
| DNAME | NEXT | * | , | | CHARACTER |
| LOC | NEXT | * | , | | CHARACTER |

Table DEPT:

4 Rows successfully loaded.
 0 Rows not loaded due to data errors.
 0 Rows not loaded because all WHEN clauses were failed.
 0 Rows not loaded because all fields were null.

Space allocated for bind array: 49536 bytes (64 rows)
 Read buffer bytes: 1048576

Total logical records skipped: 0
 Total logical records read: 4
 Total logical records rejected: 0
 Total logical records discarded: 0

Run began on Sat Mar 01 10:10:35 2014
 Run ended on Sat Mar 01 10:10:36 2014

Elapsed time was: 00:00:01.01
 CPU time was: 00:00:00.01

Журнальный файл сообщает о многих аспектах произведенной загрузки. Мы видим опцию, которая использовалась (стандартную или заданную явно), количество прочитанных записей, количество загруженных строк и т.д. В журнальном файле указаны местоположения файлов некорректных (BAD) и отброшенных (DISCARDED) записей. В нем даже сообщается время, которое заняла загрузка. Журнальные файлы незаменимы при проверке успешности загрузки, а также при диагностировании ошибок. Если загруженные данные вызывают ошибки SQL (т.е. входные данные были некорректными и привели к созданию записей в файле некорректных записей), то такие ошибки фиксируются в журнальном файле. Информация в журнальных файлах в значительной степени самоочевидна, так что мы не станем тратить время на ее объяснение.

Часто задаваемые вопросы по загрузке данных посредством SQL*Loader

Теперь давайте рассмотрим, по моему мнению, самые часто задаваемые вопросы по загрузке данных в базу Oracle с помощью SQL*Loader.

Почему получено сообщение "exceeds maximum length" (превышена максимальная длина) в файле журнала?

Вероятно, это наиболее часто повторяющийся вопрос об SQL*Loader, который мне приходилось слышать: почему файл журнала содержит вывод вроде показанного ниже?

Record 4: Rejected - Error on table DEPT, column DNAME.

Field in data file exceeds maximum length

Запись 4: отклонена — ошибка в таблице DEPT, столбец DNAME.

Поле в файле данных превышает максимальную длину

Это объясняется тем фактом, что стандартным типом данных в SQL*Loader, применяемым для обработки входной записи, является CHAR(255). Если в таблице присутствуют любые строковые типы данных, которые превосходят этот предел, то загрузчику SQL*Loader необходимо явно сообщить, что входная запись может содержать более 255 символов.

Например, предположим, что вы добавили столбец, который может хранить более 255 символов:

```
EODA@ORA12CR1> alter table dept modify dname varchar2(1000);
```

Table altered.

Таблица изменена.

Имеющийся у вас управляющий файл выглядит следующим образом (многократно повторяющийся фрагмент “дальнейший текст” находится в одной строке входного файла):

```
LOAD DATA
INFILE *
INTO TABLE DEPT
FIELDS TERMINATED BY ','
(DEPTNO, DNAME, LOC )
BEGINDATA
10,Sales,Virginia
20,Accounting,Virginia
30,Consulting,Virginia
40,Finance дальнейший текст дальнейший текст ... <повторяется многократно>
... дальнейший текст,Virginia
```

Запустив SQL*Loader, вы получите сообщение об ошибке, которое было показано ранее. Решить проблему довольно просто:

```
LOAD DATA
INFILE *
INTO TABLE DEPT
FIELDS TERMINATED BY ','
(DEPTNO, DNAME char(1000), LOC )
BEGINDATA
...
```

Вот и все! Просто укажите инструменту SQL*Loader максимальную ширину поля во входной записи, в этом случае — 1000.

Как загрузить данные с разделителями?

Данные с разделителями, или данные, которые разделены специальным символом и возможно заключены в кавычки, представляют собой самый популярный формат файлов данных на сегодняшний день. В среде мэйнфрейма, вероятно, наиболее узнаваемым файловым форматом будут файлы записей фиксированной длины, но в среде UNIX/Linux и Windows нормой являются файлы данных с разделителями. В этом разделе мы исследуем распространенные способы загрузки данных с разделителями.

Самым популярным форматом для данных с разделителями является формат со значениями, разделенными запятыми (comma-separated values — CSV). В этом файловом формате каждое поле данных отделяется от следующего запятой. Текстовые строки могут быть заключены в кавычки, что позволяет строкам содержать внутри себя запятые. Если строка должна включать также и символ кавычки, то по соглашению он дублируется (в приведенном ниже коде вместо " используется ""). Типичный управляющий файл для загрузки данных с разделителями будет выглядеть во многом похожим на тот, который применялся в первом примере, но конструкция FIELDS TERMINATED BY обычно должна задаваться так:

```
FIELDS TERMINATED BY ',' OPTIONALLY ENCLOSED BY ''''
```

Она указывает, что запятая разделяет поля данных, а каждое поле *может* быть заключено в двойные кавычки. Давайте предположим, что мы изменили окончание этого управляющего файла, как показано ниже:

```
FIELDS TERMINATED BY ',' OPTIONALLY ENCLOSED BY ''''
(DEPTNO, DNAME, LOC )
BEGINDATA
10,Sales,"Virginia,USA"
20,Accounting,"Va, ""USA""""
30,Consulting,Virginia
40,Finance,Virginia
```

Когда мы запустим SQLLDR с применением этого управляющего файла, результаты будут следующими:

```
EODA@ORA12CR1> select * from dept;

DEPTNO DNAME          LOC
-----
10 Sales             Virginia,USA
20 Accounting        Va, "USA"
30 Consulting         Virginia
40 Finance            Virginia
```

Обратите внимание на перечисленные далее моменты.

- Virginia,USA в отделе 10. Это результат обработки входных данных "Virginia,USA". Поле входных данных должно быть заключено в кавычки, чтобы сохранить запятую как часть данных. В противном случае запятая будет трактоваться как маркер конца поля, и тогда строка Virginia загрузится без текста USA.
- Va, "USA". Это получилось из входных данных "Va, ""USA""""". Внутри заключенной в кавычки строки SQLLDR трактует двойное вхождение " как одиночное. Чтобы загрузить строку, которая содержит необязательный ограничивающий символ, вы должны обеспечить дублирование этого символа.

Еще одним популярным форматом являются данные, разделенные табуляциями, при котором данные разделяются символом табуляции, а не запятой. Существуют два способа загрузки этих данных с использованием конструкции TERMINATED BY:

- TERMINATED BY X'09' (символ табуляции в шестнадцатеричном формате; в ASCII символ табуляции имеет код 9);
- TERMINATED BY WHITESPACE.

Как вы увидите ниже, эти два способа значительно отличаются по реализации. Загрузим таблицу DEPT из предыдущего примера с применением этого управляющего файла:

```
LOAD DATA
INFILE *
INTO TABLE DEPT
REPLACE
FIELDS TERMINATED BY WHITESPACE
(DEPTNO, DNAME, LOC)
BEGINDATA
10 Sales Virginia
```

На печатной странице это не особенно легко заметить, но между частями данных здесь находятся по *два* символа табуляции. Строка данных в действительности выглядит следующим образом, где \t — универсально распознаваемая управляющая последовательность для знака табуляции:

```
10\t\tSales\t\tVirginia
```

Если использовать этот управляющий файл с конструкцией TERMINATED BY WHITESPACE, как было показано ранее, то результирующие данные в таблице DEPT будут такими:

```
EODA@ORA12CR1> select * from dept;

  DEPTNO  DNAME          LOC
-----
      10   Sales          Virginia
```

Конструкция TERMINATED BY WHITESPACE разбирает строку, найдя первое вхождение пробельного символа (табуляции, пробела или новой строки), и затем продолжает, пока не столкнется с первым *непробельным* символом. Следовательно, при разборе данных полю DEPTNO присваивается значение 10, два подряд символа табуляции трактуются как пробельные, полю DNAME присваивается значение Sales и т.д.

С другой стороны, предположим, что вы применили конструкцию вида FIELDS TERMINATED BY X'09', как показано ниже в модифицированном управляющем файле:

```
...
FIELDS TERMINATED BY X'09'
(DEPTNO, DNAME, LOC )
...
```

Вы обнаружите, что в таблицу DEPT загружены следующие данные:

```
EODA@ORA12CR1> select * from dept;

  DEPTNO  DNAME          LOC
-----
      10   Sales
```

Как только инструмент SQLldr встречает табуляцию, он выводит значение. Следовательно, полю DEPTNO присваивается значение 10, а поле DNAME получает NULL, поскольку между первым и вторым символами табуляции ничего нет. Полю LOC присваивается значение Sales.

Это запланированное поведение конструкций `TERMINATED BY WHITESPACE` и `TERMINATED BY <символ>`. Выбор наиболее подходящей из них в конкретном случае определяется входными данными и тем, как вы хотите их интерпретировать.

Наконец, при загрузке данных с разделителями подобного рода часто возникает потребность пропускать разнообразные столбцы во входной записи. Например, может понадобиться загрузить поля 1, 3 и 5, пропустив поля 2 и 4. Для этого `SQLDR` предлагает ключевое слово `FILLER`. Оно позволяет установить соответствие со столбцом во входной записи, но не помещать его в базу данных. Скажем, если взять таблицу `DEPT` и последний управляющий файл из приведенных выше, то с использованием ключевого слова `FILLER` этот управляющий файл можно изменить, чтобы данные загружались корректно (с пропуском символом табуляции):

```
LOAD DATA
INFILE *
INTO TABLE DEPT
REPLACE
FIELDS TERMINATED BY x'09'
(DEPTNO, dummy1 filler, DNAME, dummy2 filler, LOC)
BEGINDATA
10      Sales      Virginia
```

Теперь результирующая таблица `DEPT` будет выглядеть так:

```
EODA@ORA12CR1> select * from dept;

  DEPTNO  DNAME          LOC
-----
      10  Sales          Virginia
```

Как загрузить данные фиксированного формата?

Часто в вашем распоряжении имеется плоский файл, сгенерированный внешней системой, который содержит записи фиксированной длины с позиционными данными. Например, поле `NAME` занимает байты с 1 по 10, поле `ADDRESS` — байты с 11 по 35 и т.д. Мы посмотрим, как `SQLDR` может импортировать данные подобного вида.

Эти позиционные данные фиксированной ширины являются оптимальным форматом для загрузки с помощью `SQLDR`. Они будут обрабатываться быстрее других, поскольку разбирать такой входной поток данных очень легко. Загрузчику `SQLDR` потребуется сохранить смещения и длины, выражающиеся фиксированным количеством байтов, в записи данных, а извлечение самих полей будет совершенно простым. Если вам необходимо загрузить данные исключительно большого объема, то их преобразование в формат с фиксированными позициями обычно будет самым лучшим подходом. Недостаток файла такого формата в том, что он будет иметь на много больший размер, чем простой файл данных с разделителями.

Для загрузки позиционных данных фиксированной ширины в управляющем файле будет применяться ключевое слово `POSITION`, например:

```
LOAD DATA
INFILE *
INTO TABLE DEPT
REPLACE
( DEPTNO position(1:2),
```

```

DNAME position(3:16),
LOC position(17:29)
)
BEGINDATA
10Accounting Virginia,USA

```

В этом управляющем файле не задействована конструкция FIELDS TERMINATED BY; взамен в нем используется ключевое слово POSITION для указания загрузчику SQLDR на то, где начинаются и заканчиваются поля. Обратите внимание, что с помощью POSITION можно задавать перекрывающиеся позиции и перемещаться вперед и назад внутри записи. Например, изменим таблицу DEPT, как показано ниже:

```

EODA@ORA12CR1> alter table dept add entire_line varchar2(29);
Table altered.
Таблица изменена.

```

Далее применим следующий управляющий файл:

```

LOAD DATA
INFILE *
INTO TABLE DEPT
REPLACE
( DEPTNO position(1:2),
  DNAME position(3:16),
  LOC position(17:29),
  ENTIRE_LINE position(1:29)
)
BEGINDATA
10Accounting Virginia,USA

```

Поле ENTIRE_LINE определено как POSITION(1:29). Оно извлекает свои данные из всех 29 байтов строки входных данных, в то время как другие поля являются подстроками этой строки. Результат работы этого управляющего файла будет таким:

```

EODA@ORA12CR1> select * from dept;

```

| DEPTNO | DNAME | LOC | ENTIRE_LINE |
|--------|------------|--------------|---------------------------|
| 10 | Accounting | Virginia,USA | 10Accounting Virginia,USA |

При использовании POSITION можно указывать относительные или абсолютные смещения. В предыдущем примере мы применяли абсолютные смещения. Мы конкретно отметили, где поля начинаются и где они заканчиваются. Предыдущий управляющий файл можно переписать следующим образом:

```

LOAD DATA
INFILE *
INTO TABLE DEPT
REPLACE
( DEPTNO position(1:2),
  DNAME position(*:16),
  LOC position(*:29),
  ENTIRE_LINE position(1:29)
)

```

```
BEGINDATA
10Accounting    Virginia,USA
```

Символ * инструктирует управляющий файл о том, что обработка должна продолжаться там, где закончилось последнее поле. Поэтому в данном случае (*:16) — это то же самое, что и (3:16). Обратите внимание на возможность смешивания относительных и абсолютных позиций в управляющем файле. Вдобавок при использовании обозначения * можно добавлять смещение. Например, если поле DNAME начинается через 2 байта *после* конца поля DEPTNO, можно было бы указать (*+2:16). В настоящем примере результат был бы идентичен применению (5:16).

Завершающая позиция в конструкции POSITION должна быть абсолютной позицией столбца, где заканчиваются данные. Иногда может быть проще указать только длину каждого поля, особенно если поля располагаются рядом друг с другом, как в предыдущем примере. В таком случае нам пришлось бы просто сообщить SQLLDR о том, что запись начинается с первого байта, и затем указывать длину каждого поля. Это избавило бы от необходимости вычислять смещения начального и конечного байтов для полей, что временами оказывается затруднительным. Чтобы сделать это, мы избавимся от конечной позиции и вместо нее укажем *длину* каждого поля в записи фиксированной длины:

```
LOAD DATA
INFILE *
INTO TABLE DEPT
REPLACE
( DEPTNO      position(1) char(2),
  DNAME       position(*) char(14),
  LOC         position(*) char(13),
  ENTIRE_LINE position(1) char(29)
)
BEGINDATA
10Accounting    Virginia,USA
```

Здесь мы должны сообщить SQLLDR только место начала первого поля и его длину. Каждое последующее поле начинается там, где заканчивается предыдущее, и продолжается в соответствие с заданной длиной. Это не касается последнего поля, в котором мы снова должны указать позицию, т.к. оно начинается с самого начала записи.

Как загружать значения дат?

Загрузка значений дат с помощью SQLLDR довольно прямолинейна, но, похоже, она является распространенным источником путаницы. Необходимо просто использовать в управляющем файле тип данных DATE и указать маску даты, которая должна применяться. Это та же самая маска даты, которая используется с функциями TO_CHAR и TO_DATE в базе данных. Инструмент SQLLDR применит эту маску к данным и загрузит их.

Например, изменим таблицу DEPT, как показано ниже:

```
EODA@ORA12CR1> alter table dept add last_updated date;
Table altered.
Таблица изменена.
```

Мы можем загрузить в эту таблицу данные с помощью следующего управляющего файла:

```
LOAD DATA
INFILE *
INTO TABLE DEPT
REPLACE
FIELDS TERMINATED BY ','
(DEPTNO,
 DNAME,
 LOC,
 LAST_UPDATED date 'dd/mm/yyyy'
)
BEGINDATA
10,Sales,Virginia,1/5/2014
20,Accounting,Virginia,21/6/2014
30,Consulting,Virginia,5/1/2013
40,Finance,Virginia,15/3/2014
```

Результирующая таблица DEPT будет выглядеть так:

```
EODA@ORA12CR1> select * from dept;
```

| DEPTNO | DNAME | LOC | LAST_UPDA |
|--------|------------|----------|-----------|
| 10 | Sales | Virginia | 01-MAY-14 |
| 20 | Accounting | Virginia | 21-JUN-14 |
| 30 | Consulting | Virginia | 05-JAN-13 |
| 40 | Finance | Virginia | 15-MAR-14 |

Как видите, ничего сложного. Достаточно предоставить формат даты в управляющем файле и SQLLDR должным образом преобразует данные. В некоторых случаях может требоваться более мощная функция SQL. Например, входной файл может содержать значения дат во множестве разных форматов: иногда с компонентом времени, иногда без него; иногда в формате DD-ММ-YYYY; иногда в формате DD/ММ/YYYY и т.д. В следующем разделе будет показано, как с использованием функций в SQLLDR преодолеть эти трудности.

Как загружать данные с использованием функций?

В этом разделе вы увидите, как ссылаться на функции во время загрузки данных. Однако имейте в виду, что применение таких функций (включая последовательности базы данных) требует участия механизма SQL, а потому при загрузке в прямом режиме работать не будет.

Использовать функции в SQLLDR будет очень легко, если разобраться в том, как SQLLDR строит свой оператор INSERT. Чтобы функция применялась к полю в сценарии SQLLDR, просто добавьте ее к управляющему файлу в двойных кавычках. Например, пусть имеется таблица DEPT из предшествующих примеров и требуется обеспечить загрузку данных в верхнем регистре. Для этого можно использовать управляющий файл следующего вида:

```
LOAD DATA
INFILE *
INTO TABLE DEPT
REPLACE
```

```

FIELDS TERMINATED BY ','
(DEPTNO,
  DNAME      "upper(:dname)",
  LOC        "upper(:loc)",
  LAST_UPDATED date 'dd/mm/yyyy'
)
BEGINDATA
10,Sales,Virginia,1/5/2014
20,Accounting,Virginia,21/6/2014
30,Consulting,Virginia,5/1/2013
40,Finance,Virginia,15/3/2014

```

Результирующие данные в базе будут такими:

```
EODA@ORA12CR1> select * from dept;
```

| DEPTNO | DNAME | LOC | LAST_UPDA | ENTIRE_LINE |
|--------|------------|----------|-----------|-------------|
| 10 | SALES | VIRGINIA | 01-MAY-14 | |
| 20 | ACCOUNTING | VIRGINIA | 21-JUN-14 | |
| 30 | CONSULTING | VIRGINIA | 05-JAN-13 | |
| 40 | FINANCE | VIRGINIA | 15-MAR-14 | |

Обратите внимание, насколько легко привести данные к верхнему регистру, просто применив функцию UPPER к переменной привязки. Следует отметить, что функции SQL могут ссылаться на любой из столбцов независимо от столбца, к которому функция действительно применена. Это означает, что столбец может быть результатом выполнения функции на двух или большем числе других столбцов. Например, если вы хотите загрузить столбец ENTIRE_LINE, то можете использовать операцию конкатенации SQL. Хотя в данном случае это несколько сложнее. Прямо сейчас набор входных данных содержит четыре элемента данных. Предположим, что вы просто добавили столбец ENTIRE_LINE к управляющему файлу, как показано ниже:

```

LOAD DATA
INFILE *
INTO TABLE DEPT
REPLACE
FIELDS TERMINATED BY ','
(DEPTNO,
  DNAME      "upper(:dname)",
  LOC        "upper(:loc)",
  LAST_UPDATED date 'dd/mm/yyyy',
  ENTIRE_LINE ":deptno||:dname||:loc||:last_updated"
)
BEGINDATA
10,Sales,Virginia,1/5/2014
20,Accounting,Virginia,21/6/2014
30,Consulting,Virginia,5/1/2013
40,Finance,Virginia,15/3/2014

```

Вы обнаружите в журнальном файле следующее сообщение об ошибке для каждой входной записи:

```

Record 1: Rejected - Error on table DEPT, column ENTIRE_LINE.
Column not found before end of logical record (use TRAILING NULLCOLS)

```


*Запись 1: отклонена — ошибка в таблице DEPT, столбец ENTIRE_LINE.
Столбец не найден до конца логической записи (используйте TRAILING NULLCOLS)*

Здесь загрузчик SQLLDR уведомляет о том, что данные в записи закончились еще до того, как были исчерпаны столбцы. Решить проблему в этом случае просто, и фактически SQLLDR даже подсказывает, что нужно сделать: применить конструкцию TRAILING NULLCOLS. Она заставит SQLLDR привязывать значение NULL к столбцам, для которых отсутствуют данные во входной записи. В рассматриваемой ситуации добавление TRAILING NULLCOLS приведет к тому, что переменная привязки :ENTIRE_LINE получит значение NULL. Давайте повторим попытку загрузки, но с показанным ниже управляющим файлом:

```
LOAD DATA
INFILE *
INTO TABLE DEPT
REPLACE
FIELDS TERMINATED BY ','
TRAILING NULLCOLS
(DEPTNO,
 DNAME          "upper(:dname)",
 LOC            "upper(:loc)",
 LAST_UPDATED   date 'dd/mm/yyyy',
 ENTIRE_LINE    ":deptno||:dname||:loc||:last_updated"
)
BEGINDATA
10,Sales,Virginia,1/5/2014
20,Accounting,Virginia,21/6/2014
30,Consulting,Virginia,5/1/2013
40,Finance,Virginia,15/3/2014
```

Теперь данные в таблице выглядят следующим образом:

```
EODA@ORA12CR1> select * from dept;
```

| DEPTNO | DNAME | LOC | LAST_UPDA | ENTIRE_LINE |
|--------|------------|----------|-----------|-------------------------------|
| 10 | SALES | VIRGINIA | 01-MAY-14 | 10SalesVirginia1/5/2014 |
| 20 | ACCOUNTING | VIRGINIA | 21-JUN-14 | 20AccountingVirginia21/6/2014 |
| 30 | CONSULTING | VIRGINIA | 05-JAN-13 | 30ConsultingVirginia5/1/2013 |
| 40 | FINANCE | VIRGINIA | 15-MAR-14 | 40FinanceVirginia15/3/2014 |

Этот трюк становится возможным благодаря способу, которым SQLLDR строит оператор INSERT. Загрузчик SQLLDR видит в управляющем файле столбцы DEPTNO, DNAME, LOC, LAST_UPDATED и ENTIRE_LINE. Он настраивает пять переменных привязки с именами как у перечисленных столбцов. Обычно при отсутствии функций оператор INSERT строится просто:

```
INSERT INTO DEPT ( DEPTNO, DNAME, LOC, LAST_UPDATED, ENTIRE_LINE )
VALUES ( :DEPTNO, :DNAME, :LOC, :LAST_UPDATED, :ENTIRE_LINE );
```

Затем SQLLDR должен проанализировать входной поток, присвоить значения переменным привязки и выполнить оператор. Когда вы начинаете использовать функции, SQLLDR включает их в оператор INSERT. В предыдущем примере построенный загрузчиком SQLLDR оператор INSERT будет выглядеть так:

```
INSERT INTO T (DEPTNO, DNAME, LOC, LAST_UPDATED, ENTIRE_LINE)
```

```
VALUES (:DEPTNO, upper(:dname), upper(:loc), :last_updated,
       :deptno||:dname||:loc||:last_updated);
```

Затем SQLldr производит подготовку и привязку входных данных к этому оператору и выполняет его. Таким образом, в сценарии SQLldr можно помещать довольно многие средства, допускаемые в SQL. С добавлением оператора CASE в SQL это может делаться исключительно легко. Предположим, что входной файл содержит данные в следующих форматах:

- HH24:MI:SS. Только время; по умолчанию датой будет первый день текущего месяца.
- DD/MM/YYYY. Только дата; по умолчанию временем будет полночь.
- HH24:MI:SS DD/MM/YYYY. Дата и время указаны явно.

Можно воспользоваться управляющим файлом такого вида:

```
LOAD DATA
INFILE *
INTO TABLE DEPT
REPLACE
FIELDS TERMINATED BY ','
TRAILING NULLCOLS
(DEPTNO,
 DNAME          "upper(:dname)",
 LOC            "upper(:loc)",
 LAST_UPDATED
"case
when length(:last_updated) > 9
then to_date(:last_updated,'hh24:mi:ss dd/mm/yyyy')
when instr(:last_updated,':') > 0
then to_date(:last_updated,'hh24:mi:ss')
else to_date(:last_updated,'dd/mm/yyyy')
end"
)
BEGINDATA
10,Sales,Virginia,12:03:03 17/10/2014
20,Accounting,Virginia,02:23:54
30,Consulting,Virginia,01:24:00 21/10/2014
40,Finance,Virginia,17/8/2014
```

Результаты будут следующими:

```
EODA@ORA12CR1> alter session set nls_date_format = 'dd-mon-yyyy hh24:mi:ss';
Session altered.
```

Сеанс изменен.

```
EODA@ORA12CR1> select deptno, dname, loc, last_updated from dept;
```

| DEPTNO | DNAME | LOC | LAST_UPDATED |
|--------|------------|----------|----------------------|
| 10 | SALES | VIRGINIA | 17-oct-2014 12:03:03 |
| 20 | ACCOUNTING | VIRGINIA | 01-mar-2014 02:23:54 |
| 30 | CONSULTING | VIRGINIA | 21-oct-2014 01:24:00 |
| 40 | FINANCE | VIRGINIA | 17-aug-2014 00:00:00 |

Теперь к входной символьной строке будет применен один из трех форматов даты (обратите внимание, что мы больше *не* загружаем значение типа DATE; мы просто загружаем строку). Оператор CASE проверит длину и содержимое строки, чтобы выяснить, какая из масок должна использоваться.

Интересно отметить, что вы можете написать *собственные* функции, предназначенные для вызова в SQL*Loader. Это очевидное следствие того факта, что код PL/SQL может быть вызван из SQL.

Как загрузить данные со встроенными символами новой строки?

Исторически сложилось так, что загрузка данных свободной формы, которые могут включать символы новой строки, в SQL*Loader была проблематичной. Символ новой строки является стандартным символом конца строки для SQL*Loader, и обходные пути, предпринимаемые в прошлом, не обеспечивали достаточной гибкости. К счастью, в Oracle 8.1.6 и последующих версиях появилось несколько новых приемов. Ниже перечислены варианты загрузки данных со встроенными символами новой строки.

- Загружать данные, в которых новая строка представлена каким-то другим символом (например, символ новой строки может обозначаться в тексте комбинацией \n), и во время загрузки заменять его значением CHR(10) с помощью SQL-функции.
- Применять атрибут FIX в директиве INFILE и загружать плоский файл фиксированной длины. В этом случае ограничители записи отсутствуют; взамен для определения начала и конца записи используется тот факт, что каждая запись в точности равна по длине любой другой записи.
- Применять атрибут VAR в директиве INFILE и загружать файл с записями переменной длины, использующий формат, при котором первые несколько байтов каждой строки указывают ее длину в байтах.
- Применять атрибут STR в директиве INFILE и загружать файл записей переменной длины с определенной последовательностью символов, которая представляет конец строки, в противоположность использованию для этого одного лишь символа новой строки.

Все эти варианты по очереди демонстрируются в последующих разделах.

Использование символа, отличного от символа новой строки

Этот метод легко реализуем, если вы управляете тем, как генерируются входные данные. Если выполнять преобразование при создании файла данных достаточно легко, то такой подход будет работать хорошо. Идея заключается в применении к данным, направляющимся в базу, SQL-функции, которая заменяет определенную символьную строку символом новой строки. Давайте добавим в таблицу DEPT еще один столбец:

```
EODA@ORA12CR1> alter table dept add comments varchar2(4000);
Table altered.
Таблица создана.
```

Мы будем использовать этот столбец для загрузки текста. Пример управляющего файла с встроенными данными может выглядеть следующим образом:

```

LOAD DATA
INFILE *
INTO TABLE DEPT
REPLACE
FIELDS TERMINATED BY ','
TRAILING NULLCOLS
(DEPTNO,
  DNAME      "upper(:dname)",
  LOC        "upper(:loc)",
  COMMENTS   "replace(:comments, '\\n', chr(10))"
)
BEGINDATA
10,Sales,Virginia,This is the Sales\nOffice in Virginia
20,Accounting,Virginia,This is the Accounting\nOffice in Virginia
30,Consulting,Virginia,This is the Consulting\nOffice in Virginia
40,Finance,Virginia,This is the Finance\nOffice in Virginia

```

Обратите внимание, что в вызове функции `replace` указано `\\n`, а не `\n`. Причина в том, что `\n` распознается `SQLldr` как символ новой строки, и загрузчик должен преобразовать `\n` в символ новой строки, а не в двухсимвольную строку. В результате запуска `SQLldr` с предыдущим управляющим файлом в таблицу `DEPT` загружается такое содержимое:

```

EODA@ORA12CR1> select deptno, dname, comments from dept;

```

| DEPTNO | DNAME | COMMENTS |
|--------|------------|--|
| 10 | SALES | This is the Sales Office in Virginia |
| 20 | ACCOUNTING | This is the Accounting Office in Virginia |
| 30 | CONSULTING | This is the Consulting Office in Virginia |
| 40 | FINANCE | This is the Finance Office in Virginia |

Использование атрибута `FIX`

Атрибут `FIX` — это еще один доступный метод. Когда он применяется, входные данные должны находиться в записях фиксированной длины. Каждая запись будет иметь точно такое же количество байтов, как и любая другая. Атрибут `FIX` особенно целесообразен для *позиционных* данных. Такие файлы обычно содержат записи фиксированной длины. При использовании данных свободной формы с разделителями менее вероятно, что мы будем иметь дело с файлом фиксированной длины, т.к. длина этих файлов обычно варьируется (в этом и состоит смысл файлов данных с разделителями: сделать каждую строку лишь настолько длинной, насколько она должна быть).

Когда применяется атрибут `FIX`, должна использоваться конструкция `INFILE`, поскольку именно к ней относится атрибут `FIX`. Кроме того, данные должны храниться внешним образом, а не в самом управляющем файле. Таким образом, исходя из того, что мы имеем записи фиксированной длины, можно применять управляющий файл вроде показанного ниже:

```

LOAD DATA
INFILE demo.dat "fix 80"
INTO TABLE DEPT
REPLACE
FIELDS TERMINATED BY ','
TRAILING NULLCOLS
(DEPTNO,
  DNAME      "upper(:dname)",
  LOC        "upper(:loc)",
  COMMENTS
)

```

Здесь указан входной файл данных с записями по 80 байтов каждая. *Это включает хвостовой символ новой строки*, который может как присутствовать, так и нет. В рассматриваемом случае символ новой строки не играет какой-то специальной роли во входном файле данных. Он является просто еще одним символом, который должен или не должен быть загружен. Важно понимать, что символ новой строки в конце записи (если он есть) станет частью записи. Для полного понимания понадобится утилита вывода содержимого файла на экран, которая поможет удостовериться, что символ новой строки действительно на месте. В среде UNIX/Linux это легко сделать посредством `od` — программы для вывода дампа файлов на экран в восьмеричном и других форматах. Мы будем использовать файл `demo.dat`, содержимое которого приведено ниже. Обратите внимание, что первый столбец в следующем выводе представлен в восьмеричном виде, поэтому число 0000012 во второй строке является восьмеричным и представляет десятичное значение 10. Оно указывает, какой байт в файле мы просматриваем. Вывод сформатирован так, чтобы в каждой строке присутствовало по десять символов (с помощью опции `-w10`), поэтому 0, 12, 24 и 36 — это на самом деле 0, 10, 20 и 30:

```

[tkyte@desktop tkyte]$ od -c -w10 -v demo.dat
0000000 1 0 , S a l e s , V
0000012 i r g i n i a , T h
0000024 i s i s t h e
0000036 S a l e s \n O f f i
0000050 c e i n V i r g
0000062 i n i a
0000074
0000106
0000120 2 0 , A c c o u n t
0000132 i n g , V i r g i n
0000144 i a , T h i s i s
0000156 t h e A c c o u
0000170 n t i n g \n O f f i
0000202 c e i n V i r g
0000214 i n i a
0000226
0000240 3 0 , C o n s u l t
0000252 i n g , V i r g i n
0000264 i a , T h i s i s
0000276 t h e C o n s u
0000310 l t i n g \n O f f i
0000322 c e i n V i r g
0000334 i n i a

```

```

0000346
0000360 4 0 , F i n a n c e
0000372 , V i r g i n i a ,
0000404 T h i s i s t h
0000416 e F i n a n c e \n
0000430 O f f i c e i n
0000442 V i r g i n i a
0000454
0000466
0000500
[tkyte@desktop tkyte]$

```

Как видите, в этом входном файле символы новой строки (\n) присутствуют не для указания загрузчику SQLldr на то, где заканчивается запись; здесь они просто представляют собой данные для загрузки. Чтобы вычислить объем данных, подлежащих чтению, SQLldr применяет ширину FIX в 80 байтов. В действительности, если просмотреть входные данные, то записи для SQLldr даже не завершаются посредством \n. Непосредственно перед записью об отделе 20 находится пробел, а не символ новой строки.

Теперь, когда известно, что каждая запись имеет длину 80 байтов, мы готовы выполнить загрузку с использованием показанного ранее управляющего файла с конст-рукцией FIX 80. Вот как выглядят результаты:

```

EODA@ORA12CR1> select ''' || comments || ''' comments from dept;
COMMENTS
-----
"This is the Sales
Office in Virginia          "
"This is the Accounting
Office in Virginia         "
"This is the Consulting
Office in Virginia         "
"This is the Finance
Office in Virginia         "

```

Вам может понадобиться *усечь* эти данные, т.к. в них сохранились хвостовые пробелы. Такое действие можно предпринять в управляющем файле посредством встроенной SQL-функции TRIM.

Одно предостережение тем, кому посчастливилось работать и в Windows, и в UNIX/Linux: маркер конца строки на этих платформах отличается. В UNIX/Linux это просто \n (CHR(10) в SQL). В Windows/DOS это сочетание \r\n (CHR(13)||chr(10) в SQL). В общем случае, если вы применяете подход с атрибутом FIX, то обязательно обеспечьте *создание и загрузку* файла на однородных платформах (UNIX/Linux и UNIX/Linux или Windows и Windows).

Использование атрибута VAR

Другой метод загрузки данных со встроенными символами новой строки предусматривает применение атрибута VAR. При использовании такого формата каждая запись начинается с некоторого фиксированного количества байтов, которые представляют общую длину поступающей записи.

С применением этого формата можно загружать записи переменной длины, которые содержат встроенные символы новой строки, но только если в начале *каждой записи* имеется *поле длины записи*. Итак, предположим, что мы используем управляющий файл такого вида:

```
LOAD DATA
INFILE demo.dat "var 3"
INTO TABLE DEPT
REPLACE
FIELDS TERMINATED BY ','
TRAILING NULLCOLS
(DEPTNO,
 DNAME      "upper(:dname)",
 LOC        "upper(:loc)",
 COMMENTS
)
```

Конструкция VAR 3 говорит о том, что первые 3 байта каждой входной записи содержат длину этой записи. Загрузим с помощью этого управляющего файла следующий файл данных:

```
[tkyte@desktop tkyte]$ cat demo.dat
05510,Sales,Virginia,This is the Sales
Office in Virginia
06520,Accounting,Virginia,This is the Accounting
Office in Virginia
06530,Consulting,Virginia,This is the Consulting
Office in Virginia
05940,Finance,Virginia,This is the Finance
Office in Virginia
[tkyte@desktop tkyte]$
```

Во входном файле данных присутствуют четыре строки. Первая строка начинается с 055, и это означает, что следующие 55 байтов представляют первую входную запись. Эти 55 байтов включают завершающий символ новой строки после слова Virginia. Следующая строка начинается с 065. Она содержит 65 байтов текста и т.д. С применением такого формата файлов данных мы можем легко загружать данные со встроенными символами новой строки.

И снова, если используются среды UNIX/Linux и Windows (предыдущий пример был ориентирован на UNIX/Linux, где новая строка обозначается одним символом), придется подкорректировать поле длины каждой записи. В среде Windows файл .dat из предыдущего примера должен содержать для длин записей значения 56, 66, 66 и 60.

Использование атрибута STR

Вероятно, это наиболее гибкий метод загрузки данных со встроенными символами новой строки. С применением атрибута STR мы можем указать новый символ конца строки (или последовательность символов). Это позволит создать файл входных данных, который в конце каждой строки содержит некоторый специальный символ, т.е. символ новой строки больше не является “специальным”.

Я предпочитаю использовать последовательность символов — обычно некоторый специальный маркер, а затем символ новой строки. Это позволяет легко замечать символы конца строки при просмотре входных данных в текстовом редакторе или другой утилите, т.к. каждая запись по-прежнему содержит в конце символ новой строки. Атрибут STR задается в шестнадцатеричной форме. Возможно, наиболее простой способ получения интересующей шестнадцатеричной строки предусматривает применение для этого кода SQL и пакета UTL_RAW. Например, если предположить, что работа ведется в среде UNIX/Linux, где маркером конца строки является CHR(10), а специальным маркером — символ '|', то можно поступить так:

```
EODA@ORA12CR1> select utl_raw.cast_to_raw( '|' || chr(10) ) from dual;
UTL_RAW.CAST_TO_RAW(''|CHR(10))
-----
7C0A
```

Вывод показывает, что в атрибуте STR в среде UNIX/Linux должна быть указана строка X'7C0A'.

На заметку! В среде Windows должен использоваться вызов UTL_RAW.CAST_TO_RAW('|' || chr(13) || chr(10)).

Управляющий файл может иметь следующий вид:

```
LOAD DATA
INFILE demo.dat "str X'7C0A'"
INTO TABLE DEPT
REPLACE
FIELDS TERMINATED BY ','
TRAILING NULLCOLS
(DEPTNO,
 DNAME      "upper(:dname)",
 LOC        "upper(:loc)",
 COMMENTS
)
```

Если входной файл выглядит так, как показано ниже, где каждая запись в файле данных заканчивается сочетанием |\n, то предыдущий управляющий файл загрузит его корректно:

```
[tkyte@desktop tkyte]$ cat demo.dat
10,Sales,Virginia,This is the Sales
Office in Virginia|
20,Accounting,Virginia,This is the Accounting
Office in Virginia|
30,Consulting,Virginia,This is the Consulting
Office in Virginia|
40,Finance,Virginia,This is the Finance
Office in Virginia|
[tkyte@desktop tkyte]$
```

Заключительные соображения по поводу встроенных символов новой строки

Мы исследовали, по меньшей мере, четыре способа загрузки данных со встроенными символами новой строки. В разделе “Выгрузка в плоский файл” далее в главе

один из этих приемов — атрибут `STR` — будет применяться в утилите выгрузки общего назначения, чтобы избежать проблем, связанных с наличием в тексте символов новой строки.

Кроме того, необходимо хорошо помнить о том, что в Windows (всех версий) строки в текстовых файлах могут заканчиваться символами `\r\n` (ASCII-значение 13 и ASCII-значение 10, т.е. возврат каретки/перевод строки). В управляющем файле должно быть учтено то, что символ `\r` является частью записи. Это касается счетчиков байтов в атрибутах `FIX` и `VAR`, а также строки, указываемой в `STR`. Например, если вы возьмете любой предшествующий файл `.dat`, который в настоящее время содержит только `\n`, и передадите его по FTP в систему Windows, используя режим ASCII (по умолчанию), то каждое вхождение `\n` превратится в `\r\n`. Тот же самый управляющий файл, который нормально работал в среде UNIX/Linux, больше не сможет загружать данные. Это то, о чем следует помнить и принимать во внимание при подготовке управляющих файлов.

Как загрузить объекты LOB?

Теперь давайте рассмотрим некоторые методы загрузки данных в LOB. Речь идет не о полях типа `LONG` или `LONG RAW`, а о более предпочтительных типах данных `BLOB` и `CLOB`. Эти типы данных появились в Oracle 8.0; как обсуждалось в главе 12, они поддерживают намного более развитый интерфейс и набор функциональности, чем унаследованные типы `LONG` или `LONG RAW`.

Мы ознакомимся с двумя методами загрузки этих полей: `SQLLDR` и `PL/SQL`. Существуют и другие способы, такие как потоки Java, Pro*C и OCI. Мы начнем исследование с метода загрузки объектов LOB посредством `PL/SQL`, а затем взглянем на их загрузку с применением инструмента `SQLLDR`.

Загрузка LOB с помощью PL/SQL

Пакет `DBMS_LOB` имеет точки входа с именами `LoadFromFile`, `LoadBLOBFromFile` и `LoadCLOBFromFile`. Эти процедуры позволяют использовать тип `BFILE` (который может применяться для чтения файлов операционной системы) для наполнения полей `BLOB` или `CLOB` в базе данных. Между процедурами `LoadFromFile` и `LoadBLOBFromFile` нет существенной разницы помимо того, что `LoadBLOBFromFile` возвращает параметры `OUT`, которые показывают, до какой позиции в столбце `BLOB` были загружены данные. Однако процедура `LoadCLOBFromFile` предоставляет важную дополнительную возможность: преобразование символьного набора. Если вы помните, в главе 12 мы обсуждали средства `NLS` (National Language Support — поддержка национальных языков) базы данных Oracle и важность символьных наборов. Процедура `LoadCLOBFromFile` позволяет сообщить базе данных, что загружаемый файл использует символьный набор, отличный от символьного набора базы данных, и должно быть выполнено преобразование этого символьного набора. Например, база данных может быть совместима с UTF8, но файлы, подлежащие загрузке, закодированы в наборе символов `WE8ISO8859P1`, или наоборот. Эта функция позволит успешно загрузить такие файлы.

На заметку! Подробные сведения о процедурах, доступных в пакете `DBMS_LOB`, а также обо всех их входных и выходных параметрах можно найти в справочнике по пакетам и типам `PL/SQL` (*Oracle Database PL/SQL Packages and Types Reference*).

Для применения этих процедур в базе данных понадобится создать объект каталога. Этот объект позволит создавать объекты BFILE (и открывать их), указывающие на файлы в файловой системе, к которой сервер базы данных имеет доступ. Последнее утверждение — “к которой сервер базы данных имеет доступ” — это ключевой момент при использовании PL/SQL для загрузки объектов LOB. Пакет DBMS_LOB выполняется целиком на сервере. Он может иметь дело только с файловой системой, которую способен видеть сервер. В частности, он не может видеть локальную файловую систему, если вы обращаетесь к базе данных Oracle через сеть.

Итак, начинать следует с создания в базе данных объекта каталога. Это простой процесс. Для этого примера мы создадим два объекта каталогов (обратите внимание, что примеры здесь выполняются в среде UNIX/Linux; для ссылки на каталоги вы должны применять синтаксис, принятый в вашей операционной системе):

```
EODA@ORA12CR1> create or replace directory dir1 as '/tmp/';
Directory created.
Каталог создан.
```

```
EODA@ORA12CR1> create or replace directory "dir2" as '/tmp/';
Directory created.
Каталог создан.
```

На заметку! Объекты DIRECTORY в Oracle представляют собой логические каталоги, т.е. являются указателями на существующие физические каталоги внутри операционной системы. Команда CREATE DIRECTORY в действительности не создает каталог в файловой системе — вы должны выполнить эту операцию отдельно.

Пользователь, который выполняет эту операцию, должен располагать привилегией CREATE ANY DIRECTORY. Мы создаем два каталога для демонстрации распространенной проблемы с регистром (верхним и нижним регистрами символов) в объектах DIRECTORY. Когда база данных Oracle создает первый каталог DIR1, она по умолчанию сохраняет имя объекта в *верхнем* регистре. Во втором примере с dir2 она создает объект DIRECTORY, сохраняя регистр, используемый в имени. Важность этого будет объясняться ниже во время применения объекта BFILE.

Теперь мы хотим загрузить какие-то данные либо в BLOB, либо в CLOB. Сделать это довольно легко, например:

```
EODA@ORA12CR1> create table demo
2 ( id          int primary key,
3   theClob     clob
4 )
5 /
Table created.
Таблица создана.

EODA@ORA12CR1> host echo 'Hello World!' > /tmp/test.txt

EODA@ORA12CR1> declare
2   l_clob     clob;
3   l_bfile    bfile;
4 begin
5   insert into demo values ( 1, empty_clob() )
6   returning theclob into l_clob;
7
```

```

8      l_bfile := bfilename( 'DIR1', 'test.txt' );
9      dbms_lob.fileopen( l_bfile );
10
11      dbms_lob.loadfromfile( l_clob, l_bfile,
12                           dbms_lob.getlength( l_bfile ) );
13
14      dbms_lob.fileclose( l_bfile );
15  end;
16  /
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.
EODA@ORA12CR1> select dbms_lob.getlength(theClob), theClob from demo
2 /
DBMS_LOB.GETLENGTH(THECLOB) THECLOB
-----
13 Hello World!

```

Пройдясь по приведенному выше коду, можно отметить следующие моменты.

- В строках 5 и 6 мы создаем в таблице строку, устанавливаем столбец CLOB в EMPTY_CLOB() и извлекаем его значение в одном вызове. Все объекты LOB кроме временных находятся в базе данных — мы не можем записывать в переменную LOB, не имея указателя либо на временный объект LOB, либо на объект LOB, который существует в базе данных. EMPTY_CLOB() — это не NULL CLOB; это допустимый отличающийся от NULL указатель на пустую структуру. Кроме того, здесь получается локатор LOB, указывающий на данные в заблокированной строке. Если бы мы должны были извлечь это значение без блокировки лежащей в основе строки, то попытка записи потерпела бы неудачу, потому что перед записью объекты LOB должны быть заблокированы (в отличие от других структурированных данных). Вставляя строку, мы ее, конечно же, блокируем. Если бы мы модифицировали существующую строку вместо вставки новой, то должны были использовать оператор SELECT FOR UPDATE для извлечения и блокировки строки.
- В строке 8 мы создаем объект BFILE. Обратите внимание, что имя DIR1 представлено в верхнем регистре — как вскоре можно будет убедиться, это важно. Причина в том, что мы передаем BFILENAME() имя объекта, а не сам объект. Следовательно, мы должны гарантировать, что имя совпадает по регистру символов с именем, которое для этого объекта было сохранено Oracle.
- В строке 9 мы открываем объект LOB. Это позволит читать его.
- В строках 11 и 12 мы загружаем все содержимое файла операционной системы /tmp/test.txt в локатор LOB, который только что был вставлен. Мы применяем DBMS_LOB.GETLENGTH(), чтобы сообщить процедуре LOADFROMFILE() количество байтов BFILE, подлежащих загрузке (все байты).
- Наконец, в строке 14 мы закрываем открытый ранее файл BFILE, и столбец CLOB загружен.

Если в предыдущем примере вместо DIR1 указать dir1, будет выдано следующее сообщение об ошибке:

```
EODA@ORA12CR1> declare
```

```
...
6     returning theclob into l_clob;
7
8     l_bfile := bfilename( 'dir1', 'test.txt' );
9     dbms_lob.fileopen( l_bfile );
...
15 end;
16 /
declare
*
```

```
ERROR at line 1:
```

```
ORA-22285: non-existent directory or file for FILEOPEN operation
```

```
ORA-06512: at "SYS.DBMS_LOB", line 523
```

```
ORA-06512: at line 9
```

```
ОШИБКА в строке 1:
```

```
ORA-22285: несуществующий каталог или файл для операции FILEOPEN
```

```
ORA-06512: в SYS.DBMS_LOB, строка 523
```

```
ORA-06512: в строке 9
```

Дело в том, что каталог `dir1` не существует, а есть только `DIR1`. Если вы предпочитаете использовать имена каталогов в смешанном регистре, то при их создании должны указывать идентификаторы в кавычках, как это делалось для `dir2`. Это позволит написать такой код:

```
EODA@ORA12CR1> declare
```

```
2     l_clob    clob;
3     l_bfile   bfile;
4 begin
5     insert into demo values ( 1, empty_clob() )
6     returning theclob into l_clob;
7
8     l_bfile := bfilename( 'dir2', 'test.txt' );
9     dbms_lob.fileopen( l_bfile );
10
11     dbms_lob.loadfromfile( l_clob, l_bfile,
12                           dbms_lob.getlength( l_bfile ) );
13
14     dbms_lob.fileclose( l_bfile );
15 end;
16 /
```

```
PL/SQL procedure successfully completed.
```

```
Процедура PL/SQL успешно завершена.
```

Ранее в этом разделе упоминалось, что процедура `LoadCLOBFromFile` позволяет сообщить базе данных о том, что загружаемый файл имеет символьный набор, отличающийся от применяемого в базе данных, поэтому должно быть выполнено необходимое преобразование символьного набора. Если в результате запуска предыдущих примеров вы получаете вывод следующего вида:

```
DBMS_LOB.GETLENGTH(THECLOB) THECLOB
```

```
-----
6 ??????
```

то наиболее вероятно, что символьные наборы базы данных и файла не совпадают. В показанном ниже примере используется процедура LoadCLOBFromFile, чтобы учесть тот факт, что в файле применяется кодировка WE8ISO8859P1:

```

EODA@ORA12CR1> declare
  2  l_clob  clob;
  3  l_bfile bfile;
  4  dest_offset integer := 1;
  5  src_offset integer := 1;
  6  src_csid number := NLS_CHARSET_ID('WE8ISO8859P1');
  7  lang_context integer := dbms_lob.default_lang_ctx;
  8  warning integer;
  9  begin
 10  insert into demo values ( 1, empty_clob() )
 11  returning theclob into l_clob;
 12  l_bfile := bfilename( 'dir2', 'test.txt' );
 13  dbms_lob.fileopen( l_bfile );
 14  dbms_lob.loadclobfromfile( l_clob, l_bfile,
    dbms_lob.getlength( l_bfile ), dest_offset, src_offset,
    src_csid, lang_context, warning );
 15  dbms_lob.fileclose( l_bfile );
 16  end;
 17  /

```

PL/SQL procedure successfully completed.

Процедура PL/SQL успешно завершена.

Теперь выборка данных из таблицы дает следующий вывод:

```

EODA@ORA12CR1> select dbms_lob.getlength(theclob), theclob from demo;
DBMS_LOB.GETLENGTH(THECLOB) THECLOB
-----
13 Hello World!

```

Помимо процедур загрузки из файлов существуют и другие методы, посредством которых можно наполнить объект LOB данными с использованием PL/SQL. Работать с пакетом DBMS_LOB и предлагаемыми в нем процедурами намного легче, если файл должен быть загружен целиком. Когда во время загрузки файла требуется обрабатывать его содержимое, для чтения данных можно также применять процедуру DBMS_LOB.READ в отношении файла BFILE.

Использовать процедуру UTL_RAW.CAST_TO_VARCHAR2 удобно, когда читается текст, а не данные RAW. Затем с помощью процедуры DBMS_LOB.WRITE или DBMS_LOB.WRITEAPPEND данные можно поместить в CLOB или BLOB.

Загрузка данных LOB с помощью SQLLDR

Теперь мы исследуем загрузку данных в объект LOB посредством SQLLDR. Для этого доступно множество методов, но мы рассмотрим две самых распространенных ситуации.

- Данные встроены вместе с другими данными.
- Данные хранятся отдельно, и входные данные содержат имя файла, подлежащего загрузке. В терминологии SQLLDR он также называется вторичным файлом данных (secondary data file — SDF).

Начнем со встроенных данных.

Загрузка встроенных данных LOB

Объекты LOB обычно будут включать символы новой строки и другие специальные символы. Следовательно, для загрузки этих данных вы почти всегда будете применять один из четырех методов, которые были описаны в разделе “Как загрузить данные со встроенными символами новой строки?”. Давайте начнем с модификации таблицы DEPT, изменив тип столбца COMMENTS с VARCHAR2 на CLOB:

```

EODA@ORA12CR1> truncate table dept;
Table truncated.
Таблица усечена.

EODA@ORA12CR1> alter table dept drop column comments;
Table altered.
Таблица изменена.

EODA@ORA12CR1> alter table dept add comments clob;
Table altered.
Таблица изменена.

```

Например, пусть имеется файл данных (demo.dat) со следующим содержимым:

```

10, Sales, Virginia, This is the Sales
Office in Virginia|
20, Accounting, Virginia, This is the Accounting
Office in Virginia|
30, Consulting, Virginia, This is the Consulting
Office in Virginia|
40, Finance, Virginia, "This is the Finance
Office in Virginia, it has embedded commas and is
much longer than the other comments field. If you
feel the need to add double quoted text in here like
this: ""You will need to double up those quotes!"" to
preserve them in the string. This field keeps going for up to
1000000 bytes (because of the control file definition I used)
or until we hit the magic end of record marker,
the | followed by an end of line - it is right here ->"|

```

Каждая запись заканчивается символом |, за которым идет маркер конца строки. Текст для отдела 40 намного длиннее остальных; в нем присутствует множество переносов строк, встроенных кавычек и запятых. Имея такой файл данных, создадим управляющий файл:

```

LOAD DATA
INFILE demo.dat "str X'7C0A'"
INTO TABLE DEPT
REPLACE
FIELDS TERMINATED BY ',' OPTIONALLY ENCLOSED BY '"'
TRAILING NULLCOLS
(DEPTNO,
  DNAME      "upper(:dname)",
  LOC        "upper(:loc)",
  COMMENTS   char(1000000)
)

```

На заметку! Этот пример ориентирован на среду UNIX/Linux, где маркер конца строки занимает 1 байт; отсюда и установка STR в приведенном выше управляющем файле. Для среды Windows установка STR имела бы вид '7C0D0A'.

Чтобы загрузить файл данных, мы указываем CHAR(1000000) для столбца COMMENTS, поскольку, как обсуждалось ранее, по умолчанию SQLLDR выбирает для любого входного поля тип CHAR(255). Тип CHAR(1000000) позволит SQLLDR обработать до 1 000 000 байтов входного текста. Вы должны указать здесь длину, которая превышает любую порцию текста, ожидаемую во входном файле. Просмотрев загруженные данные, мы увидим следующий вывод:

```
EODA@ORA12CR1> select comments from dept;
```

```
COMMENTS
```

```
-----  
This is the Consulting  
Office in Virginia
```

```
This is the Finance  
Office in Virginia, it has embedded commas and is  
much longer than the other comments field. If you  
feel the need to add double quoted text in here like  
this: "You will need to double up those quotes!" to  
preserve them in the string. This field keeps going for up to  
1000000 bytes or until we hit the magic end of record marker,  
the | followed by an end of line - it is right here ->
```

```
This is the Sales  
Office in Virginia
```

```
This is the Accounting  
Office in Virginia
```

Легко заметить, что ранее продублированные двойные кавычки больше не дублируются. Лишние кавычки были устранены инструментом SQLLDR.

Загрузка внешних данных LOB

Распространенный сценарий предусматривает наличие файла данных, который содержит имена файлов, предназначенных для загрузки в объекты LOB, вместо смешивания данных LOB со структурированными данными. Организация данных в подобной манере обеспечивает более высокую гибкость, т.к. в файле данных, передаваемом SQLLDR, не придется решать проблему существования встроенных символов новой строки с помощью одного из четырех методов, что часто случается с большими объемами текстовых или двоичных данных. В SQLLDR этот тип дополнительных файлов данных называется LOBFILE.

Инструмент SQLLDR также поддерживает загрузку структурированного файла данных, который указывает на другой одиночный файл данных. Мы можем сообщить SQLLDR, каким образом проводить разбор данных LOB из внешнего файла, чтобы каждая строка структурированных данных загружалась с порцией данных из этого файла. Я считаю этот режим довольно ограниченным в использовании (на самом деле мне до сих пор не удалось найти ему применение), и здесь он рассматриваться не будет. В рамках SQLLDR такие внешние файлы называют *составными вторичными файлами данных*.

Тип LOBFILE представляет собой относительно простой файл данных, направленный на облегчение загрузки объектов LOB. Характерной чертой, которая отличает файлы LOBFILE от основных файлов данных, является отсутствие концепции записи, поэтому *символы новой строки никогда не встретятся*. Данные в LOBFILE находятся в одном из перечисленных ниже форматов:

- поля фиксированной длины (например, нужно загружать байты 100–1000 из LOBFILE);
- поля с разделителями (чем-то ограниченные или во что-то заключенные);
- пары длина/значение, т.е. поля переменной длины.

Наиболее распространенными являются поля с разделителями — в сущности те, которые завершаются символом конца файла (EOF). Обычно есть каталог, заполненный файлами, которые необходимо загрузить в столбцы LOB, и каждый файл целиком поступит в BLOB. В этом случае вы будете использовать тип данных LOBFILE с конструкцией TERMINATED BY EOF.

Итак, предположим, что имеется каталог с файлами, подлежащими загрузке в базу данных. Мы хотели бы загрузить поля OWNER (владелец файла), TIME_STAMP (временная метка файла), NAME (имя файла) и содержимое самого файла. Таблица, в которую будет производиться загрузка, создается следующим образом:

```

EODA@ORA12CR1> create table lob_demo
2  ( owner          varchar2(255),
3    time_stamp     date,
4    filename       varchar2(255),
5    data           blob
6  )
7  /

```

Table created.

Таблица создана.

Захватив вывод простой команды `ls -l` в среде UNIX/Linux или `dir /q /n` в среде Windows, можно сгенерировать входной файл и загрузить его с применением управляющего файла вроде показанного ниже для случая UNIX/Linux:

```

LOAD DATA
INFILE *
REPLACE
INTO TABLE LOB_DEMO
( owner          position(14:19),
  time_stamp     position(31:42) date "Mon DD HH24:MI",
  filename       position(44:100),
  data           LOBFILE(filename) TERMINATED BY EOF
)
BEGINDATA
-rwxr-xr-x 1 oracle dba 14889 Jul 22 22:01 demo1.log_xt
-rwxr-xr-x 1 oracle dba  123 Jul 22 20:07 demo2.ctl
-rwxr-xr-x 1 oracle dba  712 Jul 23 12:11 demo.bad
-rwxr-xr-x 1 oracle dba  8136 Mar  9 12:36 demo.control_files
-rwxr-xr-x 1 oracle dba  825 Jul 23 12:26 demo.ctl
-rwxr-xr-x 1 oracle dba 1681 Jul 23 12:26 demo.log
-rw-r----- 1 oracle dba  118 Jul 23 12:52 dl.sql

```



```
-rwxr-xr-x 1 oracle dba 127 Jul 23 12:05 lob_demo.sql
-rwxr-xr-x 1 oracle dba 171 Mar 10 13:53 p.bsh
-rwxr-xr-x 1 oracle dba 327 Mar 10 11:10 prime.bsh
-rwxr-xr-x 1 oracle dba 24 Mar 6 12:09 run_df.sh
```

Вот как теперь выглядит содержимое таблицы LOB_DEMO после запуска SQLLDR:

```
EODA@ORA12CR1> select owner, time_stamp, filename, dbms_lob.getlength(data)
2 from lob_demo
3 /
```

| OWNER | TIME_STAM | FILENAME | DBMS_LOB.GETLENGTH(DATA) |
|--------|-----------|--------------------|--------------------------|
| oracle | 22-JUL-14 | demo1.log_xt | 14889 |
| oracle | 22-JUL-14 | demo2.ct1 | 123 |
| oracle | 23-JUL-14 | demo.bad | 712 |
| oracle | 09-MAR-14 | demo.control_files | 8136 |
| oracle | 23-JUL-14 | demo.ct1 | 825 |
| oracle | 23-JUL-14 | demo.log | 0 |
| oracle | 23-JUL-14 | dl.sql | 118 |
| oracle | 23-JUL-14 | lob_demo.sql | 127 |
| oracle | 10-MAR-14 | p.bsh | 171 |
| oracle | 10-MAR-14 | prime.bsh | 327 |
| oracle | 06-MAR-14 | run_df.sh | 24 |

11 rows selected.

На заметку! Может возникнуть вопрос: почему размер файла demo.log равен 0? На протяжении выполнения SQLLDR открывает файл demo.log для записи, что обнуляет его длину и сбрасывает этот файл. Таким образом, во время своей загрузки файл demo.log пуст.

Это работает с типами CLOB и BLOB. Загружать каталог текстовых файлов с использованием SQLLDR в подобной манере очень просто.

Загрузка данных LOB в объектные столбцы

Теперь, когда вы знаете, как загружать данные в созданную простую таблицу, может также возникнуть потребность в загрузке данных в таблицу, которая имеет сложный объектный тип с входящим в него LOB. Чаще всего это случается при работе с изображениями. Средства работы с изображениями реализованы с применением сложного объектного типа ORDSYS.ORDIMAGE. Нам нужна возможность сообщить SQLLDR, каким образом производить в него загрузку.

Для загрузки LOB в столбец типа ORDIMAGE необходимо чуть лучше понимать структуру типа ORDIMAGE. Запустив DESCRIBE в отношении таблицы, в которую будет осуществляться загрузка, в SQL*Plus, мы обнаружим столбец по имени IMAGE типа ORDSYS.ORDIMAGE, куда в конечном итоге должны быть загружены данные IMAGE.SOURCE.LOCALDATA. Приведенные далее примеры будут работать, только если в системе установлено и сконфигурировано средство Oracle Text; в противном случае тип данных ORDSYS.ORDIMAGE будет неизвестен.

```
EODA@ORA12CR1> create table image_load(
2 id number,
3 name varchar2(255),
4 image ordsys.ordimage
5 )
6 /
```

Table created.

```
EODA@ORA12CR1> desc image_load
```

| Name | Null? | Type |
|-------|-------|-----------------|
| ----- | ----- | ----- |
| ID | | NUMBER |
| NAME | | VARCHAR2 (255) |
| IMAGE | | ORDSYS.ORDIMAGE |

```
EODA@ORA12CR1> desc ordsys.ordimage
```

| Name | Null? | Type |
|---------------|-------|------------------|
| ----- | ----- | ----- |
| SOURCE | | ORDSYS.ORDSOURCE |
| HEIGHT | | NUMBER (38) |
| WIDTH | | NUMBER (38) |
| CONTENTLENGTH | | NUMBER (38) |
| FILEFORMAT | | VARCHAR2 (4000) |

```
...
```

```
EODA@ORA12CR1> desc ordsys.ordsource
```

| Name | Null? | Type |
|-------------|-------|-----------------|
| ----- | ----- | ----- |
| LOCALDATA | | BLOB |
| SRCTYPE | | VARCHAR2 (4000) |
| SRCLOCATION | | VARCHAR2 (4000) |
| SRCNAME | | VARCHAR2 (4000) |
| UPDATETIME | | DATE |

```
...
```

На заметку! Чтобы отобразить сразу всю иерархию, можно было бы выполнить в SQL*Plus команду SET DESC DEPTH ALL или SET DESC DEPTH <n>. Учитывая, что описание типа ORDSYS.ORDIMAGE займет много страниц, оно будет приводиться по частям.

Управляющий файл для загрузки этого может выглядеть следующим образом:

```
LOAD DATA
INFILE *
INTO TABLE image_load
REPLACE
FIELDS TERMINATED BY ','
( ID,
  NAME,
  file_name FILLER,
  IMAGE column object
(
  SOURCE column object
(
    LOCALDATA LOBFILE (file_name) TERMINATED BY EOF
    NULLIF file_name = 'NONE'
  )
)
)
)
BEGINDATA
1,icons,icons.gif
```

Здесь введены две новых конструкции.

- **COLUMN OBJECT.** Сообщает SQLLDR, что это не полное имя столбца, а часть имени столбца. Конструкция не сопоставляется с каким-то полем во входном файле, но служит для построения корректной ссылки на объектный столбец, используемой во время загрузки. В предыдущем файле есть два дескриптора объектных столбцов, один вложенный внутрь другого. Поэтому имя столбца, которое будет применяться, выглядит как `IMAGE.SOURCE.LOCALDATA` — так и должно быть. Обратите внимание, что мы не загружаем другие атрибуты этих двух объектных типов (скажем, `IMAGE.HEIGHT`, `IMAGE.CONTENTLENGTH` и `IMAGE.SOURCE.SRCTYPE`). Вскоре вы увидите, как их заполнять.
- **NULLIF FILE_NAME='NONE'.** Сообщает SQLLDR, что в случае, если `FILE_NAME` содержит слово `NONE`, то в объектный столбец должно быть загружено значение `NULL`.

После загрузки типа Oracle Text обычно требуется выполнить заключительную обработку загруженных данных с использованием PL/SQL, чтобы заставить средство Oracle Text работать с ним. Например, для предыдущих данных можно было бы запустить следующий код, позволяющий корректно настроить свойства изображения:

```
begin
  for c in ( select * from image_load ) loop
    c.image.setproperties;
  end loop;
end;
/
```

`SETPROPERTIES` — это объектный метод, предоставляемый типом `ORDSYS.ORDIMAGE`, который обрабатывает само изображение и обновляет остальные атрибуты объекта подходящими значениями.

Как вызвать SQLLDR из хранимой процедуры?

Если отвечать кратко, то сделать это невозможно. Инструмент SQLLDR — не API-интерфейс; он не является чем-то вызываемым. Это программа командной строки. Вы определенно можете написать на языке Java или C внешнюю процедуру, которая запустит SQLLDR, но это не будет тем же самым, что и “вызов” SQLLDR. Загрузка произойдет в другом сеансе и находится за пределами вашего контроля над транзакцией. Кроме того, понадобится проанализировать результирующий журнальный файл, чтобы определить, прошла ли загрузка успешно, и насколько успешно (т.е. сколько строк было загружено до того, как загрузка была прекращена из-за ошибки). Вызывать SQLLDR из хранимой процедуры я не рекомендую.

В прошлом, до выхода версии Oracle9i, можно было реализовывать собственный процесс, подобный SQLLDR. Ниже перечислены возможные варианты.

- Написание мини-загрузчика SQLLDR на PL/SQL. Он может применять либо типы `BFILE` для чтения двоичных данных, либо пакет `UTL_FILE` для чтения текстовых данных с целью последующего разбора и загрузки.
- Написание мини-загрузчика SQLLDR на Java. Это может быть чуть сложнее, чем реализация загрузчика на PL/SQL, и требовать использования многих доступных процедур Java.
- Написание загрузчика SQLLDR на C и его вызов в качестве внешней процедуры.

Мы завершим тему, связанную с SQLLDR, обсуждением нескольких аспектов, которые не являются интуитивно понятными.

Предостережения относительно SQLLDR

В этом разделе речь пойдет о ряде вещей, которых следует остерегаться при работе с SQLLDR.

Опция TRUNCATE работает иначе

Может показаться, что опция TRUNCATE в SQLLDR работает не так, как TRUNCATE в SQL*Plus или в любом другом инструменте. Исходя из предположения, что вы будете повторно загружать в таблицу похожий объем данных, SQLLDR применяет расширенную форму TRUNCATE, а именно — он выдает следующую команду:

```
truncate table t reuse storage
```

Конструкция REUSE STORAGE не освобождает выделенные экстеннты, а только помечает их как свободное пространство. Если это нежелательно, можете перед запуском SQLLDR очистить таблицу.

Стандартным типом входных полей в SQLLDR является CHAR (255)

С этой проблемой приходится сталкиваться настолько часто, что я решил дважды упомянуть о ней в главе. Стандартная длина входных полей составляет 255 символов. Если ваше поле длиннее 255 символов, вы получите следующее сообщение об ошибке:

```
Record N: Rejected - Error on table T, column C.
Field in data file exceeds maximum length
Запись N: отклонена - ошибка в таблице T, столбец C.
Поле в файле данных превышает максимальную длину
```

Это вовсе не означает, что данные не умещаются в столбец базы данных; сообщение указывает лишь на то, что загрузчик SQLLDR ожидал 255 или менее байтов входных данных, а получил больше. Решение предусматривает просто использование в управляющем файле типа CHAR(N), где N — достаточно большое значение, чтобы соответствовать самой большой длине поля во входном файле. Пример приводился в самом начале раздела “Часто задаваемые вопросы по загрузке данных посредством SQLLDR”.

Параметры командной строки переопределяют параметры управляющего файла

Многие параметры SQLLDR могут либо помещаться в управляющий файл, либо передаваться в командной строке. Скажем, можно применять INFILE FILENAME, а также SQLLDR...DATA=FILENAME. Командная строка переопределяет любые параметры в управляющем файле. Нельзя рассчитывать на то, что параметры в управляющем файле действительно будут использоваться, т.к. тот, кто запускает SQLLDR, вполне может переопределить их.

Заключительные соображения по поводу SQLLDR

В этом разделе мы исследовали многие области, связанные с загрузкой данных. Мы раскрыли типичные повседневные проблемы, с которыми приходится сталки-

ваться: загрузка файлов с разделителями, загрузка файлов записей фиксированной длины, загрузка каталога, заполненного файлами изображений, применение функций к входным данным для их преобразования и т.д. Мы не рассматривали подробно массивные загрузки данных с использованием загрузчика в прямом режиме, а только слегка затронули эту тему. Нашей целью были ответы на часто возникающие вопросы относительно применения SQLLDR, которые интересуют широкую аудиторию.

Выгрузка в плоский файл

Одной вещью, которую не делает SQLLDR, и для которой Oracle не предлагает инструменты командной строки, является выгрузка данных в формате, воспринимаемом SQLLDR или другими программами. Это было бы полезно для перемещения данных из системы в систему, не используя утилиты EXPDP/IMPDP технологии Data Pump. Применение EXPDP/IMPDP для перемещения данных из системы в систему хорошо работает для средних объемов данных при условии, что обе системы представляют собой Oracle.

На заметку! Средство APEX (Application Express) предоставляет возможность экспорта данных как часть SQL Workshop, как это делает и SQL Developer. Вы можете без труда экспортировать информацию в формате CSV. Это хорошо работает для нескольких мегабайтов данных, но совершенно не подходит для многих десятков мегабайтов и более.

Мы разработаем небольшую утилиту PL/SQL, которая может использоваться для выгрузки данных на сервере в дружественном к SQLLDR формате. Эквивалентные инструменты для решения той же задачи на Pro*C и SQL*Plus доступны по адресу <http://tkyte.blogspot.com/2009/10/httpasktomoraclecomtkyteflat.html>. Утилита PL/SQL будет эффективно работать в большинстве ситуаций с малым объемом данных, но с применением Pro*C можно достичь лучшей производительности. Обратите внимание, что Pro*C и SQL*Plus также удобны, если необходимо генерировать файлы на стороне клиента, а не на сервере, где их будет создавать PL/SQL.

Ниже приведена спецификация разрабатываемого пакета:

```

EODA@ORA12CR1> create or replace package unloader
2  AUTHID CURRENT_USER
3  as
4  /* Функция run -- выгружает данные из любого запроса в файл
5      и создает управляющий файл для повторной
6      загрузки данных в другую таблицу.
7
8      p_query      = SQL-запрос для выгрузки. Может быть практически любым
                      запросом.
9      p_tname      = таблица для загрузки. Будет помещена в управляющий файл.
10     p_mode       = REPLACE|APPEND|TRUNCATE -- способ повторной загрузки
                      данных.
11     p_dir        = каталог, в который будут записываться файлы ctl и dat.
12     p_filename   = имя файла для записи. К этому имени будут
13                     добавляться .ctl and .dat.
14     p_separator  = разделитель полей. По умолчанию используется запятая.
15     p_enclosure  = внутри чего будет заключено каждое поле.
```

```

16  p_terminator = символ конца строки.
    Мы применяем его для того, чтобы можно было
17  выгружать и повторно загружать данные, содержание символы новой строки.
18  По умолчанию используется "\n" (UNIX/Linux) и "\r\n" (Windows).
19  Вы должны переопределять это, только если уверены, что данные будут
20  содержать эту стандартную последовательность символов.
21  Я ВСЕГДА добавляю к этой последовательности маркер конца строки
22  операционной системы, но вам поступать так не обязательно. */
23  function run( p_query      in varchar2,
24              p_tname      in varchar2,
25              p_mode       in varchar2 default 'REPLACE',
26              p_dir       in varchar2,
27              p_filename  in varchar2,
28              p_separator in varchar2 default ',',
29              p_enclosure in varchar2 default '"',
30              p_terminator in varchar2 default '|' )
31      return number;
32  end;
33  /
Package created.
Пакет создан.

```

Обратите внимание на конструкцию AUTHID CURRENT_USER. Она позволяет этому пакету однажды установиться в базе данных и затем применяться любым пользователем для выгрузки данных. Все, что пользователю понадобится иметь — это привилегию SELECT на выгружаемой таблице (таблицах) и привилегию EXECUTE на самом пакете.

Если бы в этом случае не использовалась конструкция AUTHID CURRENT_USER, то владельцу пакета понадобились бы прямые привилегии SELECT на всех таблицах, подлежащих выгрузке. Кроме того, поскольку эта процедура вполне могла бы стать целью атак внедрением SQL (ведь она принимает на входе произвольный оператор SQL), оператор CREATE должен содержать AUTHID CURRENT_USER.

Если бы это была процедура со стандартными правами владельца, то любой имеющий привилегию EXECUTE на ней, располагал бы возможностью выполнения *любого* оператора SQL с применением прав владельца. Если известно, что процедура уязвима к атакам внедрением SQL, то лучше пусть она выполняется с правами вызывающего!

На заметку! Процедура SQL будет выполняться с привилегиями вызывающего. Тем не менее, все процедуры PL/SQL будут запускаться с привилегиями *владельца* вызываемой процедуры; следовательно, возможность использования пакета UTL_FILE для записи в каталог неявно предоставляется каждому, кто имеет привилегию EXECUTE на этом пакете.

Тело пакета показано ниже. Для записи управляющего файла и файла данных применяется пакет UTL_FILE. Для динамической обработки запроса используется пакет DBMS_SQL. В запросах присутствует один тип данных: VARCHAR2(4000). Это подразумевает, что мы не можем применять данный метод для выгрузки столбцов LOB, и это действительно так, если столбец LOB имеет длину больше 4000 байтов. Однако указанный тип можно использовать для выгрузки вплоть до 4000 байтов любого столбца LOB совместно с функцией SUBSTR. К тому же, поскольку VARCHAR2

является единственным выходным типом данных, можно обрабатывать столбцы RAW длиной до 2000 байт (4000 шестнадцатеричных символов), чего вполне достаточно для всех типов кроме LONG RAW и LOB. Вдобавок любой запрос, который ссылается на нескаллярный атрибут (сложный объектный тип, вложенная таблица и т.д.), не будет работать с этой простой реализацией. Далее представлено девяностопроцентное решение — в том смысле, что оно решает проблему на протяжении 90% времени:

```
EODA@ORA12CR1> create or replace package body unloader
2 as
3
4
5 g_theCursor      integer default dbms_sql.open_cursor;
6 g_descTbl       dbms_sql.desc_tab;
7 g_nl            varchar2(2) default chr(10);
8
```

В теле этого пакета применяется несколько глобальных переменных. Глобальный курсор открывается один раз, при первом обращении к пакету, и остается открытым вплоть до завершения сеанса. Это позволяет избежать накладных расходов на получение нового курсора при каждом вызове пакета. Глобальная переменная G_DESCDTBL — это таблица PL/SQL, которая будет содержать вывод процедуры DBMS_SQL.DESCRIBE, а G_NL — символ новой строки. Он используется в строках, которые должны содержать внутри себя символы новой строки. Код не придется подстраивать для Windows — пакет UTL_FILE будет видеть CHR(10) в строке символов и автоматически превращать его в сочетание “возврат каретки/перевод строки”.

Далее следует небольшая вспомогательная функция, служащая для преобразования символа в шестнадцатеричный формат. Для этого она применяет встроенные функции:

```
9
10 function to_hex( p_str in varchar2 ) return varchar2
11 is
12 begin
13     return to_char( ascii(p_str), 'fm0x' );
14 end;
15
```

Наконец, мы создаем еще одну вспомогательную функцию IS_WINDOWS, которая возвращает TRUE или FALSE в зависимости от того, работаем ли мы на платформе Windows, поэтому конец строки является двухсимвольной строкой, а не одиночным символом, как в большинстве других платформ. Мы используем встроенную функцию GET_PARAMETER_VALUE пакета DBMS_UTILITY, которая может применяться для чтения почти любого параметра. Мы извлекаем значение параметра CONTROL_FILES и выясняем, присутствует ли в нем символ \ — если он обнаруживается, то мы работаем в среде Windows.

```
16 function is_windows return boolean
17 is
18     l_cfiles varchar2(4000);
19     l_dummy  number;
20 begin
```

```

21  if (dbms_utility.get_parameter_value( 'control_files', l_dummy,
    l_cfiles )>0)
22  then
23      return instr( l_cfiles, '\' ) > 0;
24  else
25      return FALSE;
26  end if;
27  end;

```

На заметку! Функция IS_WINDOWS полагается на использование символа \ в параметре CONTROL_FILES. Имейте в виду, что вместо него можно применять /, но это будет весьма необычно.

Ниже представлен код процедуры создания управляющего файла для повторной загрузки выгруженных данных с использованием таблицы описаний столбцов, сгенерированной DBMS_SQL.DESCRIBE_COLUMNS. Она позаботится о специфике операционной системы, такой как применение символов возврата каретки и перевода строки (для атрибута STR).

```

28
29  procedure dump_ctl( p_dir          in varchar2,
30                    p_filename      in varchar2,
31                    p_tname         in varchar2,
32                    p_mode          in varchar2,
33                    p_separator     in varchar2,
34                    p_enclosure     in varchar2,
35                    p_terminator    in varchar2 )
36  is
37      l_output      utl_file.file_type;
38      l_sep         varchar2(5);
39      l_str         varchar2(5) := chr(10);
40
41  begin
42      if ( is_windows )
43      then
44          l_str := chr(13) || chr(10);
45      end if;
46
47      l_output := utl_file.fopen( p_dir, p_filename || '.ctl', 'w' );
48
49      utl_file.put_line( l_output, 'load data' );
50      utl_file.put_line( l_output, 'infile ' ||
51                        p_filename || '.dat' || "str x'" ||
52                        utl_raw.cast_to_raw( p_terminator ||
53                        l_str ) || "'" );
54      utl_file.put_line( l_output, 'into table ' || p_tname );
55      utl_file.put_line( l_output, p_mode );
56      utl_file.put_line( l_output, 'fields terminated by X' ||
57                        to_hex(p_separator) ||
58                        "' enclosed by X'" ||
59                        to_hex(p_enclosure) || ' ' );
60      utl_file.put_line( l_output, '(' );
61

```



```

62     for i in 1 .. g_descTbl.count
63     loop
64         if ( g_descTbl(i).col_type = 12 )
65             then
66                 utl_file.put( l_output, l_sep || g_descTbl(i).col_name ||
67                             ' date ' 'ddmmyyyyhh24miss' ' ');
68             else
69                 utl_file.put( l_output, l_sep || g_descTbl(i).col_name ||
70                             ' char(' ' ||
71                             to_char(g_descTbl(i).col_max_len*2) || ' ' ) ' );
72             end if;
73             l_sep := ', ' || g_nl ;
74         end loop;
75         utl_file.put_line( l_output, g_nl || ' ' );
76         utl_file.fclose( l_output );
77     end;
78

```

Далее показан код простой функции для возвращения строки в кавычках с использованием выбранного символа. Обратите внимание, что она не только заключает строку в пару указанных символов, но также дублирует эти символы, если они встречаются внутри строки, чтобы предохранить их:

```

79 function quote(p_str in varchar2, p_enclosure in varchar2)
80     return varchar2
81 is
82 begin
83     return p_enclosure ||
84         replace( p_str, p_enclosure, p_enclosure||p_enclosure ) ||
85         p_enclosure;
86 end;
87

```

Затем следует главная функция — RUN. Из-за того, что она достаточно велика, то будет комментироваться по частям:

```

88 function run( p_query           in varchar2,
89              p_tname           in varchar2,
90              p_mode             in varchar2 default 'REPLACE',
91              p_dir             in varchar2,
92              p_filename        in varchar2,
93              p_separator       in varchar2 default ',',
94              p_enclosure       in varchar2 default '"',
95              p_terminator     in varchar2 default '|' ) return number
96 is
97     l_output           utl_file.file_type;
98     l_columnValue     varchar2(4000);
99     l_colCnt           number default 0;
100    l_separator        varchar2(10) default '';
101    l_cnt              number default 0;
102    l_line             long;
103    l_datefmt          varchar2(255);
104    l_descTbl          dbms_sql.desc_tab;
105 begin

```

Мы сохраним значение NLS_DATE_FORMAT в переменной, чтобы иметь возможность при выгрузке данных на диск изменить его на формат, предохраняющий дату и время. В подобной манере мы сохраняем компонент времени внутри даты. Затем мы настраиваем блок исключения, чтобы можно было сбросить NLS_DATE_FORMAT при возникновении любой ошибки.

```

106      select value
107      into l_datefmt
108      from nls_session_parameters
109      where parameter = 'NLS_DATE_FORMAT';
110
111      /*
112      Установка формата даты в виде большого числа. Это позволяет
113      избежать всех проблем, связанных с NLS, и сохранить время и дату.
114      */
115      execute immediate
116      'alter session set nls_date_format='''ddmmyyyyhh24miss''' ';
117
118      /*
119      Настройка блока исключения, чтобы в случае любой ошибки
120      можно было бы, по крайней мере, сбросить формат данных.
121      */
122      begin

```

Далее мы производим разбор и получаем описание запроса. Установка G_DESC_TBL в L_DESC_TBL сделана для сброса глобальной таблицы; в противном случае в дополнение к данным по текущему запросу она может содержать описание предыдущего запроса. После этого мы вызываем процедуру DUMP_CTL для действительного создания управляющего файла.

```

123      /*
124      Разбор и описание запроса. Мы сбрасываем descTbl
125      в пустую таблицу, так что вызов .count на ней
126      будет надежным.
127      */
128      dbms_sql.parse( g_theCursor, p_query, dbms_sql.native );
129      g_descTbl := l_descTbl;
130      dbms_sql.describe_columns( g_theCursor, l_colCnt, g_descTbl );
131
132      /*
133      Создание управляющего файла для повторной загрузки
134      этих данных в желаемую таблицу.
135      */
136      dump_ctl( p_dir, p_filename, p_tname, p_mode, p_separator,
137               p_enclosure, p_terminator );
138
139      /*
140      Привязка каждого отдельного столбца к varchar2(4000).
141      Нас не заботит, что извлекается - число, дата или что-то другое.
142      Все это может быть строкой.
143      */

```

Мы готовы к сбросу действительных данных на диск. Начнем с того, что определим каждый столбец для выборки в него данных как относящийся к типу

VARCHAR2(4000). Типы NUMBER, DATE, RAW и т.д. будут преобразованы в VARCHAR2. Сразу же после этого мы выполним запрос для подготовки к фазе выборки:

```

144     for i in 1 .. l_colCnt loop
145         dbms_sql.define_column( g_theCursor, i, l_columnValue, 4000);
146     end loop;
147
148     /*
149         Запуск запроса с игнорированием вывода,
            получаемого в результате выполнения.
150         Допустимо только когда DML-оператором
            является insert/update или delete.
151     */

```

Откроем файл данных для записи, извлечем все строки из запроса и выведем их в этот файл данных:

```

152     l_cnt := dbms_sql.execute(g_theCursor);
153
154     /*
155         Открытие файла для записи вывода и затем запись
156         в него данных с разделителями.
157     */
158     l_output := utl_file.fopen( p_dir, p_filename || '.dat', 'w',
159                               32760 );
160     loop
161         exit when ( dbms_sql.fetch_rows(g_theCursor) <= 0 );
162         l_separator := '';
163         l_line := null;
164         for i in 1 .. l_colCnt loop
165             dbms_sql.column_value( g_theCursor, i,
166                                   l_columnValue );
167             l_line := l_line || l_separator ||
168                       quote( l_columnValue, p_enclosure );
169             l_separator := p_separator;
170         end loop;
171         l_line := l_line || p_terminator;
172         utl_file.put_line( l_output, l_line );
173         l_cnt := l_cnt+1;
174     end loop;
175     utl_file.fclose( l_output );
176

```

Наконец, установим формат данных в старое состояние (и блок исключения сделает то же самое, если по какой-то причине в любой части предыдущего кода возникнет ошибка) и возвратим управление:

```

177     /*
178         Сброс формата данных и возвращение количества строк,
179         записанных в выходной файл.
180     */
181     execute immediate
182         'alter session set nls_date_format='' || l_datefmt || ''';
183     return l_cnt;
184 exception

```

```

185      /*
186      В случае возникновения ЛЮБОЙ ошибки сброс формата
187      данных и повторная генерация ошибки.
188      */
189      when others then
190          execute immediate
191          'alter session set nls_date_format='' || l_datefmt || ''';
192          RAISE;
193      end;
194  end run;
195
196
197  end unloader;
198  /
Package body created.
Тело пакета создано.

```

Чтобы запустить этот пакет, можно применить следующие команды (обратите внимание, что это требует выдачи привилегии SELECT на таблице SCOTT.EMP одной из ваших ролей или непосредственно вам):

```

EODA@ORA12CR1> set serveroutput on
EODA@ORA12CR1> create or replace directory my_dir as '/tmp';
Directory created.
Каталог создан.

EODA@ORA12CR1> declare
2     l_rows    number;
3  begin
4     l_rows := unloader.run
5         ( p_query      => 'select * from scott.emp order by empno',
6           p_tname       => 'emp',
7           p_mode        => 'replace',
8           p_dir         => 'MY_DIR',
9           p_filename    => 'emp',
10          p_separator   => ',',
11          p_enclosure    => '»',
12          p_terminator  => '~' );
13
14     dbms_output.put_line( to_char(l_rows) ||
15                          ' rows extracted to ascii file' );
16 end;
17 /
14 rows extracted to ascii file
14 строк извлечено в файл ascii
PL/SQL procedure successfully completed.
Процедура PL/SQL успешно завершена.

```

Ниже показан управляющий файл, который был сгенерирован этой процедурой PL/SQL (выделенные полужирным числа справа в скобках на самом деле в файле не содержатся — они предназначены только для дальнейших ссылок):

```

load data                                     (1)
infile 'emp.dat' "str x'7E0A'"              (2)
into table emp                               (3)
replace                                     (4)

```

```

fields terminated by X'2c' enclosed by X'22'      (5)
(                                                  (6)
EMPNO char(44),                                  (7)
ENAME char(20),                                  (8)
JOB char(18),                                    (9)
MGR char(44),                                    (10)
HIREDATE date 'ddmmyyyyhh24miss',                (11)
SAL char(44),                                    (12)
COMM char(44),                                    (13)
DEPTNO char(44),                                  (14)
)                                                  (15)

```

Далее описаны важные части этого управляющего файла.

- **Строка (2).** Мы используем средство STR загрузчика SQLLDR. С его помощью можно указать символ или строку, используемую для завершения записи. Это позволит легко загружать данные, которые внутри содержат символы новой строки. Строка x'7E0A' содержит символы тильды и новой строки.
- **Строка (5).** Мы применяем наш символ разделителя и окружения. Конструкция OPTIONALLY ENCLOSED BY не используется, т.к. мы будем заключать каждое отдельное поле в пару символов окружения после дублирования любых вхождений этого символа в необработанных данных.
- **Строка (11).** Мы применяем формат даты в виде большого числа. Это позволяет избежать возникновения любых проблем с NLS в данных и предохраняет компонент времени в поле даты.

Файл необработанных данных (.dat), сгенерированный предыдущим кодом, выглядит так:

```

"7369","SMITH","CLERK","7902","17121980000000","800","","20"~
"7499","ALLEN","SALESMAN","7698","20021981000000","1600","300","30"~
"7521","WARD","SALESMAN","7698","22021981000000","1250","500","30"~
"7566","JONES","MANAGER","7839","02041981000000","2975","","20"~
"7654","MARTIN","SALESMAN","7698","28091981000000","1250","1400","30"~
"7698","BLAKE","MANAGER","7839","01051981000000","2850","","30"~
"7782","CLARK","MANAGER","7839","09061981000000","2450","","10"~
"7788","SCOTT","ANALYST","7566","19041987000000","3000","","20"~
"7839","KING","PRESIDENT","","17111981000000","5000","","10"~
"7844","TURNER","SALESMAN","7698","08091981000000","1500","0","30"~
"7876","ADAMS","CLERK","7788","23051987000000","1100","","20"~
"7900","JAMES","CLERK","7698","03121981000000","950","","30"~
"7902","FORD","ANALYST","7566","03121981000000","3000","","20"~
"7934","MILLER","CLERK","7782","23011982000000","1300","","10"~

```

С файлом .dat связаны следующие моменты:

- каждое поле заключено в символы окружения;
- даты выгружены в виде больших чисел;
- каждая строка данных в файле заканчивается символом ~, как и требовалось.

Теперь эти данные можно легко загрузить повторно, используя SQLLDR. К командной строке SQLLDR можно добавить любые необходимые параметры.

Как было указано ранее, логика пакета выгрузки может быть реализована с помощью широкого многообразия языков и инструментов. На моем веб-сайте вы найдете реализацию этого примера не только на PL/SQL, но также на Pro*C и в виде сценариев SQL*Plus. Самую быструю реализацию обеспечивает Pro*C, и она всегда выполняет запись в файловую систему клиентской рабочей станции. Язык PL/SQL дает хорошее универсальное решение (не требующее компиляции и установки на клиентских рабочих станциях), однако оно всегда производит запись в файловую систему сервера. Решение в виде сценариев SQL*Plus является удачным компромиссом, предлагая приемлемую производительность и возможность записи в файловую систему клиента.

Резюме

В этой главе были раскрыты многие детали загрузки и выгрузки данных. Сначала мы обсудили преимущества внешних таблиц по сравнению с SQLLDR. Затем мы взглянули на простые приемы, помогающие начать работу с внешними таблицами. Также были показаны примеры применения конструкции PREPROCESSOR для выполнения команд операционной системы перед загрузкой данных.

Далее мы рассмотрели средство Oracle 10g и последующих версий — выгрузку во внешние таблицы, а также возможность создания и перемещения фрагментов данных из одной базы данных в другую. Кроме того, мы исследовали способ выгрузки данных из таблицы в файл дампа, который может использоваться для переноса данных между базами данных.

Мы показали, что в большинстве сценариев вы должны применять внешние таблицы, а не инструмент SQLLDR. Тем не менее, некоторые ситуации могут требовать использования SQLLDR, например, при загрузке данных по сети. Мы представили множество базовых приемов загрузки данных с разделителями, данных фиксированной ширины, объектов LOB и тому подобного.

Наконец, мы взглянули на обратный процесс — выгрузку данных, и показали, как выводить данные из базы в формате, который может восприниматься другими инструментами, такими как электронные таблицы и т.д. В ходе обсуждения для демонстрации процесса мы разработали утилиту PL/SQL, которая выгружает данные в дружественном к SQLLDR (или внешним таблицам) формате, но может быть легко приведена в соответствие с имеющимися потребностями.

Предметный указатель

A

ADR (Automatic Diagnostic Repository), 156
ADRCI (Automatic Diagnostic Repository
Command Interpreter), 163
AMM (Automatic memory management), 205
AQ (Advanced Queuing), 277
ASH (Active Session History), 156
ASSM (Automatic Segment Space
Management), 493
ASM (Automatic Storage Management), 173
ASMB (Automatic Storage Management
Background), 290
ASMM (Automatic Shared Memory
Management), 205
AWR (Automatic Workload Repository), 156;
259; 294

C

CLOB (Character Large Object), 492
CSV (Comma-Separated Values), 917

D

Data Pump, 881
DBRM (Database Resource Manager), 289
DDL (Data Definition Language), 309; 327; 398
DML (Data Manipulating Language), 326
DRCP (Database Resident Connection
Pooling), 264; 276
DSS (Decision Support System), 188

E

ETL (extract, transform, load), 325

F

FBDA (Flashback Data Archiver), 289
FGAC (Fine-Grained Access Control), 106
FIFO (First In, First Out), 488

H

HIPAA (Health Insurance Portability and
Accountability Act), 834
HWM (High-Water Mark), 432; 494

I

IOT (Index Organized Table), 487

L

LDAP (Lightweight Directory Access Protocol), 131
LMT (Locally-Managed Tablespace), 861
LOB (Large Object), 727
LRU (Least Recently Used), 240

M

MMNL (Manageability Monitor Light), 292
MSSM (Manual Segment Space Management), 493

N

NIST (National Institute of Standards and
Technology), 90
NLS (National Language Support), 671

O

OLTP (On-Line Transaction Processing), 271
Oracle
 подключение к Oracle, 125
 через TCP/IP, 129
Oracle Automatic Storage Management Cluster
 File System (ACFS), 174
Oracle Cluster File System (OCFS), 174
Oracle Enterprise Manager (OEM), 145
Oracle Exadata, 842
Oracle Net (Сеть Oracle), 261
Oracle Parallel Server (OPS), 291
OSD (Operation System Dependent), 98

P

PDML (Parallel DML), 758; 838
PGA (Process Global Area), 114; 205
PID (Process ID), 262

Q

QCSID (Query Coordinator SID), 847

R

RAC (Real Application Clusters), 116

S

SAN (Storage Area Network), 842
SGA (System Global Area), 65; 114; 123; 205; 230
SPA (SQL Performance Analyzer), 156
SQL*Loader (SQL*Loader), 881; 883

T

Two Phase Commit (2PC), 285; 390; 418

U

UGA (User Global Area), 114; 205
Unicode (Unicode Consortium), 671

W

WLM (Workload Manager), 113

А

Автоматическая фиксация, 417
 Автоматическое управление памятью (ASM), 173
 Архитектура
 Oracle, 59
 двухпроцессная (двухзадачная), 261
 Атомарность
 на уровне оператора, 392
 на уровне процедуры, 394
 на уровне транзакции, 398
 Атрибут
 FIX, 927
 STR, 930
 VAR, 929
 Аудит, 834

Б

База данных, 114
 кеш буфера базы данных, 187
 контейнерная (container), 115; 122
 корневая (root), 115
 подключаемая, 132; 154
 подключаемая (pluggable), 115; 132; 154
 с единственным владельцем (single tenant), 114
 таблицы базы данных, 487
 Блок, 175
 отложенная очистка блоков, 481
 структура блока, 177
 Блокирование (blocking), 316
 обновления, 319
 ресурса, 77
 слияния, 319
 удаления, 319
 Блокировка (lock), 74; 299
 AE, 338
 DDL, 327; 339
 монопольная, 340
 разделяемая, 340
 DML, 326
 TM (помещение в очередь DML), 337
 TX (транзакций), 327
 взаимоблокировка, 87; 272
 внутренняя, 62
 вручную, 356
 оптимистическая, 305; 307
 выбор между оптимистической и пессимистической блокировкой, 315
 с использованием столбца версии, 308
 пессимистическая, 305
 преобразование блокировки (lock conversion), 326
 прямой загрузки (direct load), 341

разделяемого чтения, 83
 создание собственных блокировок, 357
 эскалация блокировок, 326
 Буфер повторения (redo buffer), 237

В

Взаимоблокировка (deadlock), 87; 320
 искусственная, 272
 Внедрение SQL, 66
 Вставка
 заблокированная, 316
 Выборка
 динамическая (dynamic sampling), 568

Г

Гранула (granule), 234

Д

Данные
 выгрузка, 881
 Data Pump, 881
 в плоский файл, 881; 944
 загрузка, 881
 дат, 921
 объекта LOB, 932
 с использованием функций, 922
 со встроенными символами новой строки, 926
 с разделителями, 916
 фиксированного формата, 919
 поиск по кодам, 882
 скользящее окно данных, 803
 слияние (MERGE), 882
 с низкой кардинальностью, 619
 хранилища данных, 803
 Диспетчер, 128
 OEM (Oracle Enterprise Manager), 145
 памяти (MMAN), 292
 рабочей загрузки (WLM), 113
 ресурсов базы данных (DBRM), 155; 273; 289

Ж

Журнал
 повторения транзакций, 189
 ретроспективный, 196

З

Запись
 отмены (undo), 84
 согласованность записи, 378
 Запрос
 координатор параллельного запроса, 297
 параллельный, 843
 производительность запросов, 758
 скалярные подзапросы, 101

Зашелка (latch), См. Блокировка, внутренняя, 62 62; 65; 344

И

Идентификатор процесса (PID), 262

Индекс, 804

глобальный, 794; 799; 801; 803; 813

асинхронное обслуживание глобальных индексов, 811

недействительный, 803

односекционный, 814

оперативное обслуживание глобальных индексов, 808

локальный, 795; 799

частичный, 818

Индекс (index), 490

секция индекса (index partition), 491

Индексация, 587

некоторых строк, 638

расширенных столбцов, 646

расширяемая, 642

Индекс-таблица (IOT), 509; 588

кластеризованная, 527

Индексы

битовые (bitmap index), 589; 618

соединений (bitmap join index), 589; 624

быстрое полное сканирование, 607

на основе функций (function-based index), 589; 627; 642

невидимые, 644

по реверсированным ключам (reverse key index), 588

предметной области (application domain index), 590; 642

сбалансированные по высоте, 592

сжатие ключей индекса, 593

со структурой В-дерева (B*Tree index), 588; 617

текстовые, 642

упорядоченные по убыванию (descending index), 588; 603

Инструмент

ADRCI, 163

SQLLDR, 883; 910; 942

К

Кеш

буфера базы данных, 187

буферов блоков, 238

выгрузка из кеша (ping), 291

Кластер (cluster), 487; 490

индексный, 528

однотабличный хеш-кластер, 543

Ключ секционирования (partition key), 750

Код

OSD, 98

Команда

AUDIT, 155

disconnect, 267

exit, 267

FLASHBACK DATABASE, 196

ipcs, 117; 120

ps, 116

pslist, 124

pstat, 124

pwd, 116

Копия

моментальная (snapshot), 292

Куча (heap), 505

Л

Логические идентификаторы строк, 526

Локатор LOB, 731

М

Маркер максимального уровня заполнения (HWM), 494

Масштабирование

вертикальное, 290

горизонтальное, 290

Метод

DRCP, 264

Easy Connect (легкое подключение), 131

черного ящика, 48

Многоверсионность, 78

Монитор

системный, 283

Мониторинг файловой системы посредством SQL, 898

О

Область памяти

PGA, 206

автоматическое управление памятью PGA, 215; 228

ручное управление памятью PGA, 208

SGA, 206; 230; 233

фиксированная (fixed SGA), 236

UGA, 206; 230

быстрого восстановления (Fast Recovery Area), 197

гранула (granule), 234

рабочая, 216

системная глобальная, 123

Обновление

потерянное (lost update), 303; 369

предотвращение потерянных обновлений, 75

Отмена (undo), 427
 Очередь заданий, 292
 Очистка блока
 отложенная, 481
 Ошибки
 обработка ошибок, 893
 проблемы многопользовательского
 доступа, 896

П

Память
 PGA (Глобальная область процесса), 205
 автоматическое управление памятью
 PGA, 215; 228
 ручное управление памятью PGA, 208
 SGA, 230; 233
 фиксированная (fixed SGA), 236
 SGA (Глобальная область системы), 205
 UGA (Глобальная область пользователя),
 205; 230
 автоматическое управление памятью
 (ASM), 173; 257
 контроль расхода памяти, 216
 рабочая, 216
 ручное управление памятью, 216
 Параллелизм
 процедурный, 867
 самодельный, 871; 875
 управление параллелизмом, 359
 Параллельное выполнение, 837
 Повтор (redo), 427
 Подключение (connection), 260; 264; 265
 Прагма (pragma), 421
 Представление
 V\$, 155
 V\$DIAG_INFO, 158
 Программирование
 уровневое, 96
 Протокол
 2PC, 285; 418
 LDAP, 131
 Процедурный параллелизм, 867
 Процесс, 259
 ARCP (архивный), 288
 ASMB (автоматического управления про-
 странством хранения), 290
 СКРТ (выполнения контрольных точек),
 285
 CTWR (отслеживания изменений), 295
 DBRM (диспетчера ресурсов), 289
 DBWn (записи блоков базы данных), 286
 DIAG (диагностики), 289

EMNn (монитора событий), 294
 FBDA (архивации ретроспективных
 данных), 289
 GEN0 (выполнения общей задачи), 289
 Jnnn (очереди заданий), 293
 LCK0 (блокировки), 291
 LGWR (записи журналов), 287
 LMD (демона диспетчера блокировки), 291
 LMON (монитора блокировки), 291
 LMSn (сервера диспетчера блокировки), 291
 LREG (регистрации прослушивателем), 283
 MMAN (диспетчер памяти), 294
 MMON (монитора управляемости), 294
 PMON (монитор процессов), 282
 PSP0 (генератора процессов), 291
 QMNC (расширенные очереди), 293
 RECO (распределенное восстановление
 базы данных), 285
 RVWR (записи данных восстановления), 295
 SJQ0 (координатора очередей заданий), 292
 SMC0 (координатора управления
 пространством), 292
 SMON (системный монитор), 283
 TMON (монитора перемещения), 295
 VKRM (виртуального планировщика для
 диспетчера ресурсов), 292
 виртуального хранителя времени (VKTM), 291
 подчиненный, 259; 296
 ввода-вывода, 296
 прослушивающий TNS, 131
 ребалансировки (RBAL), 290
 серверный, 259; 260
 фоновый, 123; 124; 259; 277; 278
 служебный, 292
 специализированный, 278

Пул
 Java, 232; 251
 Streams (Потоки), 232; 252
 большой (large pool), 232; 249
 неопределенный ("Null" pool), 233
 подключений, 273
 разделяемый (shared pool), 232; 246
 резидентный, 264
 рециклирующий (recycle pool), 239
 стандартный (default pool), 238
 удерживающий (keep pool), 239

Р

Раздел
 чистый (raw), 173
 Репликация, 100
 Ресурс
 блокирование ресурса, 77

С

Сеанс (session), 260–265
 Сегмент, 174; 178
 отмены (отката), 80
 управление пространством сегментов, 493
 Секционирование (partitioning), 749
 индексов, 793
 ключ секционирования (partition key), 750
 по виртуальному столбцу, 785
 по диапазонам ключей, 762
 по интервалам ключей, 771; 783
 по списку значений ключа, 769
 по ссылкам, 778; 783
 составное, 787
 хеш-, 765
 Секция (partition), 490
 Семафор (mutex), 355
 Сервер
 Pnnn (выполнения параллельного запроса), 297
 выделенный, 126; 276
 выделенный или разделяемый, 125
 параллельного выполнения, 846
 разделяемый, 128; 263; 276
 управления, 145
 Сеть Oracle, 261
 Системная глобальная область, 123
 Системный монитор, 283
 Системы OLTP, 813
 Сканирование
 быстрое полное, 607
 События Oracle, 155
 Среда SQL*Plus, 266
 Средство
 ADR, 158
 AQ, 277
 Data Pump, 295
 DRCP, 276
 FGAC, 106
 Стандарт
 SQL89, 90
 SQL99, 90; 91
 Строка
 логические идентификаторы строк, 526
 перемещение строк, 789
 перемещенная, 503
 символьная, 679

Т

Таблица (table), 490
 SQL Server
 кластеризованная, 300
 вложенная, 488; 491; 549
 хранение вложенных таблиц, 559
 внешняя, 489; 881; 908
 использование для загрузки разных файлов, 896
 настройка, 883
 предварительная обработка (preprocessing), 898
 временная, 488; 563
 индекс-таблица (IOT), 487; 509; 588
 кластеризованная, 487
 объектная, 489; 577
 подсекция (subpartition) таблицы, 490
 толстая, 607
 тонкая, 607
 традиционная, 487; 505
 хеш-таблица
 кластеризованная, 488; 536
 Табличное пространство (tablespace), 174; 177; 178
 управляемое словарем, 179
 управляемое локально, 179
 Тестирование
 с переменными привязки, 352
 Технология, 116
 RAC (Real Application Clusters), 116
 Типы данных Oracle
 BFILE, 674; 727; 744
 BINARY_DOUBLE, 672; 692; 698
 BINARY_FLOAT, 672; 692; 698
 BLOB, 673; 727
 CHAR, 671
 CLOB, 674; 727
 DATE, 673; 708; 710
 DECIMAL(p,s), 699
 DOUBLE PRECISION, 699
 FLOAT(p), 699
 INTEGER, 699
 INTERVAL, 708; 724
 INTERVAL DAY TO SECOND, 673; 726
 INTERVAL YEAR TO MONTH, 673; 725
 LONG, 672; 701
 LONG RAW, 673; 701
 NCHAR, 671
 NCLOB, 674; 727
 NUMBER, 672; 692
 NUMERIC(p,s), 699
 NVARCHAR2, 672

RAW, 672; 686
 REAL, 699
 ROWID, 674; 745
 SMALLINT, 699
 TIMESTAMP, 673; 708; 716
 TIMESTAMP WITH LOCAL TIME ZONE,
 673; 722
 TIMESTAMP WITH TIME ZONE, 673; 720
 UROWID, 674; 745
 VARCHAR2, 672; 683
 Транзакция, 302; 389
 автономная, 420
 атомарная, 390
 уровни изоляции транзакций, 361
 Триггеры событий базы данных, 156

У

Утилита
 EXPDP, 200
 Export (Экспорт), 198
 IMPDP, 200
 Import (Импорт), 198

Ф

Файл
 Data Pump (помпа данных), 138; 200
 DMP, 198
 init.ora, 143
 временный, 137; 182
 дампа, 138
 данных, 137; 172
 журнала повторения транзакций, 184
 журнальный, 137
 отслеживания изменений, 138; 194
 параметров (init.ora), 119; 137; 138
 параметров сервера (SPFILE), 146
 паролей, 137; 191
 плоский, 138; 202; 881
 ретроспективного журнала, 138
 сигнальный, 137; 168
 текстовый, 152
 генерируемый в ответ на внутренние
 ошибки, 162
 трассировочный, 137; 156
 маркирование трассировочных файлов, 161
 унаследованный, 143
 управляющий, 137; 184

Файловая система
 кластеризованная, 173
 Формат
 CSV, 917
 Data Pump, 881; 909
 Функция
 AUTOTRACE, 266
 DUMP, 678
 COUNT(*), 894
 exec(), 262
 fork(), 262

Х

Хранилища данных, 803

Ч

Число с плавающей точкой, 698
 Чтение
 согласованное, 378
 текущее, 378

Э

Экземпляр, 114; 115; 122
 Экстент, 175

Я

Язык
 DDL, 327; 398
 DML, 326

Oracle

ДЛЯ ПРОФЕССИОНАЛОВ

**Архитектура, методики
программирования и основные
особенности версий 9i, 10g, 11g и 12c**

В третьем издании продолжается исследование применения баз данных Oracle для построения масштабируемых приложений, которые эффективно функционируют и создают корректные результаты. Том Кайт и Дарл Кун придерживаются простой философии: вы можете трактовать базу данных Oracle как черный ящик и только помещать туда данные или же вы можете полностью разобраться в ее работе и эксплуатировать как мощную вычислительную среду. Выбрав второй подход, вы обнаружите, что остается совсем немного задач по управлению информацией, которые невозможно решить быстро и элегантно.

В полностью пересмотренном третьем издании раскрываются особенности разработки приложений вплоть до версии Oracle 12c. Большая часть нового материала посвящена появившемуся набору облачных средств Oracle и применению подключаемых баз данных. Каждое новое средство объясняется в манере, подкрепляемой примерами, с обсуждением не только того, как оно работает, но также способа реализации программного обеспечения с его использованием и связанных с ним распространенных заблуждений.

Не воспринимайте базу данных Oracle как черный ящик! Возьмите эту книгу. Загляните за кулисы. Ускорьте свой карьерный рост.

В КНИГЕ РАССМАТРИВАЮТСЯ СЛЕДУЮЩИЕ ТЕМЫ

- Разработка подхода к решению задач на основе фактов
- Управление транзакциями в средах с высокой степенью параллелизма
- Ускорение доступа к данным благодаря эффективному проектированию таблиц и индексов
- Управление файлами и структурами памяти, позволяющее достичь высокой производительности и надежности
- Улучшение масштабирования с помощью секционирования и параллельной обработки
- Загрузка и выгрузка данных для взаимодействия с внешними системами

ДЛЯ КОГО ПРЕДНАЗНАЧЕНА ЭТА КНИГА

Эта книга предназначена для администраторов баз данных Oracle, разработчиков на PL/SQL и Java, которые пишут код, развертываемый внутри базы данных, и разработчиков внешних приложений, использующих базы данных Oracle в качестве хранилища. Книга ориентирована на тех, кто стремится создавать эффективные и масштабируемые приложения.

НА ВЕБ-САЙТЕ


Исходные коды всех примеров, рассмотренных в книге, можно загрузить с веб-сайта издательства по адресу:

<http://www.williamspublishing.com/Books/978-5-8459-2042-3.html>

Категория: базы данных/Oracle

Предмет рассмотрения: 9i, 10g, 11g, 12c

Уровень: для пользователей средней и высокой квалификации


www.williamspublishing.com

Apress®
www.apress.com



ISBN= 978-5-8459-2042-3



9 785845 920423