



Solution Manual Database System Concepts 6th Edition

Introduction to Computer Science (Patuakhali Science and Technology University)

INSTRUCTOR'S MANUAL
TO ACCOMPANY

Database System Concepts

Sixth Edition

Abraham Silberschatz
Yale University

Henry F. Korth
Lehigh University

S. Sudarshan
Indian Institute of Technology, Bombay

Copyright © 2010 A. Silberschatz, H. Korth, and S. Sudarshan

This document is available free of charge on

StuDocu.com

?? Lee (yulim37@nate.com) ??(?) ?????

Contents

| | | |
|------------|--|-----|
| Chapter 1 | Introduction | 1 |
| Chapter 2 | Introduction to the Relational Model | 7 |
| Chapter 3 | Introduction to SQL | 11 |
| Chapter 4 | Intermediate SQL | 25 |
| Chapter 5 | Advanced SQL | 31 |
| Chapter 6 | Formal Relational Query Languages | 43 |
| Chapter 7 | Database Design and the E-R Model | 51 |
| Chapter 8 | Relational Database Design | 67 |
| Chapter 9 | Application Design and Development | 77 |
| Chapter 10 | Storage and File Structure | 91 |
| Chapter 11 | Indexing and Hashing | 97 |
| Chapter 12 | Query Processing | 103 |
| Chapter 13 | Query Optimization | 109 |
| Chapter 14 | Transactions | 115 |
| Chapter 15 | Concurrency Control | 123 |
| Chapter 16 | Recovery System | 131 |
| Chapter 17 | Database-System Architectures | 139 |
| Chapter 18 | Parallel Databases | 143 |
| Chapter 19 | Distributed Databases | 149 |
| Chapter 20 | Data Mining | 157 |
| Chapter 21 | Information Retrieval | 163 |
| Chapter 22 | Object-Based Databases | 169 |
| Chapter 23 | XML | 175 |

| | | |
|------------|--|-----|
| Chapter 24 | Advanced Application Development | 191 |
| Chapter 25 | Advanced Data Types and New Applications | 197 |
| Chapter 26 | Advanced Transaction Processing | 201 |

Preface

This volume is an instructor's manual for the 6th edition of *Database System Concepts* by Abraham Silberschatz, Henry F. Korth and S. Sudarshan. It contains answers to the exercises at the end of each chapter of the book. (Beginning with the 5th edition, solutions for Practice Exercises have been made available on the Web; to avoid duplication, these are not included in the instructors manual.)

Before providing answers to the exercises for each chapter, we include a few remarks about the chapter. The nature of these remarks vary. They include explanations of the inclusion or omission of certain material, and remarks on how we teach the chapter in our own courses. The remarks also include suggestions on material to skip if time is at a premium, and tips on software and supplementary material that can be used for programming exercises.

The Web home page of the book, at <http://www.db-book.com>, contains a variety of useful information, including laboratory relation information such as sample data, lab exercises, and links to database software, online appendices describing the network data model, the hierarchical data model, and advanced relational database design, model course syllabi, and last but not least, up-to-date errata. We will periodically update the page with supplementary material that may be of use to teachers and students.

We would appreciate it if you would notify us of any errors or omissions in the book, as well as in the instructor's manual. Internet electronic mail should be addressed to db-book-authors@cs.yale.edu. Physical mail may be sent to Avi Silberschatz, Department of Computer Science, Yale University, 51 Prospect Street, P.O. Box 208285, New Haven, CT, 06520, USA.

Although we have tried to produce an instructor's manual which will aid all of the users of our book as much as possible, there can always be improvements. These could include improved answers, additional questions, sample test questions, programming projects, suggestions on alternative orders of presentation of the material, additional references, and so on. If you would like to suggest any such improvements to the book or the instructor's manual, we would be glad to hear from you. All contributions that we make use of will, of course, be properly credited to their contributor.

Several students at IIT Bombay contributed to the instructor manual for the 6th edition, including Mahendra Chavan, Karthik Ramachandra, Bikmal Hari Krishna, Ankush Jain, Manas Joglekar, Parakram Majumdar, Prashant Sachdeva, and Nisarg Shah. This manual is derived from the manuals for the earlier editions. John Corwin and Swathi Yadlapalli did the bulk of the work in preparing the instructors manual for the 5th edition. The manual for the 4th edition was prepared by Nilesh Dalvi, Sumit Sanghai, Gaurav Bhalotia and Arvind Hulgeri. The manual for the 3th edition was prepared by K. V. Raghavan with help from Prateek R. Kapadia. Sara Strandtman helped with the instructor manual for the 2nd and 3rd editions, while Greg Speegle and Dawn Bezviner helped us to prepare the instructor's manual for the 1th edition.

A. S.
H. F. K.
S. S.

CHAPTER 1



Introduction

Chapter 1 provides a general overview of the nature and purpose of database systems. The most important concept in this chapter is that database systems allow data to be treated at a high level of abstraction. Thus, database systems differ significantly from the file systems and general purpose programming environments with which students are already familiar. Another important aspect of the chapter is to provide motivation for the use of database systems as opposed to application programs built on top of file systems. Thus, the chapter motivates what the student will be studying in the rest of the course.

The idea of abstraction in database systems deserves emphasis throughout, not just in discussion of Section 1.3. The overview of the structure of databases is, of necessity, rather brief, and is meant only to give the student a rough idea of some of the concepts. The student may not initially be able to fully appreciate the concepts described here, but should be able to do so by the end of the course.

The specifics of the E-R, relational, and object-oriented models are covered in later chapters. These models can be used in Chapter 1 to reinforce the concept of abstraction, with syntactic details deferred to later in the course.

If students have already had a course in operating systems, it is worthwhile to point out how the OS and DBMS are related. It is useful also to differentiate between concurrency as it is taught in operating systems courses (with an orientation towards files, processes, and physical resources) and database concurrency control (with an orientation towards granularity finer than the file level, recoverable transactions, and resources accessed associatively rather than physically). If students are familiar with a particular operating system, that OS's approach to concurrent file access may be used for illustration.

Exercises

- 1.7 List four applications you have used that most likely employed a database system to store persistent data.

Answer:

- Banking: For account information, transfer of funds, banking transactions.
- Universities: For student information, online assignment submissions, course registrations, and grades.
- Airlines: For reservation of tickets, and schedule information.
- Online news sites: For updating new, maintenance of archives.
- Online-trade: For product data, availability and pricing informations, order-tracking facilities, and generating recommendation lists.

1.8 List four significant differences between a file-processing system and a DBMS.

Answer: Some main differences between a database management system and a file-processing system are:

- Both systems contain a collection of data and a set of programs which access that data. A database management system coordinates both the physical and the logical access to the data, whereas a file-processing system coordinates only the physical access.
- A database management system reduces the amount of data duplication by ensuring that a physical piece of data is available to all programs authorized to have access to it, whereas data written by one program in a file-processing system may not be readable by another program.
- A database management system is designed to allow flexible access to data (i.e., queries), whereas a file-processing system is designed to allow pre-determined access to data (i.e., compiled programs).
- A database management system is designed to coordinate multiple users accessing the same data at the same time. A file-processing system is usually designed to allow one or more programs to access different data files at the same time. In a file-processing system, a file can be accessed by two programs concurrently only if both programs have read-only access to the file.

1.9 Explain the concept of physical data independence, and its importance in database systems.

Answer: Physical data independence is the ability to modify the physical scheme without making it necessary to rewrite application programs. Such modifications include changing from unblocked to blocked record storage, or from sequential to random access files. Such a modification might be adding a field to a record; an application program's view hides this change from the program.

1.10 List five responsibilities of a database-management system. For each responsibility, explain the problems that would arise if the responsibility were not discharged.

Answer: A general purpose database-management system (DBMS) has five responsibilities:

- a. interaction with the file manager.
- b. integrity enforcement.
- c. security enforcement.
- d. backup and recovery.
- e. concurrency control.

If these responsibilities were not met by a given DBMS (and the text points out that sometimes a responsibility is omitted by design, such as concurrency control on a single-user DBMS for a micro computer) the following problems can occur, respectively:

- a. No DBMS can do without this, if there is no file manager interaction then nothing stored in the files can be retrieved.
- b. Consistency constraints may not be satisfied, for example an instructor may belong to a non-existent department, two students may have the same ID, account balances could go below the minimum allowed, and so on.
- c. Unauthorized users may access the database, or users authorized to access part of the database may be able to access parts of the database for which they lack authority. For example, a low-level user could get access to national defense secret codes, or employees could find out what their supervisors earn (which is presumably a secret).
- d. Data could be lost permanently, rather than at least being available in a consistent state that existed prior to a failure.
- e. Consistency constraints may be violated despite proper integrity enforcement in each transaction. For example, incorrect bank balances might be reflected due to simultaneous withdrawals and deposits on the same account, and so on.

- 1.11** List at least two reasons why database systems support data manipulation using a declarative query language such as SQL, instead of just providing a library of C or C++ functions to carry out data manipulation.

Answer:

- a. Declarative languages are easier for programmers to learn and use (and even more so for non-programmers).
- b. The programmer does not have to worry about how to write queries to ensure that they will execute efficiently; the choice of an efficient execution technique is left to the database system. The declarative specification makes it easier for the database system to make a proper choice of execution technique.

1.12 Explain what problems are caused by the design of the table in Figure 1.4.

Answer:

- If a department has more than one instructor, the building name and budget get repeated multiple times. Updates to the building name and budget may get performed on some of the copies but not others, resulting in an inconsistent state where it is not clear what is the actual building name and budget of a department.
- A department needs to have at least one instructor in order for building and budget information to be included in the table. Nulls can be used when there is no instructor, but null values are rather difficult to handle.
- If all instructors in a department are deleted, the building and budget information are also lost. Ideally, we would like to have the department information in the database irrespective of whether the department has an associated instructor or not, without resorting to null values.

1.13 What are five main functions of a database administrator?

Answer:

- To backup data
- In some cases, to create the schema definition
- To define the storage structure and access methods
- To modify the schema and/or physical organization when necessary
- To grant authorization for data access
- To specify integrity constraints

1.14 Explain the difference between two-tier and three-tier architectures. Which is better suited for Web applications? Why?

Answer: In a two-tier application architecture, the application runs on the client machine, and directly communicates with the database system running on server. In contrast, in a three-tier architecture, application code running on the client's machine communicates with an application server at the server, and never directly communicates with the database. The three-tier architecture is better suited for Web applications.

1.15 Describe at least 3 tables that might be used to store information in a social-networking system such as Facebook.

Answer: Some possible tables are:

- a. A *users* table containing users, with attributes such as account name, real name, age, gender, location, and other profile information.
- b. A *content* table containing user provided content, such as text and images, associated with the user who uploaded the content.

- c. A *friends* table recording for each user which other users are connected to that user. The kind of connection may also be recorded in this table.
- d. A *permissions* table, recording which category of friends are allowed to view which content uploaded by a user. For example, a user may share some photos with family but not with all friends.

CHAPTER 2



Introduction to the Relational Model

This chapter presents the relational model and a brief introduction to the relational-algebra query language. The short introduction to relational algebra is sufficient for courses that focus on application development, without going into database internals. In particular, the chapters on SQL do not require any further knowledge of relational algebra. However, courses that cover internals, in particular query processing, require a more detailed coverage of relational algebra, which is provided in Chapter 6.

Exercises

- 2.9 Consider the bank database of Figure 2.15.
- What are the appropriate primary keys?

employee (*person_name*, *street*, *city*)
works (*person_name*, *company_name*, *salary*)
company (*company_name*, *city*)

Figure 2.14 Relational database for Exercises 2.1, 2.7, and 2.12.

branch (*branch_name*, *branch_city*, *assets*)
customer (*customer_name*, *customer_street*, *customer_city*)
loan (*loan_number*, *branch_name*, *amount*)
borrower (*customer_name*, *loan_number*)
account (*account_number*, *branch_name*, *balance*)
depositor (*customer_name*, *account_number*)

Figure 2.15 Banking database for Exercises 2.8, 2.9, and 2.13.

b. Given your choice of primary keys, identify appropriate foreign keys.

Answer:

a. The primary keys of the various schema are underlined. Although in a real bank the customer name is unlikely to be a primary key, since two customers could have the same name, we use a simplified schema where we assume that names are unique. We allow customers to have more than one account, and more than one loan.

```
branch(branch_name, branch_city, assets)
customer (customer_name, customer_street, customer_city)
loan (loan_number, branch_name, amount)
borrower (customer_name, loan_number)
account (account_number, branch_name, balance)
depositor (customer_name, account_number)
```

- b. The foreign keys are as follows
- i. For *loan*: *branch_name* referencing *branch*.
 - ii. For *borrower*: Attribute *customer_name* referencing *customer* and *loan_number* referencing *loan*
 - iii. For *account*: *branch_name* referencing *branch*.
 - iv. For *depositor*: Attribute *customer_name* referencing *customer* and *account_number* referencing *account*

2.10 Consider the *advisor* relation shown in Figure 2.8, with *s_id* as the primary key of *advisor*. Suppose a student can have more than one advisor. Then, would *s_id* still be a primary key of the *advisor* relation? If not, what should the primary key of *advisor* be?

Answer: No, *s_id* would not be a primary key, since there may be two (or more) tuples for a single student, corresponding to two (or more) advisors. The primary key should then be *s_id*, *i_id*.

2.11 Describe the differences in meaning between the terms *relation* and *relation schema*.

Answer: A relation schema is a type definition, and a relation is an instance of that schema. For example, *student* (*ss#*, *name*) is a relation schema and

| | |
|-------------|------|
| 123-456-222 | John |
| 234-567-999 | Mary |

is a relation based on that schema.

2.12 Consider the relational database of Figure 2.14. Give an expression in the relational algebra to express each of the following queries:

- a. Find the names of all employees who work for “First Bank Corporation”.

- b. Find the names and cities of residence of all employees who work for “First Bank Corporation”.
- c. Find the names, street address, and cities of residence of all employees who work for “First Bank Corporation” and earn more than \$10,000.

Answer:

- a. $\Pi_{person_name} (\sigma_{company_name = \text{“First Bank Corporation”}} (works))$
- b. $\Pi_{person_name, city} (employee \bowtie (\sigma_{company_name = \text{“First Bank Corporation”}} (works)))$
- c. $\Pi_{person_name, street, city} (\sigma_{(company_name = \text{“First Bank Corporation”} \wedge salary > 10000)} (works \bowtie employee))$

2.13 Consider the bank database of Figure 2.15. Give an expression in the relational algebra for each of the following queries:

- a. Find all loan numbers with a loan value greater than \$10,000.
- b. Find the names of all depositors who have an account with a value greater than \$6,000.
- c. Find the names of all depositors who have an account with a value greater than \$6,000 at the “Uptown” branch.

Answer:

- a. $\Pi_{loan_number} (\sigma_{amount > 10000} (loan))$
- b. $\Pi_{customer_name} (\sigma_{balance > 6000} (depositor \bowtie account))$
- c. $\Pi_{customer_name} (\sigma_{balance > 6000 \wedge branch_name = \text{“Uptown”}} (depositor \bowtie account))$

2.14 List two reasons why null values might be introduced into the database.

Answer: Nulls may be introduced into the database because the actual value is either unknown or does not exist. For example, an employee whose address has changed and whose new address is not yet known should be retained with a null address. If employee tuples have a composite attribute *dependents*, and a particular employee has no dependents, then that tuple’s *dependents* attribute should be given a null value.

2.15 Discuss the relative merits of procedural and nonprocedural languages.

Answer: Nonprocedural languages greatly simplify the specification of queries (at least, the types of queries they are designed to handle). The free the user from having to worry about how the query is to be evaluated; not only does this reduce programming effort, but in fact in most situations the query optimizer can do a much better task of choosing the best way to evaluate a query than a programmer working by trial and error. On the other hand, procedural languages are far more powerful in terms of what computations they can perform. Some tasks can either not be

done using nonprocedural languages, or are very hard to express using nonprocedural languages, or execute very inefficiently if specified in a nonprocedural manner.

CHAPTER 3



Introduction to SQL

Chapter 3 introduces the relational language SQL. Further details of the SQL language are provided in Chapters 4 and 5.

Although our discussion is based on SQL standards, no database system implements the standards exactly as specified, and there are a number of minor syntactic differences that need to be kept in mind. Although we point out some of these differences where required, the system manuals of the database system you use should be used as supplements.

Although it is possible to cover this chapter using only handwritten exercises, we strongly recommend providing access to an actual database system that supports SQL. A style of exercise we have used is to create a moderately large database and give students a list of queries in English to write and run using SQL. We publish the actual answers (that is the result relations they should get, not the SQL they must enter). By using a moderately large database, the probability that a “wrong” SQL query will just happen to return the “right” result relation can be made very small. This approach allows students to check their own answers for correctness immediately rather than wait for grading and thereby it speeds up the learning process. A few such example databases are available on the Web home page of this book, <http://db-book.com>.

Exercises that pertain to database design are best deferred until after Chapter 8.

Exercises

3.11 Write the following queries in SQL, using the university schema.

- a. Find the names of all students who have taken at least one Comp. Sci. course; make sure there are no duplicate names in the result.
- b. Find the IDs and names of all students who have not taken any course offering before Spring 2009.

- c. For each department, find the maximum salary of instructors in that department. You may assume that every department has at least one instructor.
- d. Find the lowest, across all departments, of the per-department maximum salary computed by the preceding query.

Answer:

- a. SQL query:

```
select  name
from    student natural join takes natural join course
where   course.dept = 'Comp. Sci.'
```

- b. SQL query:

```
select  id, name
from    student
except
select  id, name
from    student natural join takes
where   year < 2009
```

Since the **except** operator eliminates duplicates, there is no need to use a **select distinct** clause, although doing so would not affect correctness of the query.

- c. SQL query:

```
select  dept, max(salary)
from    instructor
group by dept
```

- d. SQL query:

```
select  min(maxsalary)
from    (select dept, max(salary) as maxsalary
from    instructor
group by dept)
```

3.12 Write the following queries in SQL, using the university schema.

- a. Create a new course “CS-001”, titled “Weekly Seminar”, with 0 credits.

- b. Create a section of this course in Autumn 2009, with *section_id* of 1.
- c. Enroll every student in the Comp. Sci. department in the above section.
- d. Delete enrollments in the above section where the student's name is Chavez.
- e. Delete the course CS-001. What will happen if you run this delete statement without first deleting offerings (sections) of this course.
- f. Delete all *takes* tuples corresponding to any section of any course with the word "database" as a part of the title; ignore case when matching the word with the title.

Answer:

- a. SQL query:

```
insert into course  
values ('CS-001', 'Weekly Seminar', 'Comp. Sci.', 0)
```

- b. SQL query:

```
insert into section  
values ('CS-001', 1, 'Autumn', 2009, null, null, null)
```

Note that the building, roomnumber and slot were not specified in the question, and we have set them to null. The same effect would be obtained if they were specified to default to null, and we simply omitted values for these attributes in the above insert statement. (Many database systems implicitly set the default value to null, even if not explicitly specified.)

- c. SQL query:

```
insert into takes  
select id, 'CS-001', 1, 'Autumn', 2009, null  
from student  
where dept_name = 'Comp. Sci.'
```

- d. SQL query:

```

delete from takes
where course_id = 'CS-001' and section_id = 1 and
      year = 2009 and semester = 'Autumn' and
      id in (select id
             from student
             where name = 'Chavez')

```

Note that if there is more than one student named Chavez, all such students would have their enrollments deleted. If we had used = instead of **in**, an error would have resulted if there were more than one student named Chavez.

e. SQL query:

```

delete from takes
where course_id = 'CS-001'

delete from section
where course_id = 'CS-001'

delete from course
where course_id = 'CS-001'

```

If we try to delete the course directly, there will be a foreign key violation because *section* has a foreign key reference to *course*; similarly, we have to delete corresponding tuples from *takes* before deleting sections, since there is a foreign key reference from *takes* to *section*. As a result of the foreign key violation, the transaction that performs the delete would be rolled back.

f. SQL query:

```

delete from takes
where course_id in
      (select course_id
       from course
       where lower(title) like '%database%')

```

3.13 Write SQL DDL corresponding to the schema in Figure 3.18. Make any reasonable assumptions about data types, and be sure to declare primary and foreign keys.

Answer:

a. SQL query:

person (*driver_id*, *name*, *address*)
car (*license*, *model*, *year*)
accident (*report_number*, *date*, *location*)
owns (*driver_id*, *license*)
participated (*report_number*, *license*, *driver_id*, *damage_amount*)

Figure 3.18 Insurance database for Exercises 3.4 and 3.14.

```

create table person
  (driver_id varchar(50),
   name      varchar(50),
   address   varchar(50),
   primary key (driver_id))
  
```

b. SQL query:

```

create table car
  (license varchar(50),
   model   varchar(50),
   year    integer,
   primary key (license))
  
```

c. SQL query:

```

create table accident
  (report_number integer,
   date          date,
   location      varchar(50),
   primary key (report_number))
  
```

d. SQL query:

```

create table owns
  (driver_id varchar(50),
   license   varchar(50),
   primary key (driver_id, license)
   foreign key (driver_id) references person
   foreign key (license) references car)
  
```

e. SQL query:

```

create table participated
  (report_number integer,
   license varchar(50),
   driver_id varchar(50),
   damage_amount integer,
   primary key (report_number, license)
   foreign key (license) references car
   foreign key (report_number) references accident))

```

3.14 Consider the insurance database of Figure 3.18, where the primary keys are underlined. Construct the following SQL queries for this relational database.

- Find the number of accidents in which the cars belonging to “John Smith” were involved.
- Update the damage amount for the car with license number “AABB2000” in the accident with report number “AR2197” to \$3000.

Answer: Note: The *participated* relation relates drivers, cars, and accidents.

- SQL query:

```

select    count (*)
from      accident
where      exists
            (select *
             from participated, owns, person
             where owns.driver_id = person.driver_id
                   and person.name = 'John Smith'
                   and owns.license = participated.license
                   and accident.report_number = participated.report_number)

```

The query can be written in other ways too; for example without a subquery, by using a join and selecting **count**(**distinct** *report_number*) to get a count of number of accidents involving the car.

- SQL query:

```

update participated
set damage_amount = 3000
where report_number = “AR2197” and
       license = “AABB2000”)

```

3.15 Consider the bank database of Figure 3.19, where the primary keys are underlined. Construct the following SQL queries for this relational database.

branch(*branch_name*, *branch_city*, *assets*)
customer(*customer_name*, *customer_street*, *customer_city*)
loan(*loan_number*, *branch_name*, *amount*)
borrower(*customer_name*, *loan_number*)
account(*account_number*, *branch_name*, *balance*)
depositor(*customer_name*, *account_number*)

Figure 3.19 Banking database for Exercises 3.8 and 3.15.

- Find all customers who have an account at *all* the branches located in “Brooklyn”.
- Find out the total sum of all loan amounts in the bank.
- Find the names of all branches that have assets greater than those of at least one branch located in “Brooklyn”.

Answer:

- SQL query:

```

with branchcount as
  (select count(*)
   branch
   where branch_city = 'Brooklyn')
select customer_name
from customer c
where branchcount =
  (select count(distinct branch_name)
   from (customer natural join depositor natural join account
        natural join branch) as d
   where d.customer_name = c.customer_name)
  
```

There are other ways of writing this query, for example by first finding customers who do not have an account at some branch in Brooklyn, and then removing these customers from the set of all customers by using an **except** clause.

- SQL query:

```

select sum(amount)
from loan
  
```

- SQL query:

employee (*employee_name*, *street*, *city*)
works (*employee_name*, *company_name*, *salary*)
company (*company_name*, *city*)
manages (*employee_name*, *manager_name*)

Figure 3.20 Employee database for Exercises 3.9, 3.10, 3.16, 3.17, and 3.20.

```

select branch_name
from branch
where assets > some
      (select assets
from branch
where branch_city = 'Brooklyn')
  
```

The keyword **any** could be used in place of **some** above.

- 3.16** Consider the employee database of Figure 3.20, where the primary keys are underlined. Give an expression in SQL for each of the following queries.
- Find the names of all employees who work for First Bank Corporation.
 - Find all employees in the database who live in the same cities as the companies for which they work.
 - Find all employees in the database who live in the same cities and on the same streets as do their managers.
 - Find all employees who earn more than the average salary of all employees of their company.
 - Find the company that has the smallest payroll.

Answer:

- Find the names of all employees who work for First Bank Corporation.

```

select employee_name
from works
where company_name = 'First Bank Corporation'
  
```

- Find all employees in the database who live in the same cities as the companies for which they work.

```

select e.employee_name
from employee e, works w, company c
where e.employee_name = w.employee_name and e.city = c.city and
      w.company_name = c.company_name
  
```

- c. Find all employees in the database who live in the same cities and on the same streets as do their managers.

```
select P.employee_name
from employee P, employee R, manages M
where P.employee_name = M.employee_name and
      M.manager_name = R.employee_name and
      P.street = R.street and P.city = R.city
```

- d. Find all employees who earn more than the average salary of all employees of their company.

```
select employee_name
from works T
where salary > (select avg (salary)
                from works S
                where T.company_name = S.company_name)
```

The primary key constraint on *works* ensures that each person works for at most one company.

- e. Find the company that has the smallest payroll.

```
select company_name
from works
group by company_name
having sum (salary) <= all (select sum (salary)
                            from works
                            group by company_name)
```

3.17 Consider the relational database of Figure 3.20. Give an expression in SQL for each of the following queries.

- Give all employees of First Bank Corporation a 10 percent raise.
- Give all managers of First Bank Corporation a 10 percent raise.
- Delete all tuples in the *works* relation for employees of Small Bank Corporation.

Answer:

- Give all employees of First Bank Corporation a 10-percent raise. (the solution assumes that each person works for at most one company.)

```
update works
set salary = salary * 1.1
where company_name = 'First Bank Corporation'
```

- Give all managers of First Bank Corporation a 10-percent raise.

```

update works
set salary = salary * 1.1
where employee_name in (select manager_name
                           from manages)
                           and company_name = 'First Bank Corporation'

```

- c. Delete all tuples in the *works* relation for employees of Small Bank Corporation.

```

delete from works
where company_name = 'Small Bank Corporation'

```

3.18 List two reasons why null values might be introduced into the database.

Answer:

- “null” signifies an unknown value.
- “null” is also used when a value does not exist.

3.19 Show that, in SQL, $\langle \rangle$ **all** is identical to **not in**.

Answer: Let the set S denote the result of an SQL subquery. We compare $(x \langle \rangle \text{all } S)$ with $(x \text{ not in } S)$. If a particular value x_1 satisfies $(x_1 \langle \rangle \text{all } S)$ then for all elements y of S $x_1 \neq y$. Thus x_1 is not a member of S and must satisfy $(x_1 \text{ not in } S)$. Similarly, suppose there is a particular value x_2 which satisfies $(x_2 \text{ not in } S)$. It cannot be equal to any element w belonging to S , and hence $(x_2 \langle \rangle \text{all } S)$ will be satisfied. Therefore the two expressions are equivalent.

3.20 Give an SQL schema definition for the employee database of Figure 3.20. Choose an appropriate domain for each attribute and an appropriate primary key for each relation schema.

Answer:

```

create table    employee
(employee_name varchar(20),
 street        char(30),
 city          varchar(20),
 primary key   (employee_name))

```

```

create table    works
(employee_name person_names,
 company_name  varchar(20),
 salary        numeric(8, 2),
 primary key   (employee_name))

```

```

create table    company
(company_name  varchar(20),
 city          varchar(20),
 primary key   (company_name))

```

member(*memb_no*, *name*, *age*)
book(*isbn*, *title*, *authors*, *publisher*)
borrowed(*memb_no*, *isbn*, *date*)

Figure 3.21 Library database for Exercise 3.21.

```
create table manages
(employee_name varchar(20),
manager_name varchar(20),
primary key (employee_name))
```

3.21 Consider the library database of Figure 3.21. Write the following queries in SQL.

- Print the names of members who have borrowed any book published by “McGraw-Hill”.
- Print the names of members who have borrowed all books published by “McGraw-Hill”.
- For each publisher, print the names of members who have borrowed more than five books of that publisher.
- Print the average number of books borrowed per member. Take into account that if an member does not borrow any books, then that member does not appear in the *borrowed* relation at all.

Answer:

- Print the names of members who have borrowed any book published by McGraw-Hill.

```
select name
from member m, book b, borrowed l
where m.memb_no = l.memb_no
and l.isbn = b.isbn and
b.publisher = 'McGrawHill'
```

- Print the names of members who have borrowed all books published by McGraw-Hill. (We assume that all books above refers to all books in the *book* relation.)

```

select distinct m.name
from member m
where not exists
    ((select isbn
      from book
      where publisher = 'McGrawHill')
     except
     (select isbn
      from borrowed l
      where l.memb_no = m.memb_no))

```

- c. For each publisher, print the names of members who have borrowed more than five books of that publisher.

```

select publisher, name
from (select publisher, name, count (isbn)
      from member m, book b, borrowed l
      where m.memb_no = l.memb_no
      and l.isbn = b.isbn
      group by publisher, name) as
      membpub(publisher, name, count_books)
where count_books > 5

```

The above query could alternatively be written using the **having** clause.

- d. Print the average number of books borrowed per member.

```

with memcount as
    (select count(*)
     from member)
select count(*) / memcount
from borrowed

```

Note that the above query ensures that members who have not borrowed any books are also counted. If we instead used **count(distinct memb_no)** from *borrowed*, we would not account for such members.

3.22 Rewrite the **where** clause

where unique (select title from course)

without using the **unique** construct.

Answer:

```

where(
  (select count(title)
   from course) =
  (select count(distinct title)
   from course))

```

3.23 Consider the query:

```

select course_id, semester, year, section_id, avg(credits_earned)
from takes natural join student
where year = 2009
group by course_id, semester, year, section_id
having count(ID) >= 2;

```

Explain why joining *section* as well in the **from** clause would not change the result.

Answer: The common attributes of *takes* and *section* form a foreign key of *takes*, referencing *section*. As a result, each *takes* tuple would match at most one *section* tuple, and there would not be any extra tuples in any group. Further, these attributes cannot take on the null value, since they are part of the primary key of *takes*. Thus, joining *section* in the **from** clause would not cause any loss of tuples in any group. As a result, there would be no change in the result.

3.24 Consider the query:

```

with dept_total (dept_name, value) as
  (select dept_name, sum(salary)
   from instructor
   group by dept_name),
dept_total_avg(value) as
  (select avg(value)
   from dept_total)
select dept_name
from dept_total, dept_total_avg
where dept_total.value >= dept_total_avg.value;

```

Rewrite this query without using the **with** construct.

Answer:

There are several ways to write this query. One way is to use subqueries in the where clause, with one of the subqueries having a second level subquery in the from clause as below.

```

select distinct dept_name d
from instructor i
where
    (select sum(salary)
     from instructor
     where department = d)
    >=
    (select avg(s)
     from
       (select sum(salary) as s
        from instructor
        group by department))

```

Note that the original query did not use the *department* relation, and any department with no instructors would not appear in the query result. If we had written the above query using *department* in the outer **from** clause, a department without any instructors could appear in the result if the condition were \leq instead of \geq , which would not be possible in the original query.

As an alternative, the two subqueries in the where clause could be moved into the from clause, and a join condition (using \geq) added.

CHAPTER 4



Intermediate SQL

In this chapter, we continue our study of SQL. We consider more complex forms of SQL queries, view definition, transactions, integrity constraints, more details regarding SQL data definition, and authorization.

As in Chapter 3 exercises, students should be encouraged to execute their queries on the sample database provided on <http://db-book.com>, to check if they generate the expected answers. Students could even be encouraged to create sample databases that can expose errors in queries, for example where an inner join operation is used erroneously in place of an outerjoin operation.

Exercises

- 4.12 For the database of Figure 4.11, write a query to find those employees with no manager. Note that an employee may simply have no manager listed or may have a *null* manager. Write your query using an outer join and then write it again using no outer join at all.

Answer:

a.

```
select employee_name
from employee natural left outer join manages
where manager_name is null
```

```
employee (employee_name, street, city)
works (employee_name, company_name, salary)
company (company_name, city)
manages (employee_name, manager_name)
```

Figure 4.11 Employee database for Figure 4.7 and 4.12.

b.

```

select employee_name
from employee e
where not exists
  (select employee_name
   from manages m
   where e.employee_name = m.employee_name and
         m.manager_name is not null)

```

4.13 Under what circumstances would the query

```

select *
from student natural full outer join takes
      natural full outer join course

```

include tuples with null values for the *title* attribute?

Answer: We first rewrite the expression with parentheses to make clear the order of the left outer join operations (the SQL standard specifies that the join operations are left associative).

```

select *
from (student natural full outer join
      takes) natural full outer join course

```

Given the above query, there are 2 cases for which the *title* attribute is null

- Since *course_id* is a foreign key in the *takes* table referencing the *course* table, the title attribute in any tuple obtained from the above query can be null if there is a course in *course* table that has a null title.
- If a student has not taken any course, as it is a **natural full outer join**, such a student's entry would appear in the result with a **null title** entry.

4.14 Show how to define a view *tot_credits* (*year*, *num_credits*), giving the total number of credits taken by students in each year.**Answer:**

```

create view tot_credits(year, tot_credits)
as
  (select year, sum(credits)
   from takes natural join course
   group by year)

```

Note that this solution assumes that there is no year where students didn't take any course, even though sections were offered.

salaried_worker (*name*, *office*, *phone*, *salary*)
hourly_worker (*name*, *hourly_wage*)
address (*name*, *street*, *city*)

Figure 4.12 Employee database for Exercise 4.16.

- 4.15** Show how to express the **coalesce** operation from Exercise 4.10 using the **case** operation.

Answer:

```
select
  case Result
    when ( $A_1$  is not null) then  $A_1$ 
    when ( $A_2$  is not null) then  $A_2$ 
    .
    .
    .
    when ( $A_n$  is not null) then  $A_n$ 
    else null
  end
from A
```

- 4.16** Referential-integrity constraints as defined in this chapter involve exactly two relations. Consider a database that includes the relations shown in Figure 4.12. Suppose that we wish to require that every name that appears in *address* appears in either *salaried_worker* or *hourly_worker*, but not necessarily in both.
- Propose a syntax for expressing such constraints.
 - Discuss the actions that the system must take to enforce a constraint of this form.

Answer:

- For simplicity, we present a variant of the SQL syntax. As part of the **create table** expression for *address* we include

foreign key (*name*) **references** *salaried_worker* **or** *hourly_worker*
 - To enforce this constraint, whenever a tuple is inserted into the *address* relation, a lookup on the *name* value must be made on the *salaried_worker* relation and (if that lookup failed) on the *hourly_worker* relation (or vice-versa).
- 4.17** Explain why, when a manager, say Satoshi, grants an authorization, the grant should be done by the manager role, rather than by the user Satoshi.

Answer: Consider the case where the authorization is provided by the user Satoshi and not the manager role. If we revoke the authorization from Satoshi, for example because Satoshi left the company, all authorizations that Satoshi had granted would also be revoked, even if the grant was to an employee whose job has not changed.

If the grant is done by the manager role, revoking authorizations from Satoshi will not result in such cascading revocation.

In terms of the authorization graph, we can treat Satoshi and the role manager as nodes. When the grant is from the manager role, revoking the manager role from Satoshi has no effect on the grants from the manager role.

- 4.18 Suppose user A , who has all authorizations on a relation r , grants select on relation r to **public** with grant option. Suppose user B then grants select on r to A . Does this cause a cycle in the authorization graph? Explain why.

Answer: Yes, it does cause a cycle in the authorization graph. The grant to public results in an edge from A to public. The grant to the **public** operator provides authorization to everyone, B is now authorized. For each privilege granted to *public*, an edge must therefore be placed between *public* and all users in the system. If this is not done, then the user will not have a path from the root (DBA). And given the with grant option, B can grant select on r to A result in an edge from B to A in the authorization graph. Thus, there is now a cycle from A to public, from public to B , and from B back to A .

- 4.19 Database systems that store each relation in a separate operating-system file may use the operating system's authorization scheme, instead of defining a special scheme themselves. Discuss an advantage and a disadvantage of such an approach.

Answer:

- Advantages:
 - Operations on the database are speeded up as the authorization procedure is carried out at the OS level.
 - No need to implement the security and authorization inside the DBMS. This may result in reduced cost.
 - Administrators need not learn new commands or use of a new UI. They can create and administer user accounts on the OS they may already be familiar with.
 - No worry of unauthorized database users having direct access to the OS files and thus bypassing the database security.
- Disadvantages:
 - Database users must correspond to operating system users.
 - Fine control on authorizations is limited by what the operating system provides and is dependent on it, For example, most operating systems

do not distinguish between insert, update and delete, they just have a coarse level privilege called "modify". Privileges such as "references" cannot be provided.

- Columnwise control is not possible. You cannot differentiate update/delete and insert authorizations.
- Cannot store more than one relation in a file.
- The with grant option is limited to what the OS provides (if any: most OS's don't provide such options) and cannot be controlled by the user or administrator.

CHAPTER 5



Advanced SQL

In this chapter we address the issue of how to access SQL from a general-purpose programming language, which is very important for building applications that use a database to store and retrieve data. We describe how procedural code can be executed within the database, either by extending the SQL language to support procedural actions, or by allowing functions defined in procedural languages to be executed within the database. We describe triggers, which can be used to specify actions that are to be carried out automatically on certain events such as insertion, deletion, or update of tuples in a specified relation. We discuss recursive queries and advanced aggregation features supported by SQL. Finally, we describe online analytic processing (OLAP) systems, which support interactive analysis of very large datasets.

Given the fact that the JDBC and ODBC protocols (and variants such as ADO.NET) are have become the primary means of accessing databases, we have significantly extended our coverage of these two protocols, including some examples. However, our coverage is only introductory, and omits many details that are useful in practise. Online tutorials/manuals or textbooks covering these protocols should be used as supplements, to help students make full use of the protocols.

Exercises

5.12 Consider the following relations for a company database:

- *emp* (*ename*, *dname*, *salary*)
- *mgr* (*ename*, *mname*)

and the Java code in Figure 5.26, which uses the JDBC API. Assume that the *userid*, *password*, *machine name*, etc. are all okay. Describe in concise English what the Java program does. (That is, produce an English sentence like “It finds the manager of the toy department,” not a line-by-line description of what each Java statement does.)

```

import java.sql.*;
public class Mystery {
    public static void main(String[] args) {
        try {
            Connection con=null;
            Class.forName("oracle.jdbc.driver.OracleDriver");
            con=DriverManager.getConnection(
                "jdbc:oracle:thin:star/X@//edgar.cse.lehigh.edu:1521/XE");
            Statement s=con.createStatement();
            String q;
            String empName = "dog";
            boolean more;
            ResultSet result;
            do {
                q = "select mname from mgr where ename = '" + empName + "'";
                result = s.executeQuery(q);
                more = result.next();
                if (more) {
                    empName = result.getString("mname");
                    System.out.println (empName);
                }
            } while (more);
            s.close();
            con.close();
        } catch(Exception e){e.printStackTrace();} }}

```

Figure 5.26 Java code for Exercise 5.12.

Answer: It prints out the manager of “dog.” that manager’s manager, etc. until we reach a manager who has no manager (presumably, the CEO, who most certainly is a cat.) NOTE: if you try to run this, use your OWN Oracle ID and password, since Star, crafty cat that she is, changes her password.

- 5.13** Suppose you were asked to define a class `MetaDisplay` in Java, containing a method `static void printTable(String r)`; the method takes a relation name r as input, executes the query “**select * from r** ”, and prints the result out in nice tabular format, with the attribute names displayed in the header of the table.
- What do you need to know about relation r to be able to print the result in the specified tabular format.
 - What JDBC methods(s) can get you the required information?
 - Write the method `printTable(String r)` using the JDBC API.

Answer:

- a. We need to know the number of attributes and names of attributes of *r* to decide the number and names of columns in the table.
- b. We can use the JDBC methods `getColumnCount()` and `getColumnName(int)` to get the required information.
- c. The method is shown below.

```
static void printTable(String r)
{
    try
    {
        Class.forName("oracle.jdbc.driver.OracleDriver");
        Connection conn = DriverManager.getConnection(
            "jdbc:oracle:thin:@db.yale.edu:2000:univdb",user,passwd);
        Statement stmt = conn.createStatement();
        ResultSet rs = stmt.executeQuery(r);
        ResultSetMetaData rsmd = rs.getMetaData();
        int count = rsmd.getColumnCount();
        System.out.println("<tr>");
        for(int i=1;i<=count;i++){
            System.out.println("<td>"+rsmd.getColumnName(i)+"</td>");
        }
        System.out.println("</tr>");
        while(rs.next()){
            System.out.println("<tr>");
            for(int i=1;i<=count;i++){
                System.out.println("<td>"+rs.getString(i)+"</td>");
            }
            System.out.println("</tr>");
        }
        stmt.close();
        conn.close();
    }
    catch(SQLException sqle)
    {
        System.out.println("SQLException : " + sqle);
    }
}
```

5.14 Repeat Exercise 5.13 using ODBC, defining `void printTable(char *r)` as a function instead of a method.

Answer:

- a. Same as for JDBC.

- b. The function `SQLNumResultCols(hstmt, &numColumn)` can be used to find the number of columns in a statement, while the function `SQLColAttribute()` can be used to find the name, type and other information about any column of a result set. set, and the names
- c. The ODBC code is similar to the JDBC code, but significantly longer. ODBC code that carries out this task may be found online at the URL <http://msdn.microsoft.com/en-us/library/ms713558.aspx> (look at the bottom of the page).

5.15 Consider an employee database with two relations

employee (employee_name, street, city)
works (employee_name, company_name, salary)

where the primary keys are underlined. Write a query to find companies whose employees earn a higher salary, on average, than the average salary at “First Bank Corporation”.

- a. Using SQL functions as appropriate.
- b. Without using SQL functions.

Answer:

- a.


```
create function avg_salary(cname varchar(15))
returns integer
declare result integer;
select avg(salary) into result
from works
where works.company_name = cname
return result;
end
select company_name
from works
where avg_salary(company_name) > avg_salary("First Bank Corporation")
```
- b.


```
select company_name
from works
group by company_name
having avg(salary) > (select avg(salary)
from works
where company_name="First Bank Corporation")
```

5.16 Rewrite the query in Section 5.2.1 that returns the name and budget of all departments with more than 12 instructors, using the **with** clause instead of using a function call.

Answer:

```

with instr_count (dept_name, number) as
    (select dept_name, count (ID)
     from instructor
     group by dept_name)
select dept_name, budget
from department, instr_count
where department.dept_name = instr_count.dept_name
and number > 12

```

- 5.17 Compare the use of embedded SQL with the use in SQL of functions defined in a general-purpose programming language. Under what circumstances would you use each of these features?

Answer: SQL functions are primarily a mechanism for extending the power of SQL to handle attributes of complex data types (like images), or to perform complex and non-standard operations. Embedded SQL is useful when imperative actions like displaying results and interacting with the user are needed. These cannot be done conveniently in an SQL only environment. Embedded SQL can be used instead of SQL functions by retrieving data and then performing the function's operations on the SQL result. However a drawback is that a lot of query-evaluation functionality may end up getting repeated in the host language code.

- 5.18 Modify the recursive query in Figure 5.15 to define a relation

prereq_depth(course_id, prereq_id, depth)

where the attribute *depth* indicates how many levels of intermediate prerequisites are there between the course and the prerequisite. Direct prerequisites have a depth of 0.

Answer:

```

with recursive prereq_depth(course_id, prereq_id, depth) as
    (select course_id, prereq_id, 0
     from prereq
    union
    select prereq.course_id, prereq_depth.prereq_id, (prereq_depth.depth + 1)
     from prereq, prereq_depth
     where prereq.prereq_id= prereq_depth.course_id)

select *
from prereq_depth

```

- 5.19 Consider the relational schema

```

part(part_id, name, cost)
subpart(part_id, subpart_id, count)

```

A tuple $(p_1, p_2, 3)$ in the *subpart* relation denotes that the part with part-id p_2 is a direct subpart of the part with part-id p_1 , and p_1 has 3 copies of p_2 . Note that p_2 may itself have further subparts. Write a recursive SQL query that outputs the names of all subparts of the part with part-id “P-100”.

Answer:

```

with recursive total_part(name) as
  (select part.name
   from subpart, part
   where subpart.part_id = "P-100" and
        subpart.part_id = part.part_id
   union
   select p2.name
   from subpart s, part p1, part p2
   where s.part_id = p1.part_id
        and p1.name total_part.name
        and s.subpart_id = p2.part_id)

select *
  from total_part

```

- 5.20 Consider again the relational schema from Exercise 5.19. Write a JDBC function using non-recursive SQL to find the total cost of part “P-100”, including the costs of all its subparts. Be sure to take into account the fact that a part may have multiple occurrences of a subpart. You may use recursion in Java if you wish.

Answer: The SQL function ‘total_cost’ is called from within the JDBC code.

SQL function:

```

create function total_cost(id char(10))

returns table(number integer)

begin
  create temporary table result (name char(10), number integer);
  create temporary table newpart (name char(10), number integer);
  create temporary table temp (name char(10), number integer);
  create temporary table final_cost(number integer);

  insert into newpart
    select subpart_id, count

```

```

    from subpart
    where part_id = id
repeat
    insert into result
    select name, number
    from newpart;

    insert into temp
    (select subpart.subpart_id, count
     from newpart, subpart
     where newpart.subpart_id = subpart.part_id;
    )
except(
    select subpart_id, count
    from result;
);

delete from newpart;
insert into newpart
select *
from temp;
delete from temp;

until not exists(select * from newpart)
end repeat;

with part_cost(number) as
select (count*cost)
    from result, part
    where result.subpart_id = part.part_id);
insert into final_cost
select *
    from part_cost;
return table final_cost;
end

```

JDBC function:

```

Connection conn = DriverManager.getConnection(
    "jdbc:oracle:thin:@db.yale.edu:2000:bankdb",
    userid,passwd);
Statement stmt = conn.createStatement();

ResultSet rset = stmt.executeQuery(
    "SELECT SUM(number) FROM TABLE(total_cost('P-100'))");

System.out.println(rset.getFloat(2));

```

5.21 Suppose there are two relations r and s , such that the foreign key B of r references the primary key A of s . Describe how the trigger mechanism can be used to implement the **on delete cascade** option, when a tuple is deleted from s .

Answer: We define triggers for each relation whose primary-key is referred to by the foreign-key of some other relation. The trigger would be activated whenever a tuple is deleted from the referred-to relation. The action performed by the trigger would be to visit all the referring relations, and delete all the tuples in them whose foreign-key attribute value is the same as the primary-key attribute value of the deleted tuple in the referred-to relation. These set of triggers will take care of the **on delete cascade** operation.

5.22 The execution of a trigger can cause another action to be triggered. Most database systems place a limit on how deep the nesting can be. Explain why they might place such a limit.

Answer: It is possible that a trigger body is written in such a way that a non-terminating recursion may result. An example of such a trigger is a *before insert* triggered on a relation that tries to insert another record into the same relation.

In general, it is extremely difficult to statically identify and prohibit such triggers from being created. Hence database systems, at runtime, put a limit on the depth of nested trigger calls.

5.23 Consider the relation, r , shown in Figure 5.27. Give the result of the following query:

```
select building, room_number, time_slot_id, count(*)
from r
group by rollup (building, room_number, time_slot_id)
```

Answer:

| | | | |
|----------|-------------|-------------|---|
| Garfield | 359 | P | 1 |
| Garfield | 359 | <i>null</i> | 1 |
| Garfield | <i>null</i> | <i>null</i> | 1 |
| Painter | 705 | N | 1 |
| Painter | 705 | <i>null</i> | 1 |
| Painter | <i>null</i> | <i>null</i> | 1 |
| Saucon | 550 | D | 1 |
| Saucon | 550 | <i>null</i> | 1 |
| Saucon | 651 | N | 1 |
| Saucon | 651 | <i>null</i> | 1 |
| Saucon | <i>null</i> | <i>null</i> | 2 |

5.24 For each of the SQL aggregate functions **sum**, **count**, **min**, and **max**, show how to compute the aggregate value on a multiset $S_1 \cup S_2$, given the aggregate values on multisets S_1 and S_2 .

On the basis of the above, give expressions to compute aggregate values with grouping on a subset S of the attributes of a relation $r(A, B, C, D, E)$, given aggregate values for grouping on attributes $T \supseteq S$, for the following aggregate functions:

- sum, count, min, and max**
- avg**
- Standard deviation

Answer: Given aggregate values on multisets S_1 and S_2 , we can calculate the corresponding aggregate values on multiset $S_1 \cup S_2$ as follows:

- $\text{sum}(S_1 \cup S_2) = \text{sum}(S_1) + \text{sum}(S_2)$
- $\text{count}(S_1 \cup S_2) = \text{count}(S_1) + \text{count}(S_2)$
- $\text{min}(S_1 \cup S_2) = \text{min}(\text{min}(S_1), \text{min}(S_2))$
- $\text{max}(S_1 \cup S_2) = \text{max}(\text{max}(S_1), \text{max}(S_2))$

Let the attribute set $T = (A, B, C, D)$ and the attribute set $S = (A, B)$. Let the aggregation on the attribute set T be stored in table *aggregation_on_t* with aggregation columns *sum_t*, *count_t*, *min_t*, and *max_t* storing **sum**, **count**, **min** and **max** resp.

- The aggregations *sum_s*, *count_s*, *min_s*, and *max_s* on the attribute set S are computed by the query:

```
select A, B, sum(sum_t) as sum_s, sum(count_t) as count_s,
       min(min_t) as min_s, max(max_t) as max_s
from aggregation_on_t
groupby A, B
```

- The aggregation *avg* on the attribute set S is computed by the query:

```
select A, B, sum(sum_t)/sum(count_t) as avg_s
from aggregation_on_t
groupby A, B
```

- For calculating standard deviation we use an alternative formula:

$$\text{stddev}(S) = \frac{\sum_{s \in S} s^2}{|S|} - \text{avg}(S)^2$$

which we get by expanding the formula

$$\text{stddev}(S) = \frac{\sum_{s \in S} (s^2 - \text{avg}(S)^2)}{|S|}$$

If S is partitioned into n sets S_1, S_2, \dots, S_n then the following relation holds:

$$\text{stddev}(S) = \frac{\sum_{S_i} |S_i| (\text{stddev}(S_i)^2 + \text{avg}(S_i)^2)}{|S|} - \text{avg}(S)^2$$

Using this formula, the aggregation **stddev** is computed by the query:

```
select A, B,
       (sum(count_t * (stddev_t*stddev_t+ avg_t* avg_t))/sum(count_t)) -
       (sum(sum_t)/sum(count_t))
from aggregation_on_t
groupby A, B
```

- 5.25 In Section 5.5.1, we used the *student_grades* view of Exercise 4.5 to write a query to find the rank of each student based on grade-point average. Modify that query to show only the top 10 students (that is, those students whose rank is 1 through 10).

Answer:

```
with s_grades as
  select ID,rank() over (order by (GPA)desc) as s_rank
  from student_grades
select ID,s_rank
from s_grades
where s_rank <= 10
```

- 5.26 Give an example of a pair of groupings that cannot be expressed by using a single **group by** clause with **cube** and **rollup**.

Answer: Consider an example of hierarchies on dimensions from Figure 5.19. We can not express a query to seek aggregation on groups (*City, Hour of day*) and (*City, Date*) using a single **group by** clause with **cube** and **rollup**.

Any single **groupby** clause with **cube** and **rollup** that computes these two groups would also compute other groups also.

- 5.27 Given relation $s(a, b, c)$, show how to use the extended SQL features to generate a histogram of c versus a , dividing a into 20 equal-sized partitions (that is, where each partition contains 5 percent of the tuples in s , sorted by a).

Answer:

```
select tile20, sum(c)
from (select c, ntile(20) over (order by (a)) as tile20
      from r) as s
groupby tile20
```

- 5.28 Consider the bank database of Figure 5.25 and the *balance* attribute of the *account* relation. Write an SQL query to compute a histogram of *balance* values, dividing the range 0 to the maximum account balance present, into three equal ranges.

Answer:

```
(select 1, count(*)
  from account
 where 3* balance <= (select max(balance)
                      from account)
)
union
(select 2, count(*)
  from account
 where 3* balance > (select max(balance)
                     from account)
    and 1.5* balance <= (select max(balance)
                        from account)
)
union
(select 3, count(*)
  from account
 where 1.5* balance > (select max(balance)
                       from account)
)
)
```


CHAPTER 6



Formal Relational Query Languages

In this chapter we study three additional formal relational languages. Relational Algebra, tuple relational calculus and domain relational calculus.

Of these three formal languages, we suggest placing an emphasis on relational algebra, which is used extensively in the chapters on query processing and optimization, as well as in several other chapters. The relational calculi generally do not merit as much emphasis.

Our notation for the tuple relational calculus makes it easy to present the concept of a safe query. The concept of safety for the domain relational calculus, though identical to that for the tuple calculus, is much more cumbersome notationally and requires careful presentation. This consideration may suggest placing somewhat less emphasis on the domain calculus for classes not focusing on database theory.

Exercises

- 6.10** Write the following queries in relational algebra, using the university schema.
- Find the names of all students who have taken at least one Comp. Sci. course.
 - Find the IDs and names of all students who have not taken any course offering before Spring 2009.
 - For each department, find the maximum salary of instructors in that department. You may assume that every department has at least one instructor.
 - Find the lowest, across all departments, of the per-department maximum salary computed by the preceding query.

Answer:

$employee(\underline{person_name}, street, city)$
 $works(\underline{person_name}, company_name, salary)$
 $company(\underline{company_name}, city)$
 $manages(\underline{person_name}, \underline{manager_name})$

Figure 6.22 Relational database for Exercises 6.2, 6.8, 6.11, 6.13, and 6.15

- a. $\Pi_{name}(student \bowtie takes \bowtie \Pi_{course_id}(\sigma_{dept_name = 'Comp.Sci.'}(course)))$
 Note that if we join *student*, *takes*, and *course*, only students from the Comp. Sci. department would be present in the result; students from other departments would be eliminated even if they had taken a Comp. Sci. course since the attribute *dept_name* appears in both *student* and *course*.
- b. $\Pi_{ID, name}(student) - \Pi_{ID, name}(\sigma_{year < 2009}(student \bowtie takes))$ Note that Spring is the first semester of the year, so we do not need to perform a comparison on *semester*.
- c. $dept_name \mathcal{G}_{\max(salary)}(instructor)$
- d. $\mathcal{G}_{\min(maxsal)}(dept_name \mathcal{G}_{\max(salary)} \text{ as } maxsal(instructor))$

6.11 Consider the relational database of Figure 6.22, where the primary keys are underlined. Give an expression in the relational algebra to express each of the following queries:

- a. Find the names of all employees who work for “First Bank Corporation”.
- b. Find the names and cities of residence of all employees who work for “First Bank Corporation”.
- c. Find the names, street addresses, and cities of residence of all employees who work for “First Bank Corporation” and earn more than \$10,000.
- d. Find the names of all employees in this database who live in the same city as the company for which they work.
- e. Assume the companies may be located in several cities. Find all companies located in every city in which “Small Bank Corporation” is located.

Answer:

- a. $\Pi_{person_name}(\sigma_{company_name = \text{“First Bank Corporation”}}(works))$
- b. $\Pi_{person_name, city}(employee \bowtie (\sigma_{company_name = \text{“First Bank Corporation”}}(works)))$

- c. $\Pi_{person_name, street, city}$
 $(\sigma_{(company_name = \text{"First Bank Corporation"} \wedge salary > 10000)}$
 $works \bowtie employee)$
- d. $\Pi_{person_name} (employee \bowtie works \bowtie company)$
- e. Note: Small Bank Corporation will be included in each answer.
 $\Pi_{company_name} (company \div$
 $(\Pi_{city} (\sigma_{company_name = \text{"Small Bank Corporation"} (company))))$
- 6.12 Using the university example, write relational-algebra queries to find the course sections taught by more than one instructor in the following ways:
- Using an aggregate function.
 - Without using any aggregate functions.

Answer:

- $\sigma_{instruct > 1} (course_id, section_id, year, semester \mathcal{G}_{count(*)} \text{ as } instruct (teaches))$
 - $\Pi_{course_id, section_id, year, semester} (\sigma_{ID < > ID2} (takes \bowtie$
 $\rho_{takes1(ID2, course_id, section_id, year, semester)} (takes)))$
- 6.13 Consider the relational database of Figure 6.22. Give a relational-algebra expression for each of the following queries:
- Find the company with the most employees.
 - Find the company with the smallest payroll.
 - Find those companies whose employees earn a higher salary, on average, than the average salary at First Bank Corporation.

Answer:

- $t_1 \leftarrow company_name \mathcal{G}_{count-distinct}(person_name)(works)$
 $t_2 \leftarrow \mathcal{G}_{max}(num_employees)(\rho_{company_strength}(company_name, num_employees)(t_1))$
 $\Pi_{company_name} (\rho_{t_3}(company_name, num_employees)(t_1) \bowtie \rho_{t_4}(num_employees)(t_2))$
- $t_1 \leftarrow company_name \mathcal{G}_{sum}(salary)(works)$
 $t_2 \leftarrow \mathcal{G}_{min}(payroll)(\rho_{company_payroll}(company_name, payroll)(t_1))$
 $\Pi_{company_name} (\rho_{t_3}(company_name, payroll)(t_1) \bowtie \rho_{t_4}(payroll)(t_2))$
- $t_1 \leftarrow company_name \mathcal{G}_{avg}(salary)(works)$
 $t_2 \leftarrow \sigma_{company_name = \text{"First Bank Corporation"}}(t_1)$
 $\Pi_{t_3.company_name} ((\rho_{t_3}(company_name, avg_salary)(t_1))$
 $\bowtie_{t_3.avg_salary > first_bank.avg_salary} (\rho_{first_bank}(company_name, avg_salary)(t_2)))$

- 6.14 Consider the following relational schema for a library:

$member(\underline{memb_no}, name, dob)$
 $books(\underline{isbn}, title, authors, publisher)$
 $borrowed(\underline{memb_no}, \underline{isbn}, date)$

Write the following queries in relational algebra.

- Find the names of members who have borrowed any book published by “McGraw-Hill”.
- Find the name of members who have borrowed all books published by “McGraw-Hill”.
- Find the name and membership number of members who have borrowed more than five different books published by “McGraw-Hill”.
- For each publisher, find the name and membership number of members who have borrowed more than five books of that publisher.
- Find the average number of books borrowed per member. Take into account that if an member does not borrow any books, then that member does not appear in the *borrowed* relation at all.

Answer:

- $$t_1 \leftarrow \Pi_{isbn}(\sigma_{publisher="McGraw-Hill"}(books))$$

$$\Pi_{name}((member \bowtie borrowed) \bowtie t_1)$$
- $$t_1 \leftarrow \Pi_{isbn}(\sigma_{publisher="McGraw-Hill"}(books))$$

$$\Pi_{name, isbn}(member \bowtie borrowed) \div t_1$$
- $$t_1 \leftarrow member \bowtie borrowed \bowtie (\sigma_{publisher="McGraw-Hill"}(books))$$

$$\Pi_{name}(\sigma_{count isbn > 5}((memb_no \text{ } G_{count-distinct(isbn) \text{ as } count isbn}(t_1))))$$
- $$t_1 \leftarrow member \bowtie borrowed \bowtie books$$

$$\Pi_{publisher, name}(\sigma_{count isbn > 5}((publisher, memb_no \text{ } G_{count-distinct(isbn) \text{ as } count isbn}(t_1))))$$

6.15 Consider the employee database of Figure 6.22. Give expressions in tuple relational calculus and domain relational calculus for each of the following queries:

- Find the names of all employees who work for “First Bank Corporation”.
- Find the names and cities of residence of all employees who work for “First Bank Corporation”.
- Find the names, street addresses, and cities of residence of all employees who work for “First Bank Corporation” and earn more than \$10,000.

- d. Find all employees who live in the same city as that in which the company for which they work is located.
- e. Find all employees who live in the same city and on the same street as their managers.
- f. Find all employees in the database who do not work for “First Bank Corporation”.
- g. Find all employees who earn more than every employee of “Small Bank Corporation”.
- h. Assume that the companies may be located in several cities. Find all companies located in every city in which “Small Bank Corporation” is located.

Answer:

- a. Find the names of all employees who work for First Bank Corporation:
 - i. $\{t \mid \exists s \in \text{works} (t[\text{person_name}] = s[\text{person_name}] \wedge s[\text{company_name}] = \text{“First Bank Corporation”})\}$
 - ii. $\{ \langle p \rangle \mid \exists c, s (\langle p, c, s \rangle \in \text{works} \wedge c = \text{“First Bank Corporation”})\}$
- b. Find the names and cities of residence of all employees who work for First Bank Corporation:
 - i. $\{t \mid \exists r \in \text{employee} \exists s \in \text{works} (t[\text{person_name}] = r[\text{person_name}] \wedge t[\text{city}] = r[\text{city}] \wedge r[\text{person_name}] = s[\text{person_name}] \wedge s[\text{company_name}] = \text{“First Bank Corporation”})\}$
 - ii. $\{ \langle p, c \rangle \mid \exists co, sa, st (\langle p, co, sa \rangle \in \text{works} \wedge \langle p, st, c \rangle \in \text{employee} \wedge co = \text{“First Bank Corporation”})\}$
- c. Find the names, street address, and cities of residence of all employees who work for First Bank Corporation and earn more than \$10,000 per annum:
 - i. $\{t \mid t \in \text{employee} \wedge (\exists s \in \text{works} (s[\text{person_name}] = t[\text{person_name}] \wedge s[\text{company_name}] = \text{“First Bank Corporation”} \wedge s[\text{salary}] > 10000))\}$
 - ii. $\{ \langle p, s, c \rangle \mid \langle p, s, c \rangle \in \text{employee} \wedge \exists co, sa (\langle p, co, sa \rangle \in \text{works} \wedge co = \text{“First Bank Corporation”} \wedge sa > 10000)\}$
- d. Find the names of all employees in this database who live in the same city as the company for which they work:

- i. $\{t \mid \exists e \in \text{employee} \exists w \in \text{works} \exists c \in \text{company}$
 $(t[\text{person_name}] = e[\text{person_name}]$
 $\wedge e[\text{person_name}] = w[\text{person_name}]$
 $\wedge w[\text{company_name}] = c[\text{company_name}] \wedge e[\text{city}] =$
 $c[\text{city}])\}$
- ii. $\{ \langle p \rangle \mid \exists st, c, co, sa (\langle p, st, c \rangle \in \text{employee}$
 $\wedge \langle p, co, sa \rangle \in \text{works} \wedge \langle co, c \rangle \in \text{company})\}$
- e. Find the names of all employees who live in the same city and on the same street as do their managers:
 - i. $\{t \mid \exists l \in \text{employee} \exists m \in \text{manages} \exists r \in \text{employee}$
 $(l[\text{person_name}] = m[\text{person_name}] \wedge m[\text{manager_name}] =$
 $r[\text{person_name}]$
 $\wedge l[\text{street}] = r[\text{street}] \wedge l[\text{city}] = r[\text{city}] \wedge t[\text{person_name}] =$
 $l[\text{person_name}])\}$
 - ii. $\{ \langle t \rangle \mid \exists s, c, m (\langle t, s, c \rangle \in \text{employee} \wedge \langle t, m \rangle \in$
 $\text{manages} \wedge \langle m, s, c \rangle \in \text{employee})\}$
- f. Find the names of all employees in this database who do not work for First Bank Corporation:
 If one allows people to appear in the database (e.g. in *employee*) but not appear in *works*, the problem is more complicated. We give solutions for this more realistic case later.
 - i. $\{t \mid \exists w \in \text{works} (w[\text{company_name}] \neq \text{"First Bank Corporation"}$
 $\wedge t[\text{person_name}] = w[\text{person_name}])\}$
 - ii. $\{ \langle p \rangle \mid \exists c, s (\langle p, c, s \rangle \in \text{works} \wedge c \neq \text{"First Bank Corporation"})\}$
 If people may not work for any company:
 - i. $\{t \mid \exists e \in \text{employee} (t[\text{person_name}] = e[\text{person_name}] \wedge \neg \exists w \in$
 works
 $(w[\text{company_name}] = \text{"First Bank Corporation"}$
 $\wedge w[\text{person_name}] = t[\text{person_name}])\}$
 - ii. $\{ \langle p \rangle \mid \exists s, c (\langle p, s, c \rangle \in \text{employee}) \wedge \neg \exists x, y$
 $(y = \text{"First Bank Corporation"} \wedge \langle p, y, x \rangle \in \text{works})\}$
- g. Find the names of all employees who earn more than every employee of Small Bank Corporation:
 - i. $\{t \mid \exists w \in \text{works} (t[\text{person_name}] = w[\text{person_name}] \wedge \forall s \in$
 works
 $(s[\text{company_name}] = \text{"Small Bank Corporation"} \Rightarrow w[\text{salary}] >$
 $s[\text{salary}])\}$

- ii. $\{ \langle p \rangle \mid \exists c, s (\langle p, c, s \rangle \in \text{works} \wedge \forall p_2, c_2, s_2 (\langle p_2, c_2, s_2 \rangle \notin \text{works} \vee c_2 \neq \text{"Small Bank Corporation"} \vee s_2 \neq s)) \}$
- h. Assume the companies may be located in several cities. Find all companies located in every city in which Small Bank Corporation is located.
Note: Small Bank Corporation will be included in each answer.
- i. $\{ t \mid \forall s \in \text{company} (s[\text{company_name}] = \text{"Small Bank Corporation"} \Rightarrow \exists r \in \text{company} (t[\text{company_name}] = r[\text{company_name}] \wedge r[\text{city}] = s[\text{city}])) \}$
- ii. $\{ \langle co \rangle \mid \forall co_2, ci_2 (\langle co_2, ci_2 \rangle \notin \text{company} \vee co_2 \neq \text{"Small Bank Corporation"} \vee \langle co, ci_2 \rangle \in \text{company}) \}$
- 6.16 Let $R = (A, B)$ and $S = (A, C)$, and let $r(R)$ and $s(S)$ be relations. Write relational-algebra expressions equivalent to the following domain-relational-calculus expressions:
- a. $\{ \langle a \rangle \mid \exists b (\langle a, b \rangle \in r \wedge b = 17) \}$
- b. $\{ \langle a, b, c \rangle \mid \langle a, b \rangle \in r \wedge \langle a, c \rangle \in s \}$
- c. $\{ \langle a \rangle \mid \exists b (\langle a, b \rangle \in r) \vee \forall c (\exists d (\langle d, c \rangle \in s) \Rightarrow \langle a, c \rangle \in s) \}$
- d. $\{ \langle a \rangle \mid \exists c (\langle a, c \rangle \in s \wedge \exists b_1, b_2 (\langle a, b_1 \rangle \in r \wedge \langle c, b_2 \rangle \in r \wedge b_1 > b_2)) \}$

Answer:

- a. $\Pi_A (\sigma_{B=17} (r))$
- b. $r \bowtie s$
- c. $\Pi_A(r) \cup (r \div \sigma_B(\Pi_C(s)))$
- d. $\Pi_{r.A} ((r \bowtie s) \bowtie_{c=r2.A \wedge r.B > r2.B} (\rho_{r2}(r)))$
It is interesting to note that (d) is an abstraction of the notorious query "Find all employees who earn more than their manager." Let $R = (\text{emp}, \text{sal})$, $S = (\text{emp}, \text{mgr})$ to observe this.
- 6.17 Repeat Exercise 6.16, writing SQL queries instead of relational-algebra expressions.

Answer:

- a. **select** a
from r
where $b = 17$

- b. **select** a, b, c
from r, s
where $r.a = s.a$
- c. **(select** a
from r)
union
(select a
from s)
- d. **select** a
from r **as** $r1, r$ **as** $r2, s$
where $r1.a = s.a$ **and** $r2.a = s.c$ **and** $r1.b > r2.b$

6.18 Let $R = (A, B)$ and $S = (A, C)$, and let $r(R)$ and $s(S)$ be relations. Using the special constant *null*, write tuple-relational-calculus expressions equivalent to each of the following:

- a. $r \bowtie s$
- b. $r \bowtie s$
- c. $r \bowtie s$

Answer:

- a. $\{t \mid \exists r \in R \exists s \in S (r[A] = s[A] \wedge t[A] = r[A] \wedge t[B] = r[B] \wedge t[C] = s[C]) \vee \exists s \in S (\neg \exists r \in R (r[A] = s[A]) \wedge t[A] = s[A] \wedge t[C] = s[C] \wedge t[B] = \text{null})\}$
- b. $\{t \mid \exists r \in R \exists s \in S (r[A] = s[A] \wedge t[A] = r[A] \wedge t[B] = r[B] \wedge t[C] = s[C]) \vee \exists r \in R (\neg \exists s \in S (r[A] = s[A]) \wedge t[A] = r[A] \wedge t[B] = r[B] \wedge t[C] = \text{null}) \vee \exists s \in S (\neg \exists r \in R (r[A] = s[A]) \wedge t[A] = s[A] \wedge t[C] = s[C] \wedge t[B] = \text{null})\}$
- c. $\{t \mid \exists r \in R \exists s \in S (r[A] = s[A] \wedge t[A] = r[A] \wedge t[B] = r[B] \wedge t[C] = s[C]) \vee \exists r \in R (\neg \exists s \in S (r[A] = s[A]) \wedge t[A] = r[A] \wedge t[B] = r[B] \wedge t[C] = \text{null})\}$

6.19 Give a tuple-relational-calculus expression to find the maximum value in relation $r(A)$.

Answer: $\{ \langle a \rangle \mid \langle a \rangle \in r \wedge \forall \langle b \rangle \in R \langle a \rangle \geq b \}$

CHAPTER 7



Database Design and the E-R Model

This chapter introduces the entity-relationship model in detail. A significant change in the 6th edition is the change in the E-R notation used in the book. There are several alternative E-R notations used in the industry. Increasingly, however, the UML class diagram notation is used instead of the traditional E-R notation used in earlier editions of the book. Among the reasons is the wide availability of UML tools, as well as the conciseness of the UML notation compared to the notation with separate ovals to represent attributes. In keeping with this trend, we have changed our E-R notation to be more compatible with UML.

The chapter covers numerous features of the model, several of which can be omitted depending on the planned coverage of the course. Extended E-R features (Section 7.8) and all subsequent sections may be omitted in case of time constraints in the course, without compromising the students understanding of basic E-R modeling. However, we recommend covering specialization (Section 7.8.1) at least in some detail, since it is widely used in object-oriented modeling.

The E-R model itself and E-R diagrams are used often in the text. It is important that students become comfortable with them. The E-R model is an excellent context for the introduction of students to the complexity of database design. For a given enterprise there are often a wide variety of E-R designs. Although some choices are arbitrary, it is often the case that one design is inherently superior to another. Several of the exercises illustrate this point. The evaluation of the goodness of an E-R design requires an understanding of the enterprise being modeled and the applications to be run. It is often possible to lead students into a debate of the relative merits of competing designs and thus illustrate by example that understanding the application is often the hardest part of database design.

Among the points that are worth discussing when coming up with an E-R design are:

1. Naming of attributes: this is a key aspect of a good design. One approach to design ensures that no two attributes share a name by accident; thus, if ID appears as an attribute of person, it should not appear as an attribute of

another relation, unless it references the ID of person. When natural joins are used in queries, this approach avoids accidental equation of attributes to some extent, although not always; for example, students and instructors share attributes ID and name (presumably inherited from a generalization *person*), so a query that joins the student and instructor relations would equate the respective attribute names.

2. Primary keys: one approach to design creates identifier values for every entity, which are internal to the system and not normally made visible to end users. These internal values are often declared in SQL as **auto increment**, meaning that whenever a tuple is inserted to the relation, a unique value is given to the attribute automatically.

In contrast, the alternative approach, which we have used in this book, avoids creating artificial internal identifiers, and instead uses externally visible attributes as primary key values wherever possible.

As an example, in any university employees and students have externally visible identifiers. These could be used as the primary keys, or alternatively, the application can create identifiers that are not externally visible, and use them as the value for the primary key.

As another example, the *section* table, which has the combination of (*course_id*, *section_id*, *semester*, *year*) as primary key, could instead have a section identifier that is unique across all sections as primary key, with the *course_id*, *section_id*, *semester*, *year* as non-primary key attributes. The difference would be that the relations that refer to *section*, namely *teaches* and *takes*, would have a single unique section id attribute as a foreign key referring to *section*, and would not need to store *course_id*, *section_id*, *semester*, and *year*.

Considerable emphasis is placed on the construction of tables from E-R diagrams. This serves to build intuition for the discussion of the relational model in the subsequent chapters. It also serves to ground abstract concepts of entities and relationships into the more concrete concepts of relations. Several other texts place this material along with the relational data model, rather than in the E-R model chapter. Our motivation for placing this material here is help students to appreciate how E-R data models get used in reality, while studying the E-R model rather than later on.

Exercises

- 7.14 Explain the distinctions among the terms primary key, candidate key, and superkey.

Answer: A *superkey* is a set of one or more attributes that, taken collectively, allows us to identify uniquely an entity in the entity set. A superkey may contain extraneous attributes. If *K* is a superkey, then so is any superset of *K*. A superkey for which no proper subset is also a superkey is called a *candidate key*. It is possible that several distinct sets of attributes could

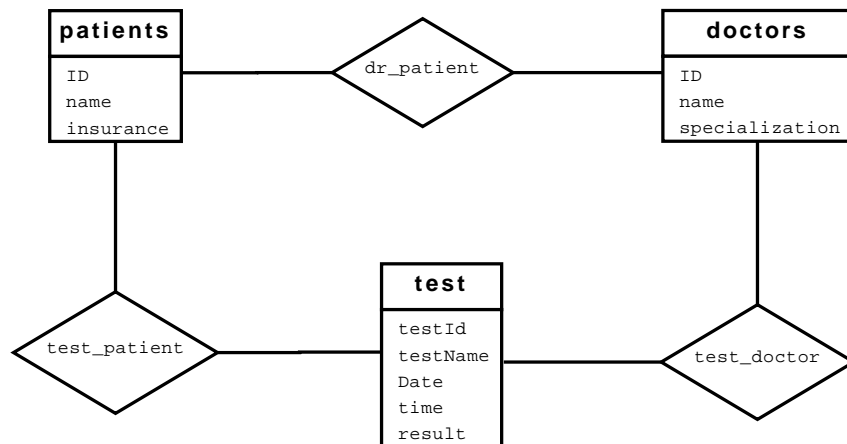


Figure 7.1 E-R diagram for a hospital.

serve as candidate keys. The *primary key* is one of the candidate keys that is chosen by the database designer as the principal means of identifying entities within an entity set.

- 7.15 Construct an E-R diagram for a hospital with a set of patients and a set of medical doctors. Associate with each patient a log of the various tests and examinations conducted.

Answer:

An E diagram for the hospital is shown in Figure 7.1. Although the diagram meets the specifications of the question, a real-world hospital would have many more requirements, such as tracking patient admissions and visits, including which doctor sees a patient on each visit, recording results of tests in a more structured manner, and so on.

- 7.16 Construct appropriate relation schemas for each of the E-R diagrams in Practice Exercises 7.1 to 7.3.

Answer:

- a. Car insurance tables:

customer (customer_id, name, address)

car (license, model)

owns(customer_id, license_no)

accident (report_id, date, place)

participated(license_no, report_id) *policy*(policy_id)

covers(policy_id, license_no)

premium_payment(policy_id, payment_no, due_date, amount, received_on)

Note that a more realistic database design would include details of who was driving the car when an accident happened, and the damage amount for each car that participated in the accident.

- b. Student Exam tables:
- i. Ternary Relationship:

student (*student_id*, *name*, *dept_name*, *tot_cred*)
course (*course_id*, *title*, *credits*)
section (*course_id*, *section_id*, *semester*, *year*)
exam (*exam_id*, *name*, *place*, *time*)
exam_marks (*student_id*, *course_id*, *section_id*, *semester*, *year*, *exam_id*, *marks*)

- ii. Binary relationship:

student (*ID*, *name*, *dept_name*, *tot_cred*)
course (*course_id*, *title*, *credits*)
section (*course_id*, *section_id*, *semester*, *year*)
exam_marks (*student_id*, *course_id*, *sec_id*, *semester*, *year*, *exam_id*, *marks*)

- c. Player Match tables:

match (*match_id*, *date*, *stadium*, *opponent*, *own_score*, *opp_score*)
player (*player_id*, *name*, *age*, *season_score*)
played (*match_id*, *player_id*, *score*)

- 7.17 Extend the E-R diagram of Practice Exercise 7.3 to track the same information for all teams in a league.

Answer: See Figure 7.2. Note that we assume a player can play in only one team; if a player may switch teams, we would have to track for each match which team the player was in, which we could do by turning the relationship *played* into a ternary relationship.

- 7.18 Explain the difference between a weak and a strong entity set.

Answer: A strong entity set has a primary key. All tuples in the set are distinguishable by that key. A weak entity set has no primary key unless attributes of the strong entity set on which it depends are included. Tuples in a weak entity set are partitioned according to their relationship with tuples in a strong entity set. Tuples within each partition are distinguishable by a discriminator, which is a set of attributes.

- 7.19 We can convert any weak entity set to a strong entity set by simply adding appropriate attributes. Why, then, do we have weak entity sets?

Answer: We have weak entities for several reasons:

- We want to avoid the data duplication and consequent possible inconsistencies caused by duplicating the key of the strong entity.
- Weak entities reflect the logical structure of an entity being dependent on another entity.
- Weak entities can be deleted automatically when their strong entity is deleted.

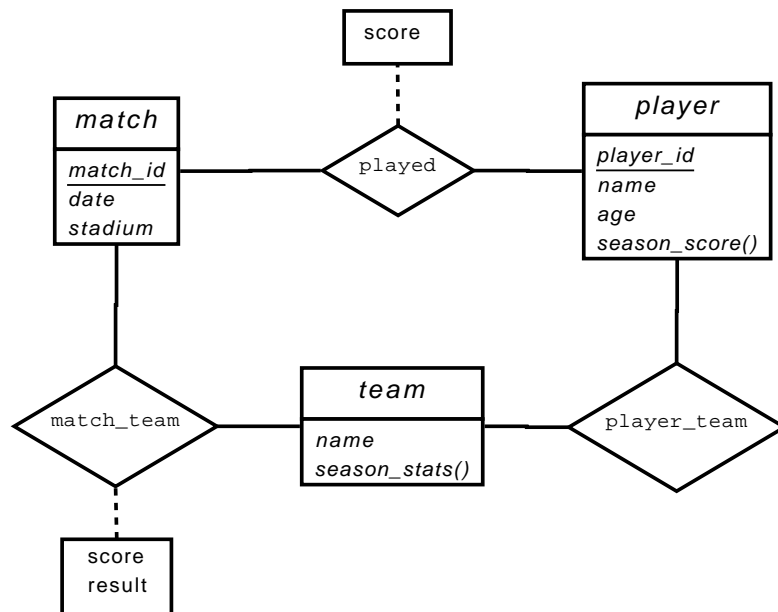


Figure 7.2 E-R diagram for all teams statistics.

- Weak entities can be stored physically with their strong entities.

7.20 Consider the E-R diagram in Figure 7.29, which models an online bookstore.

- List the entity sets and their primary keys.
- Suppose the bookstore adds Blu-ray discs and downloadable video to its collection. The same item may be present in one or both formats, with differing prices. Extend the E-R diagram to model this addition, ignoring the effect on shopping baskets.
- Now extend the E-R diagram, using generalization, to model the case where a shopping basket may contain any combination of books, Blu-ray discs, or downloadable video.

Answer: Interpret the second part of the question as the bookstore adds videos, which may be in Blu-ray disk format or in downloadable format; the same video may be present in both formats.

- Entity sets:

author(name, address, URL)
publisher(name, address, phone, URL)
book(ISBN, title, year, price)
customer(email, name, address, phone)
shopping-basket(basket-id)
warehouse(code, address, phone)

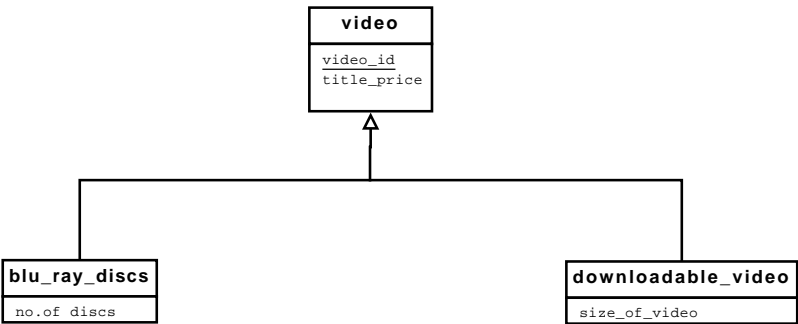


Figure 7.3 ER Diagram for Exercise 7.20 (b)

- b. The ER-diagram portion related to videos is shown in Figure 7.3.
- c. The E-R diagram shown in Figure 7.4 should be added to the E-R diagram of Figure 7.29. Entities that are shown already in Figure 7.29 are shown with only their names, omitting the attributes. The *contains* relationship in Figure 7.29 should be replaced by the version in Figure 7.4. All other parts of Figure 7.29 remain unchanged.

7.21 Design a database for an automobile company to provide to its dealers to assist them in maintaining customer records and dealer inventory and to assist sales staff in ordering cars.

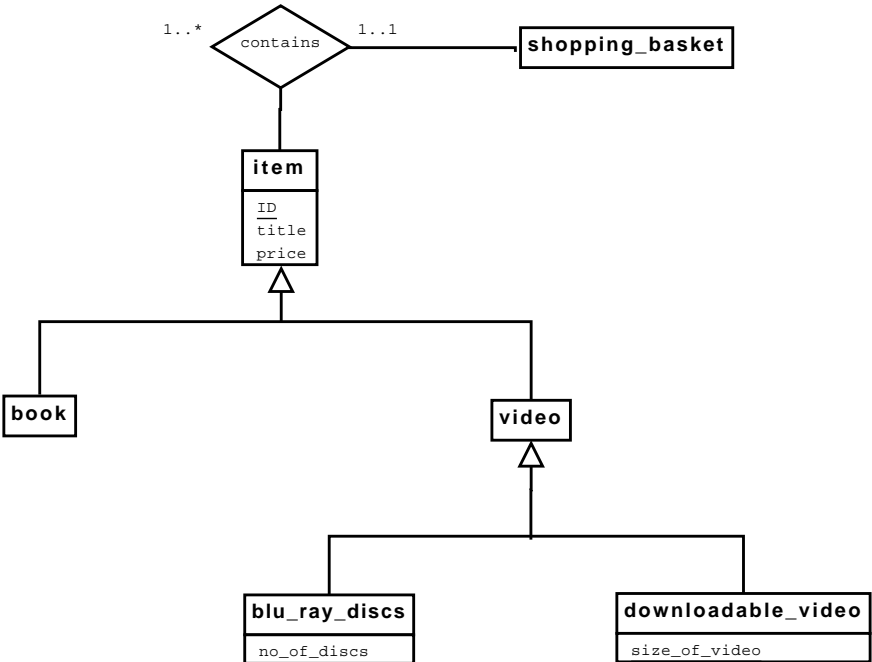


Figure 7.4 ER Diagram for Exercise 7.20 (c)

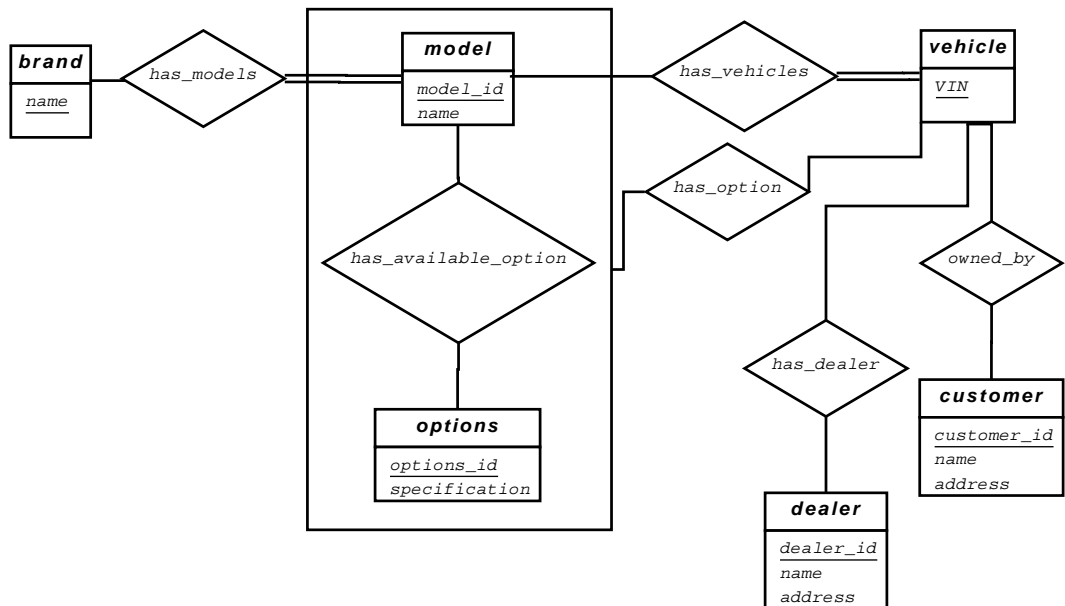


Figure 7.5 ER Diagram for Exercise 7.21

Each vehicle is identified by a vehicle identification number (VIN). Each individual vehicle is a particular model of a particular brand offered by the company (e.g., the XF is a model of the car brand Jaguar of Tata Motors). Each model can be offered with a variety of options, but an individual car may have only some (or none) of the available options. The database needs to store information about models, brands, and options, as well as information about individual dealers, customers, and cars.

Your design should include an E-R diagram, a set of relational schemas, and a list of constraints, including primary-key and foreign-key constraints.

Answer:

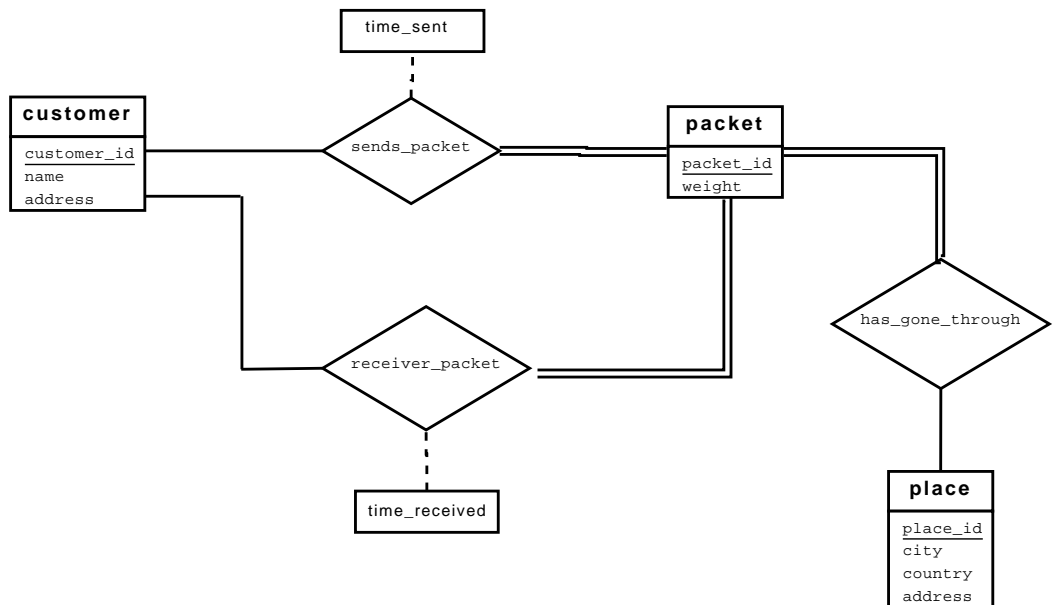
The E-R diagram is shown in Figure 7.5. Note that the *has_option* relationship links a vehicle to an aggregation on the relationship *has_available_option*, instead of directly to the entity set *options*, to ensure that a particular vehicle instance cannot get an option that does not correspond to its model. The alternative of directly linking to *options* is acceptable if ensuring the above integrity constraint is not critical.

The relational schema, along with primary-key and foreign-key constraints is shown below.


```

brand(name)
model(model_id,
      name)
vehicle(VIN)
option(option_id,
       specification)
customer(customer_id,
        name,
        address)
dealer(dealer_id,
       name,
       address)
has_models(name,
           model_id ,
           foreign key name references brand ,
           foreign key model_id references model
)
has_vehicles(model_id,
            VIN,
            foreign key VIN references vehicle,
            foreign key model_id references model
)
available_options(model_id,
                 option_id,
                 foreign key option_id references option,
                 foreign key model_id references model
)
has_options(VIN,
            model_id,
            option_id,
            foreign key VIN references vehicle,
            foreign key (model_id, option_id) references available_options
)
has_dealer(VIN,
           dealer_id ,
           foreign key dealer_id references dealer,
           foreign key VIN references vehicle
)
owned_by(VIN,
         customer_id,
         foreign key customer_id references customer,
         foreign key VIN references vehicle
)

```



As an alternative:

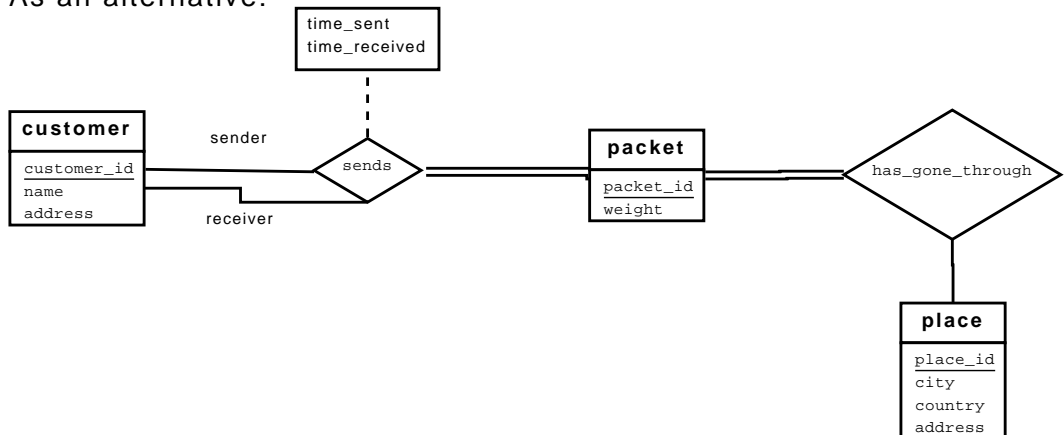


Figure 7.6 ER Diagram Alternatives for Exercise 7.22

- 7.22** Design a database for a world-wide package delivery company (e.g., DHL or FedEx). The database must be able to keep track of customers (who ship items) and customers (who receive items); some customers may do both. Each package must be identifiable and trackable, so the database must be able to store the location of the package and its history of locations. Locations include trucks, planes, airports, and warehouses.

Your design should include an E-R diagram, a set of relational schemas, and a list of constraints, including primary-key and foreign-key constraints.

Answer:

Two alternative E-R diagrams are shown in Figure 7.6. The relational schema, including primary-key and foreign-key constraints, corresponding to the second alternative is shown below.

```

customer(customer_id,
        name,
        address)
packet(packet_id,
       weight)
place(place_id,
      city,
      country,
      address)
sends(sender_id,
      receiver_id,
      packet_id,
      time_received,
      time_sent
      foreign key sender_id references customer,
      foreign key receiver_id references customer,
      foreign key packet_id references packet
)
has_gone_through(
  packet_id,
  place_id
  foreign key packet_id references packet,
  foreign key place_id references place
)

```

- 7.23** Design a database for an airline. The database must keep track of customers and their reservations, flights and their status, seat assignments on individual flights, and the schedule and routing of future flights.

Your design should include an E-R diagram, a set of relational schemas, and a list of constraints, including primary-key and foreign-key constraints.

Answer:

The E-R diagram is shown in Figure 7.7. We assume that the schedule of a flight is fixed across time, although we allow specification of on which days a flight is scheduled. For a particular instance of a flight however we record actual times of departure and arrival. In reality, schedules change with time, so the schedule and routing should be for a particular flight for specified dates, or for a specified range of dates; we ignore this complexity.

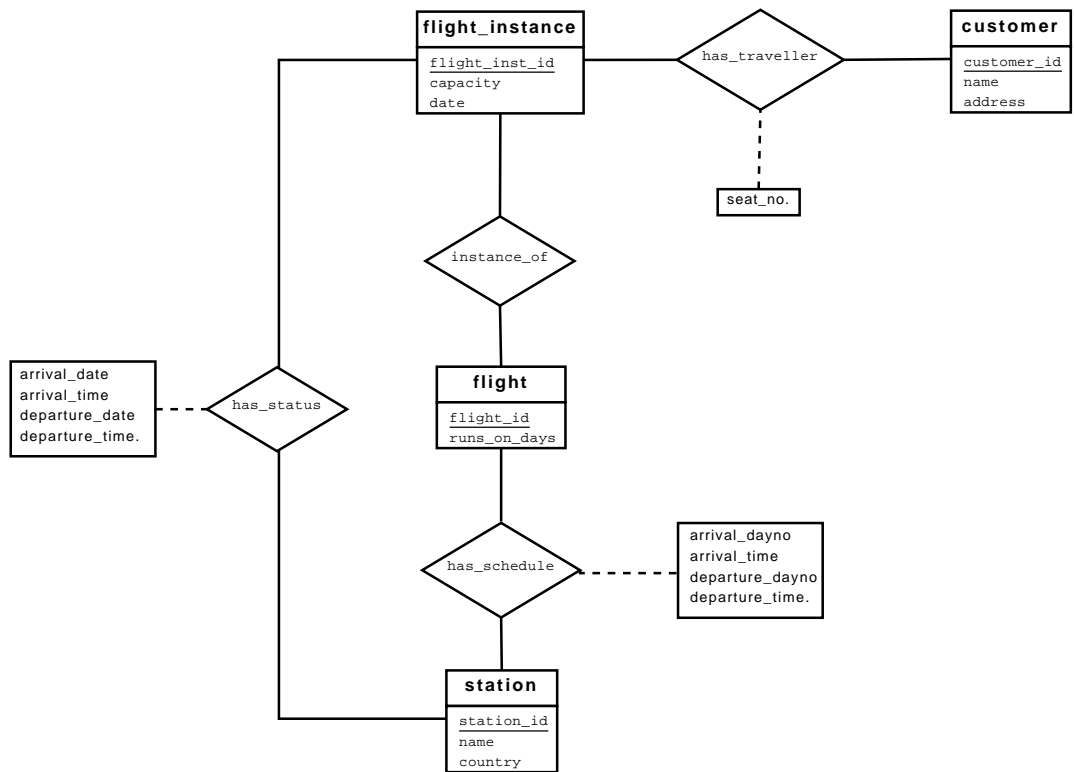


Figure 7.7 ER Diagram for Exercise 7.23

```

flight_instance(flight_inst_id, capacity, date)
customer(customer_id, name, address)
flight(flight_id, runs_on_days)
station(station_id, name, country)
has_traveller(
    flight_inst_id,
    customer_id,
    seat_number,
    foreign key flight_inst_id references flight_instance,
    foreign key customer_id references customer
)
instance_of(
    flight_inst_id,
    flight_id,
    foreign key flight_inst_id references flight_instance,
    foreign key flight_id references flight
)

```

```
has_schedule(  
    flight_id,  
    station_id,  
    order,  
    arrival_dayno,  
    arrival_time,  
    departure_dayno,  
    departure_time,  
    foreign key flight_id references flight,  
    foreign key station_id references station  
)  
has_status(  
    flight_inst_id,  
    station_id,  
    arrival_date,  
    arrival_time,  
    departure_date,  
    departure_time,  
    foreign key flight_inst_id references flight_instance,  
    foreign key station_id references station  
)
```

7.24 In Section 7.7.3, we represented a ternary relationship (repeated in Figure 7.27a) using binary relationships, as shown in Figure 7.27b. Consider the alternative shown in Figure 7.27c. Discuss the relative merits of these two alternative representations of a ternary relationship by binary relationships.

Answer: In the model of Figure 7.27b, there can be instances where E , A , B , C , R_A , R_B and R_C cannot correspond to any instance of A , B , C and R .

The model of Figure 7.27c will not be able to represent all ternary relationships. Consider the ABC relationship set below.

| A | B | C |
|---|---|---|
| 1 | 2 | 3 |
| 4 | 2 | 7 |
| 4 | 8 | 3 |

If ABC is broken into three relationships sets AB , BC and AC , the three will imply that the relation $(4, 2, 3)$ is a part of ABC .

7.25 Consider the relation schemas shown in Section 7.6, which were generated from the E-R diagram in Figure 7.15. For each schema, specify what foreign-key constraints, if any, should be created.

Answer: The foreign-key constraints are as specified below.

```

teaches(
    foreign key ID references instructor,
    foreign key (course_id, sec_id, semester, year) references sec_course
)

takes(
    foreign key ID references student,
    foreign key (course_id, sec_id, semester, year) references sec_course
)

prereq(
    foreign key course_id references course,
    foreign key prereq_id references course
)

advisor(
    foreign key s_ID references student,
    foreign key i_id references instructor
)

sec_course(
    foreign key course_id references course,
    foreign key (sec_id, semester, year) references section
)

sec_time_slot(
    foreign key (course_id, sec_id, semester, year) references sec_course
    foreign key time_slot_id references time_slot
)

sec_class(
    foreign key (course_id, sec_id, semester, year) references sec_course
    foreign key (building, room_number) references classroom
)

inst_dept(
    foreign key ID references instructor
    foreign key dept_name references department
)

stud_dept(
    foreign key ID references student
    foreign key dept_name references department
)

```

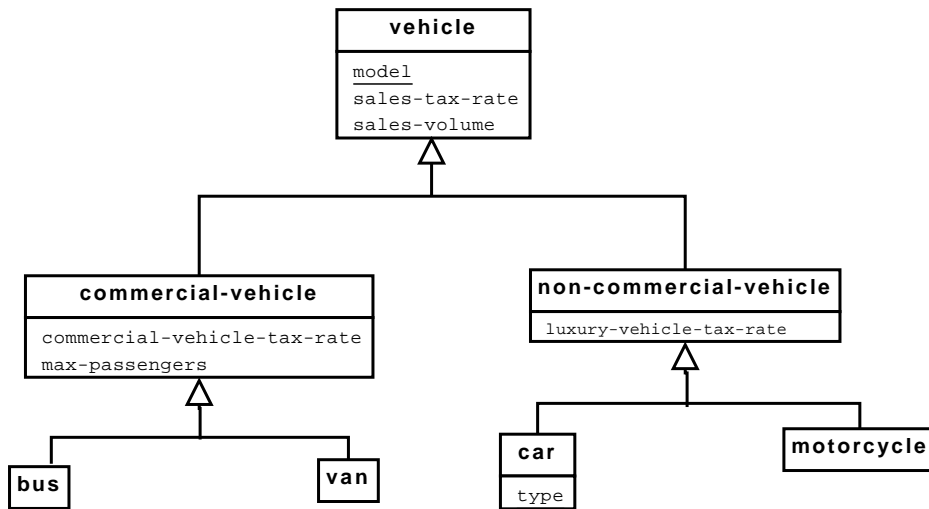


Figure 7.8 E-R diagram of motor-vehicle sales company.

```

course_dept(
    foreign key course_id references course
    foreign key dept_name references department
)

```

- 7.26 Design a generalization–specialization hierarchy for a motor vehicle sales company. The company sells motorcycles, passenger cars, vans, and buses. Justify your placement of attributes at each level of the hierarchy. Explain why they should not be placed at a higher or lower level.

Answer: Figure 7.8 gives one possible hierarchy; note that there could be many alternative solutions. The generalization–specialization hierarchy for the motor-vehicle company is given in the figure. *model*, *sales-tax-rate* and *sales-volume* are attributes necessary for all types of vehicles. Commercial vehicles attract commercial vehicle tax, and each kind of commercial vehicle has a passenger carrying capacity specified for it. Some kinds of non-commercial vehicles attract luxury vehicle tax. Cars alone can be of several types, such as sports-car, sedan, wagon etc., hence the attribute *type*.

- 7.27 Explain the distinction between condition-defined and user-defined constraints. Which of these constraints can the system check automatically? Explain your answer.

Answer: In a generalization–specialization hierarchy, it must be possible to decide which entities are members of which lower level entity sets. In a condition-defined design constraint, membership in the lower level entity-sets is evaluated on the basis of whether or not an entity satisfies an explicit condition or predicate. User-defined lower-level entity sets are not constrained by a membership condition; rather, entities are assigned to a given entity set by the database user.

Condition-defined constraints alone can be automatically handled by the system. Whenever any tuple is inserted into the database, its membership in the various lower level entity-sets can be automatically decided by evaluating the respective membership predicates. Similarly when a tuple is updated, its membership in the various entity sets can be re-evaluated automatically.

- 7.28 Explain the distinction between disjoint and overlapping constraints.

Answer: In a disjointness design constraint, an entity can belong to not more than one lower-level entity set. In overlapping generalizations, the same entity may belong to more than one lower-level entity sets. For example, in the employee-workteam example of the book, a manager may participate in more than one work-team.

- 7.29 Explain the distinction between total and partial constraints.

Answer: In a generalization–specialization hierarchy, a total constraint means that an entity belonging to the higher level entity set must belong to the lower level entity set. A partial constraint means that an entity belonging to the higher level entity set may or may not belong to the lower level entity set.

CHAPTER 8



Relational Database Design

This chapter presents the principles of relational database design. Undergraduates frequently find this chapter difficult. It is acceptable to cover only Sections 8.1 and 8.3 for classes that find the material particularly difficult. However, a careful study of data dependencies and normalization is a good way to introduce students to the formal aspects of relational database theory.

There are many ways of stating the definitions of the normal forms. We have chosen a style which we think is the easiest to present and which most clearly conveys the intuition of the normal forms.

Exercises

- 8.19 Give a lossless-join decomposition into BCNF of schema R of Exercise 8.1.
Answer: From Exercise 8.6, we know that $B \rightarrow D$ is nontrivial and the left hand side is not a superkey. By the algorithm of Figure 8.11 we derive the relations $\{(A, B, C, E), (B, D)\}$. This is in BCNF.
- 8.20 Give a lossless-join, dependency-preserving decomposition into 3NF of schema R of Practice Exercise 8.1.
Answer: First we note that the dependencies given in Practice Exercise 8.1 form a canonical cover. Generating the schema from the algorithm of Figure 8.12 we get

$$R' = \{(A, B, C), (C, D, E), (B, D), (E, A)\}.$$

Schema (A, B, C) contains a candidate key. Therefore R' is a third normal form dependency-preserving lossless-join decomposition.

Note that the original schema $R = (A, B, C, D, E)$ is already in 3NF. Thus, it was not necessary to apply the algorithm as we have done above. The single original schema is trivially a lossless join, dependency-preserving decomposition.

8.21 Normalize the following schema, with given constraints, to 4NF.

```
books(accessionno, isbn, title, author, publisher)
users(userid, name, deptid, deptname)
accessionno  $\rightarrow$  isbn
isbn  $\rightarrow$  title
isbn  $\rightarrow$  publisher
isbn  $\twoheadrightarrow$  author
userid  $\rightarrow$  name
userid  $\rightarrow$  deptid
deptid  $\rightarrow$  deptname
```

Answer: In *books*, we see that

$$isbn \twoheadrightarrow title, publisher, author$$

and yet, *isbn* is not a super key. Thus, we break *books* into

```
books_accnno(accessionno, isbn)
books_details(isbn, title, publisher, author)
```

After this, we still have

$$isbn \twoheadrightarrow author$$

but neither is *isbn* a primary key of *book_details*, nor are the attributes of *book_details* equal to $\{isbn\} \cup \{author\}$. Therefore we decompose *book_details* again into

```
books_details1(isbn, title, publisher)
books_authors(isbn, author)
```

Similarly, in *users*,

$$deptid \rightarrow deptname$$

and yet, *deptid* is not a super key. Hence, we break *users* to

```
users(userid, name, deptid)
departments(deptid, deptname)
```

We verify that there are no further functional or multivalued dependencies that cause violation of 4NF, so the final set of relations are:

```
books_accno(accessionno, isbn)
books_details1(isbn, title, publisher)
books_authors(isbn, author) users(userid, name, deptid)
departments(deptid, deptname)
```

- 8.22 Explain what is meant by *repetition of information* and *inability to represent information*. Explain why each of these properties may indicate a bad relational-database design.

Answer:

- Repetition of information is a condition in a relational database where the values of one attribute are determined by the values of another attribute in the same relation, and both values are repeated throughout the relation. This is a bad relational database design because it increases the storage required for the relation and it makes updating the relation more difficult, and can lead to inconsistent data if updates are done to one copy of the value, but not to another.
- Inability to represent information is a condition where a relationship exists among only a proper subset of the attributes in a relation. This is bad relational database design because all the unrelated attributes must be filled with null values otherwise a tuple without the unrelated information cannot be inserted into the relation.
- Inability to represent information can also occur because of loss of information which results from the decomposition of one relation into two relations, which cannot be combined to recreate the original relation. Such a lossy decomposition may happen implicitly, even without explicitly carrying out decomposition, if the initial relational schema itself corresponds to the decomposition.

- 8.23 Why are certain functional dependencies called *trivial* functional dependencies?

Answer: Certain functional dependencies are called trivial functional dependencies because they are satisfied by all relations.

- 8.24 Use the definition of functional dependency to argue that each of Armstrong's axioms (reflexivity, augmentation, and transitivity) is sound.

Answer: The definition of functional dependency is: $\alpha \rightarrow \beta$ holds on R if in any legal relation $r(R)$, for all pairs of tuples t_1 and t_2 in r such that $t_1[\alpha] = t_2[\alpha]$, it is also the case that $t_1[\beta] = t_2[\beta]$.

Reflexivity rule: if α is a set of attributes, and $\beta \subseteq \alpha$, then $\alpha \rightarrow \beta$.
Assume $\exists t_1$ and t_2 such that $t_1[\alpha] = t_2[\alpha]$

$t_1[\beta] = t_2[\beta]$
 $\alpha \rightarrow \beta$

since $\beta \subseteq \alpha$
definition of FD

Augmentation rule: if $\alpha \rightarrow \beta$, and γ is a set of attributes, then $\gamma \alpha \rightarrow \gamma \beta$.
Assume $\exists t_1, t_2$ such that $t_1[\gamma \alpha] = t_2[\gamma \alpha]$

$t_1[\gamma] = t_2[\gamma]$
 $t_1[\alpha] = t_2[\alpha]$
 $t_1[\beta] = t_2[\beta]$
 $t_1[\gamma \beta] = t_2[\gamma \beta]$
 $\gamma \alpha \rightarrow \gamma \beta$

$\gamma \subseteq \gamma \alpha$
 $\alpha \subseteq \gamma \alpha$
definition of $\alpha \rightarrow \beta$
 $\gamma \beta = \gamma \cup \beta$
definition of FD

Transitivity rule: if $\alpha \rightarrow \beta$ and $\beta \rightarrow \gamma$, then $\alpha \rightarrow \gamma$.
Assume $\exists t_1, t_2$ such that $t_1[\alpha] = t_2[\alpha]$

$t_1[\beta] = t_2[\beta]$
 $t_1[\gamma] = t_2[\gamma]$
 $\alpha \rightarrow \gamma$

definition of $\alpha \rightarrow \beta$
definition of $\beta \rightarrow \gamma$
definition of FD

8.25 Consider the following proposed rule for functional dependencies: If $\alpha \rightarrow \beta$ and $\gamma \rightarrow \beta$, then $\alpha \rightarrow \gamma$. Prove that this rule is *not* sound by showing a relation r that satisfies $\alpha \rightarrow \beta$ and $\gamma \rightarrow \beta$, but does not satisfy $\alpha \rightarrow \gamma$.

Answer: Consider the following rule: if $A \rightarrow B$ and $C \rightarrow B$, then $A \rightarrow C$. That is, $\alpha = A, \beta = B, \gamma = C$. The following relation r is a counterexample to the rule.

r:

| A | B | C |
|----------------|----------------|----------------|
| a ₁ | b ₁ | c ₁ |
| a ₁ | b ₁ | c ₂ |

Note: $A \rightarrow B$ and $C \rightarrow B$, (since no 2 tuples have the same C value, $C \rightarrow B$ is true trivially). However, it is not the case that $A \rightarrow C$ since the same A value is in two tuples, but the C value in those tuples disagree.

8.26 Use Armstrong’s axioms to prove the soundness of the decomposition rule.

Answer: The decomposition rule, and its derivation from Armstrong’s axioms are given below:

if $\alpha \rightarrow \beta\gamma$, then $\alpha \rightarrow \beta$ and $\alpha \rightarrow \gamma$.

$\alpha \rightarrow \beta\gamma$
 $\beta\gamma \rightarrow \beta$
 $\alpha \rightarrow \beta$
 $\beta\gamma \rightarrow \gamma$
 $\alpha \rightarrow \gamma$

given
reflexivity rule
transitivity rule
reflexive rule
transitive rule

8.27 Using the functional dependencies of Practice Exercise 8.6, compute B^+ .

Answer: Computing B^+ by the algorithm in Figure 8.8 we start with $result = \{B\}$. Considering FDs of the form $\beta \rightarrow \gamma$ in F , we find that the only dependencies satisfying $\beta \subseteq result$ are $B \rightarrow B$ and $B \rightarrow D$. Therefore $result = \{B, D\}$. No more dependencies in F apply now. Therefore $B^+ = \{B, D\}$

- 8.28** Show that the following decomposition of the schema R of Practice Exercise 8.1 is not a lossless-join decomposition:

(A, B, C)
 $(C, D, E).$

Hint: Give an example of a relation r on schema R such that

$$\Pi_{A, B, C}(r) \bowtie \Pi_{C, D, E}(r) \neq r$$

Answer: Following the hint, use the following example of r :

| A | B | C | D | E |
|-------|-------|-------|-------|-------|
| a_1 | b_1 | c_1 | d_1 | e_1 |
| a_2 | b_2 | c_1 | d_2 | e_2 |

With $R_1 = (A, B, C)$, $R_2 = (C, D, E)$:

- a. $\Pi_{R_1}(r)$ would be:

| A | B | C |
|-------|-------|-------|
| a_1 | b_1 | c_1 |
| a_2 | b_2 | c_1 |

- b. $\Pi_{R_2}(r)$ would be:

| C | D | E |
|-------|-------|-------|
| c_1 | d_1 | e_1 |
| c_1 | d_2 | e_2 |

- c. $\Pi_{R_1}(r) \bowtie \Pi_{R_2}(r)$ would be:

| A | B | C | D | E |
|-------|-------|-------|-------|-------|
| a_1 | b_1 | c_1 | d_1 | e_1 |
| a_1 | b_1 | c_1 | d_2 | e_2 |
| a_2 | b_2 | c_1 | d_1 | e_1 |
| a_2 | b_2 | c_1 | d_2 | e_2 |

Clearly, $\Pi_{R_1}(r) \bowtie \Pi_{R_2}(r) \neq r$. Therefore, this is a lossy join.

8.29 Consider the following set F of functional dependencies on the relation schema $r(A, B, C, D, E, F)$:

$$\begin{aligned} A &\rightarrow BCD \\ BC &\rightarrow DE \\ B &\rightarrow D \\ D &\rightarrow A \end{aligned}$$

- Compute B^+ .
- Prove (using Armstrong's axioms) that AF is a superkey.
- Compute a canonical cover for the above set of functional dependencies F ; give each step of your derivation with an explanation.
- Give a 3NF decomposition of r based on the canonical cover.
- Give a BCNF decomposition of r using the original set of functional dependencies.
- Can you get the same BCNF decomposition of r as above, using the canonical cover?

Answer:

- $B \rightarrow BD$ (third dependency)
 $BD \rightarrow ABD$ (fourth dependency)
 $ABD \rightarrow ABCD$ (first dependency)
 $ABCD \rightarrow ABCDE$ (second dependency)

Thus, $B^+ = ABCDE$

- Prove (using Armstrong's axioms) that AF is a superkey.

$$\begin{aligned} A &\rightarrow BCD \text{ (Given)} \\ A &\rightarrow ABCD \text{ (Augmentation with A)} \\ BC &\rightarrow DE \text{ (Given)} \\ ABCD &\rightarrow ABCDE \text{ (Augmentation with ABCD)} \\ A &\rightarrow ABCDE \text{ (Transitivity)} \\ AF &\rightarrow ABCDEF \text{ (Augmentation with F)} \end{aligned}$$

- We see that D is extraneous in dep. 1 and 2, because of dep. 3. Removing these two, we get the new set of rules

$$\begin{aligned} A &\rightarrow BC \\ BC &\rightarrow E \\ B &\rightarrow D \\ D &\rightarrow A \end{aligned}$$

Now notice that B^+ is $ABCDE$, and in particular, the FD $B \rightarrow E$ can be determined from this set. Thus, the attribute C is extraneous in

the third dependency. Removing this attribute, and combining with the third FD, we get the final canonical cover as :

$$\begin{aligned} A &\rightarrow BC \\ B &\rightarrow DE \\ D &\rightarrow A \end{aligned}$$

Here, no attribute is extraneous in any FD.

- d. We see that there is no FD in the canonical cover such that the set of attributes is a subset of any other FD in the canonical cover. Thus, each each FD gives rise to its own relation, giving

$$\begin{aligned} r_1(A, B, C) \\ r_2(B, D, E) \\ r_3(D, A) \end{aligned}$$

Now the attribute F is not dependent on any attribute. Thus, it must be a part of every superkey. Also, none of the relations in the above schema have F , and hence, none of them have a superkey. Thus, we need to add a new relation with a superkey.

$$r_4(A, F)$$

- e. We start with

$$r(A, B, C, D, E, F)$$

We see that the relation is not in BCNF because of the first FD. Hence, we decompose it accordingly to get

$$r_1(A, B, C, D) \ r_2(A, E, F)$$

Now we notice that $A \rightarrow E$ is an FD in F^+ , and causes r_2 to violate BCNF. Once again, decomposing r_2 gives

$$r_1(A, B, C, D) \ r_2(A, F) \ r_3(A, E)$$

This schema is now in BCNF.

- f. Can you get the same BCNF decomposition of r as above, using the canonical cover?

If we use the functional dependencies in the preceding canonical cover directly, we cannot get the above decomposition. However, we can infer the original dependencies from the canonical cover, and if we use those for BCNF decomposition, we would be able to get the same decomposition.

- 8.30 List the three design goals for relational databases, and explain why each is desirable.

Answer: The three design goals are lossless-join decompositions, dependency preserving decompositions, and minimization of repetition of

information. They are desirable so we can maintain an accurate database, check correctness of updates quickly, and use the smallest amount of space possible.

- 8.31** In designing a relational database, why might we choose a non-BCNF design?

Answer: BCNF is not always dependency preserving. Therefore, we may want to choose another normal form (specifically, 3NF) in order to make checking dependencies easier during updates. This would avoid joins to check dependencies and increase system performance.

- 8.32** Given the three goals of relational-database design, is there any reason to design a database schema that is in 2NF, but is in no higher-order normal form? (See Practice Exercise 8.17 for the definition of 2NF.)

Answer: The three design goals of relational databases are to avoid

- Repetition of information
- Inability to represent information
- Loss of information.

2NF does not prohibit as much repetition of information since the schema (A, B, C) with dependencies $A \rightarrow B$ and $B \rightarrow C$ is allowed under 2NF, although the same (B, C) pair could be associated with many A values, needlessly duplicating C values. To avoid this we must go to 3NF. Repetition of information is allowed in 3NF in some but not all of the cases where it is allowed in 2NF. Thus, in general, 3NF reduces repetition of information. Since we can always achieve a lossless join 3NF decomposition, there is no loss of information needed in going from 2NF to 3NF.

Note that the decomposition $\{(A, B), (B, C)\}$ is a dependency-preserving and lossless-join 3NF decomposition of the schema (A, B, C) . However, in case we choose this decomposition, retrieving information about the relationship between A, B and C requires a join of two relations, which is avoided in the corresponding 2NF decomposition.

Thus, the decision of which normal form to choose depends upon how the cost of dependency checking compares with the cost of the joins. Usually, the 3NF would be preferred. Dependency checks need to be made with *every* insert or update to the instances of a 2NF schema, whereas, only some queries will require the join of instances of a 3NF schema.

- 8.33** Given a relational schema $r(A, B, C, D)$, does $A \twoheadrightarrow BC$ logically imply $A \twoheadrightarrow B$ and $A \twoheadrightarrow C$? If yes prove it, else give a counter example.

Answer: $A \twoheadrightarrow BC$ holds on the following table:

$r :$

| A | B | C | D |
|-------|-------|-------|-------|
| a_1 | b_1 | c_1 | d_1 |
| a_1 | b_2 | c_2 | d_2 |
| a_1 | b_1 | c_1 | d_2 |
| a_1 | b_2 | c_2 | d_1 |

If $A \twoheadrightarrow B$, then we know that there exists t_1 and t_3 such that $t_1[B] = t_3[B]$. Thus, we must choose one of the following for t_1 and t_3 :

- $t_1 = r_1$ and $t_3 = r_3$, or $t_1 = r_3$ and $t_3 = r_1$:
Choosing either $t_2 = r_2$ or $t_2 = r_4$, $t_3[C] \neq t_2[C]$.
- $t_1 = r_2$ and $t_3 = r_4$, or $t_1 = r_4$ and $t_3 = r_2$:
Choosing either $t_2 = r_1$ or $t_2 = r_3$, $t_3[C] \neq t_2[C]$.

Therefore, the condition $t_3[C] = t_2[C]$ can not be satisfied, so the conjecture is false.

8.34 Explain why 4NF is a normal form more desirable than BCNF.

Answer: 4NF is more desirable than BCNF because it reduces the repetition of information. If we consider a BCNF schema not in 4NF (see Practice Exercise 7.16), we observe that decomposition into 4NF does not lose information provided that a lossless join decomposition is used, yet redundancy is reduced.

CHAPTER 9



Application Design and Development

Exercises

- 9.12 Write a servlet and associated HTML code for the following very simple application: A user is allowed to submit a form containing a value, say n , and should get a response containing n "*" symbols.

Answer:

HTML form

```
<html>
  <head>
    <title>DB Book Exercise 8.8 </title>
  </head>
  <form action="/servlet/StarServlet" method="get">
    Enter the value for "n"
    <br>
    <input type="text" size=5 name="n">
    <input type="submit" value="submit">
  </form>
</html>
```

Servlet Code

```

import java.io.*;
import javax.servlet.*;
import javax.servlet.http.*;

public class StarServlet extends HttpServlet {

    public void doGet(HttpServletRequest request, HttpServletResponse response)
        throws ServletException, IOException
    {
        int n = Integer.parseInt(request.getParameter("n"));
        response.setContentType("text/html");
        PrintWriter out = response.getWriter();
        out.println("<HEAD><TITLE>Exercise 8.8</TITLE> </HEAD>");
        out.println("<BODY>");
        for (int i = 0; i < n; i++) {
            out.print("**");
        }
        out.println("</BODY>");
        out.close();
    }
}

```

- 9.13** Write a servlet and associated HTML code for the following simple application: A user is allowed to submit a form containing a number, say n , and should get a response saying how many times the value n has been submitted previously. The number of times each value has been submitted previously should be stored in a database.

Answer: HTML form

```

<html>
  <head>
    <title>DB Book Exercise 9.13 </title>
  </head>
  <form action="/servlet/KeepCountServlet" method=get>
    Enter the value for "n"
    <br>
    <input type=text size=5 name="n">
    <input type=submit value="submit">
  </form>
</html>

```

Schema

```

CREATE TABLE SUBMISSION.COUNT (
  value integer unique,
  scount integer not null);

```

Servlet Code

```

import java.io.*; import java.sql.*; import java.util.*;
import javax.servlet.*;
import javax.servlet.http.*;

public class KeepCountServlet extends HttpServlet {
    private static final String query =
        "SELECT scount FROM SUBMISSION_COUNT WHERE value=?";

    public void doGet(HttpServletRequest request, HttpServletResponse response)
        throws ServletException, IOException
    {
        int n = Integer.parseInt(request.getParameter("n"));
        int count = 0;
        try {
            Connection conn = getConnection();
            PreparedStatement pstmt = conn.prepareStatement(query);
            pstmt.setInt(1, n);
            ResultSet rs = pstmt.executeQuery();
            if (rs.next()) {
                count = rs.getInt(1);
            }
            pstmt.close();

            Statement stmt = conn.createStatement();
            if (count == 0) {
                stmt.executeUpdate("INSERT INTO SUBMISSION_COUNT VALUES ("
                    + n + ", 1)");
            } else {
                stmt.executeUpdate("UPDATE SUBMISSION_COUNT SET "
                    + "scount=scount+1 WHERE value=" + n);
            }
            stmt.close();
            conn.close();
        }
        catch (Exception e) {
            throw new ServletException(e.getMessage());
        }
        response.setContentType("text/html");
        PrintWriter out = response.getWriter();
        out.println("<HEAD> <TITLE>Exercise 9.13</TITLE> </HEAD>");
        out.println("<BODY>");
        out.println("The value " + n + " has been submitted " + count + " times previously.");
        out.println("</BODY>");
        out.close();
    }
}

```

- 9.14** Write a servlet that authenticates a user (based on user names and passwords stored in a database relation), and sets a session variable called *userid* after authentication.

Answer: HTML form

```
<html>
  <head>
    <title>DB Book Exercise 9.14 </title>
  </head>
  <form action="servlet/SimpleAuthServlet" method=get>
    User Name:
    <input type=text size=20 name="user">
    <BR>
    <BR>
    Password :
    <input type=password size=20 name="passwd">
    <BR>
    <input type=submit value="submit">
  </form>
</html>
```

Schema

```
CREATE TABLE USERAUTH(
  userid integer primary key,
  username varchar(100) unique,
  password varchar(20)
);
```

Servlet Code

```

import java.io.*; import java.sql.*; import java.util.*;
import javax.servlet.*;
import javax.servlet.http.*;

public class SimpleAuthServlet extends HttpServlet {
    private static final String query =
        "SELECT userid, password FROM USERAUTH WHERE username=?";

    public void doGet(HttpServletRequest request, HttpServletResponse response)
        throws ServletException, IOException
    {
        String user = request.getParameter("user");
        String passwd = request.getParameter("passwd");
        String dbPass = null;
        int userId = -1;
        try {
            Connection conn = getConnection();
            PreparedStatement pstmt = conn.prepareStatement(query);
            pstmt.setString(1, user);
            ResultSet rs = pstmt.executeQuery();
            if (rs.next()) {
                userId = rs.getInt(1);
                dbPass = rs.getString(2);
            }
            pstmt.close();
            conn.close();
        }
        catch(Exception e) {
            throw new ServletException(e.getMessage());
        }
        String message;
        if(passwd.equals(dbPass)) {
            message = "Authentication successful";
            getServletContext().setAttribute("userid", new Integer(userId));
        } else {
            message = "Authentication failed! Please check the username " +
                "and password.";
        }

        response.setContentType("text/html");
        PrintWriter out = response.getWriter();
        out.println("<HEAD><TITLE>Exercise 9.14</TITLE></HEAD>");
        out.println("<BODY>");
        out.println(message);
        out.println("</BODY>");
        out.close();
    }
}

```


- 9.15 What is an SQL injection attack? Explain how it works, and what precautions must be taken to prevent SQL injection attacks.

Answer:

SQL injection attack occurs when a malicious user (attacker) manages to get an application to execute an SQL query created by the attacker. If an application constructs an SQL query string by concatenating the user supplied parameters, the application is prone to SQL injection attacks. For example, suppose an application constructs and executes a query to retrieve a user's password in the following way:

```
String userid = request.getParameter("userid");
executeQuery("SELECT password FROM userinfo WHERE userid= ' " + userid + " '");
```

Now, if a user types the value for the parameter as:

john' OR userid= 'admin

the query constructed will be:

```
SELECT password FROM userinfo WHERE userid='john' OR userid='admin';
```

This can reveal unauthorized information to the attacker.

Prevention:

Use prepared statements, with any value that is taken as user input (not just text fields, but even options in drop-down menus) passed as a parameter; user input should **never** be concatenated directly into a query string. The JDBC, ODBC, ADO.NET, or other libraries that provide prepared statements ensure that special characters like quotes are escaped as appropriate for the target database, so that SQL injection attempts will fail.

- 9.16 Write pseudocode to manage a connection pool. Your pseudocode must include a function to create a pool (providing a database connection string, database user name, and password as parameters), a function to request a connection from the pool, a function to release a connection to the pool, and a function to close the connection pool.

Answer:

```
// This connection pool manager is NOT thread-safe.

INITIAL_POOL_SIZE = 20;
POOL_SIZE_INCREMENT = 5;
MAX_POOL_SIZE = 100;
Queue freeConnections = empty queue;
Queue activeConnections = empty queue;
String poolConnURL;
String poolConnUserId;
String poolConnPasswd;

createPool(connString, userid, password) {
    poolConnURL = connString;
    poolConnUserId = userid;
    poolConnPasswd = password;
    for (i = 0; i < INITIAL_POOL_SIZE; i++) {
        conn = createConnection(connString, userid, password);
        freeConnections.add(conn);
    }
}

Connection getConnection() {
    if(freeConnections.size() != 0){
        conn = freeConnections.remove();
        activeConnections.add(conn);
        return conn;
    }
    activeConns = activeConnections.size();
    if (activeConns == MAX_POOL_SIZE)
        ERROR("Max pool size reached");
    if (MAX_POOL_SIZE - activeConns > POOL_SIZE_INCREMENT)
        connsToCreate = POOL_SIZE_INCREMENT;
    else
        connsToCreate = MAX_POOL_SIZE - activeConns;

    for (i = 0; i < connsToCreate; i++) {
        conn = createConnection(poolConnURL, poolConnUserId,
                                poolConnPasswd);
        freeConnections.add(conn);
    }
    return getConnection();
}

releaseConnection(conn) {
    activeConnections.remove(conn);
    freeConnections.add(conn);
}
```

```

closePool() {
    if(activeConnections.size() != 0)
        WARNING("Connections active. Will force close.");
    for (i=0; i < freeConnections.size(); i++) {
        conn = freeConnections.elementAt(i);
        freeConnections.removeElementAt(i);
        conn.close();
    }

    for (i=0; i < activeConnections.size(); i++) {
        conn = activeConnections.elementAt(i);
        activeConnections.removeElementAt(i);
        conn.close();
    }
}

```

9.17 Explain the terms CRUD and REST.

Answer:

The term CRUD refers to simple user interfaces to a relation (or an object model), that provide the ability to Create tuples, Read tuples, Update tuples and Delete tuples (or objects).

In Representation State Transfer (or REST), Web service function calls are executed by a standard HTTP request to a URL at an application server, with parameters sent as standard HTTP request parameters. The application server executes the request (which may involve updating the database at the server), generates and encodes the result, and returns the result as the result of the HTTP request. The server can use any encoding for a particular requested URL; XML and the JavaScript Object Notation (JSON) encoding are widely used. The requestor parses the returned page to access the returned data.

9.18 Many Web sites today provide rich user-interfaces using Ajax. List two features each of which reveals if a site uses Ajax, without having to look at the source code. Using the above features, find three sites which use Ajax; you can view the HTML source of the page to check if the site is actually using Ajax.

Answer:

- a. A Web site with a form that allows you to select a value from one menu (e.g. country), and once a value has been selected, you are allowed to select a value from another menu (e.g. state from a list of states in that country) with the list of values for the second menu (e.g. state) empty until the first value (e.g. country) is selected, probably uses Ajax. If the number of countries is small, the site may well send all country-state pairs ahead of time, and then simply use Javascript without Ajax; however, if you notice a small delay in populating the second menu, the site probably uses Ajax.

- b. A Web site that supports autocompletion for text that you are typing almost surely uses Ajax. For example, a search Web site that suggests possible completions of your query as you type the query in, or a Web-based email site that suggests possible completions of an email address as you type in the address almost surely use Ajax to communicate with a server after you type in each character (sometimes starting after the 3rd or 4th character), and respond with possible completions.
- c. A Web form that, on filling in one piece of data, such as your email address or employee code, fills in other fields such as your name and contact information, without refreshing the page, almost surely uses Ajax to retrieve required information using the information (such as the email address or employee code) provided by the user.

Popular Web sites that use Ajax include almost all current generation Web email interfaces (such as GMail, Yahoo! mail, or Windows Live mail), and almost all search engines, which provide autocompletion. Online document management systems such as Google Docs or Office Live use Ajax extensively to push your updates to the server, and to fetch concurrent updates (to different parts of the document or spreadsheet) transparently. Check your organizations Web applications to find more local examples.

9.19 XSS attacks:

- a. What is an XSS attack?
- b. How can the referer field be used to detect some XSS attacks?

Answer:

- a. In an XSS attack, a malicious user enters code written in a client-side scripting language such as JavaScript or Flash instead of entering a valid name or comment. When a different user views the entered text, the browser would execute the script, which can carry out actions such as sending private cookie information back to the malicious user, or execute an action on a different Web site, such as a bank Web site, that the user may be logged into.
- b. The HTTP protocol allows a server to check the **referer** of a page access, that is, the URL of the page that had the link that the user clicked on to initiate the page access. By checking that the referer is valid, for example, that the referer URL is a page on the same Web site (say the bank Web site in the previous example), XSS attacks that originated on a different Web site accessed by the user can be prevented. The referer field is set by the browser, so a malicious or compromised browser can spoof the referer field, but a basic XSS attack can be detected and prevented.

- 9.20 What is multi-factor authentication? How does it help safeguard against stolen passwords?

Answer: In multi-factor authentication (with two-factor authentication as a special case), where multiple independent *factors* (that is, pieces of information or processes) are used to identify a user. The factors should not share a common vulnerability; for example, if a system merely required two passwords, both could be vulnerable to leakage in the same manner (by network sniffing, or by a virus on the computer used by the user, for example). In addition to secret passwords, which serve as one factor, one-time passwords sent by SMS to a mobile phone, or a hardware token that generates a (time-based) sequence of one-time passwords, are examples of extra factors.

- 9.21 Consider the Oracle **Virtual Private Database** (VPD) feature described in Section 9.7.5, and an application based on our university schema.

- What predicate (using a subquery) should be generated to allow each faculty member to see only *takes* tuples corresponding to course sections that they have taught?
- Give an SQL query such that the query with the predicate added gives a result that is a subset of the original query result without the added predicate.
- Give an SQL query such that the query with the predicate added gives a result containing a tuple that is not in the result of the original query without the added predicate.

Answer:

- The following predicate can be added to queries on *takes*, to ensure that each faculty member only sees *takes* tuples corresponding to course sections that they have taught; the predicate assumes that *syscontext.user_id()* returns the instructor identifier.

```
exists (select *
from teaches
where teaches.ID=syscontext.user_name() and
teaches.course_id= takes.course_id and
teaches.section_id= takes.section_id and
teaches.year= takes.year and
teaches.semester= takes.semester)
```

- The following query retrieves a subset of the answers, due to the above predicate:

```
select * from takes;
```

- The following query gives a result tuple that is not in the result of the original query:

```
select count(*) from takes;
```

The aggregated function above can be any of the aggregate functions, such as sum, average, min or max on any attribute; in the case of min or max the result could be the same if the person executing the query is authorized to see the tuple corresponding to the min or max value. For count, sum, and average, the values are likely to be different as long as there is more than one section.

9.22 What are two advantages of encrypting data stored in the database?

Answer:

- Unauthorized users who gain access to the OS files in which the DBMS stores the data cannot read the data.
- If the application encrypts the data before it reaches the database, it is possible to ensure privacy for the user's data such that even privileged users like database administrators cannot access other users' data.

9.23 Suppose you wish to create an audit trail of changes to the *takes* relation.

- Define triggers to create an audit trail, logging the information into a relation called, for example, *takes_trail*. The logged information should include the user-id (assume a function *user_id()* provides this information) and a timestamp, in addition to old and new values. You must also provide the schema of the *takes_trail* relation.
- Can the above implementation guarantee that updates made by a malicious database administrator (or someone who manages to get the administrator's password) will be in the audit trail? Explain your answer.

Answer:

- Schema for the *takes_trail* table**

```
takes_trail(user_id integer, timestamp datetime, operation_code integer,
            new_ID, new_course_id, new_sec_id, new_year, new_sem, new_grade
            old_ID, old_course_id, old_sec_id, old_year, old_sem, old_grade)
```

Trigger for INSERT

```
create trigger takes_insert after insert on takes
referencing new row as nrow
for each row
begin
    insert into takes_trail values (user_id(), systime(), 1,
                                   nrow.ID, nrow.course_id, nrow.sec_id, nrow.year, nrow.sem, nrow.grade
                                   null, null, null, null, null, null);
end
```

Trigger for UPDATE

```

create trigger takes_update after update on takes
referencing old row as orow, referencing new row as nrow
for each row
begin
    insert into takes_trail values (user_id(), systime(), 2,
                                   nrow.ID, nrow.course_id, nrow.sec_id, nrow.year, nrow.sem, nrow.grade
                                   orow.ID, orow.course_id, orow.sec_id, orow.year, orow.sem, orow.grade);
end

```

Trigger for DELETE

```

create trigger takes_delete after delete on takes
referencing old row as orow
for each row
begin
    insert into account_trail values (user_id(), systime(), 3,
                                     null, null, null, null, null, null);
    orow.ID, orow.course_id, orow.sec_id, orow.year, orow.sem, orow.grade);
end

```

- b. No. Someone who has the administrator privileges can disable the trigger and thus bypass the trigger based audit trail.

9.24 Hackers may be able to fool you into believing that their Web site is actually a Web site (such as a bank or credit card Web site) that you trust. This may be done by misleading email, or even by breaking into the network infrastructure and rerouting network traffic destined for, say mybank.com, to the hacker's site. If you enter your user name and password on the hacker's site, the site can record it, and use it later to break into your account at the real site. When you use a URL such as https://mybank.com, the HTTPS protocol is used to prevent such attacks. Explain how the protocol might use digital certificates to verify authenticity of the site.

Answer: In the HTTPS protocol, a Web site first sends a digital certificate to the user's browser. The browser decrypts the digital certificate using the stored public key of the trusted certification authority and displays the site's name from the decrypted message. The user can then verify if the site name matches the one he/she intended to visit (in this example mybank.com) and accept the certificate. The browser then uses the site's public key (that is part of the digital certificate) to encrypt user's data. Note that it is possible for a malicious user to gain access to the digital certificate of mybank.com, but since the user's data (such as password) is encrypted using the public key of mybank.com, the malicious user will not be able to decrypt the data.

- 9.25 Explain what is a challenge–response system for authentication. Why is it more secure than a traditional password-based system?

Answer:

In a challenge-response system, a secret password is issued to the user and is also stored on the database system. When a user has to be authenticated, the database system sends a challenge string to the user. The user encrypts the challenge string using his/her secret password and returns the result. The database system can verify the authenticity of the user by decrypting the string with the same secret password and checking the result with the original challenge string.

The challenge-response system is more secure than a traditional password-based system since the password is not transferred over the network during authentication.

CHAPTER 10



Storage and File Structure

This chapter presents basic file structure concepts. The chapter really consists of two parts—the first dealing with relational databases, and the second with object-oriented databases. The second part can be omitted without loss of continuity for later chapters.

Many computer science undergraduates have covered some of the material in this chapter in a prior course on data structures or on file structures. Even if students' backgrounds are primarily in data structures, this chapter is still important since it addresses data structure issues as they pertain to disk storage. Buffer management issues, covered in Section 10.8.1 should be familiar to students who have taken an operating systems course. However, there are database-specific aspects of buffer management that make this section worthwhile even for students with an operating system background.

Exercises

10.10 List the physical storage media available on the computers you use routinely. Give the speed with which data can be accessed on each medium.

Answer: Your answer will be based on the computers and storage media that you use. Typical examples would be hard disk, CD and DVD disks, and flash memory in the form of USB keys, memory cards or solid state disks.

The following table shows the typical transfer speeds for the above mentioned storage media, as of early 2010.

| Storage Media | Speed (in MB/s) |
|-------------------|-----------------|
| CD Drive | 8 |
| DVD Drive | 20 |
| USB Keys | 30 |
| Memory Cards | 1 - 40 |
| Hard Disk | 100 |
| Solid State Disks | > 100 |

Note that speeds of flash memory cards can vary significantly, with some low end cards giving low transfer speeds, although better ones give much higher transfer speeds.

- 10.11** How does the remapping of bad sectors by disk controllers affect data-retrieval rates?

Answer: Remapping of bad sectors by disk controllers does reduce data retrieval rates because of the loss of sequentiality amongst the sectors. But that is better than the loss of data in case of no remapping!

- 10.12** RAID systems typically allow you to replace failed disks without stopping access to the system. Thus, the data in the failed disk must be rebuilt and written to the replacement disk while the system is in operation. Which of the RAID levels yields the least amount of interference between the rebuild and ongoing disk accesses? Explain your answer.

Answer: RAID level 1 (mirroring) is the one which facilitates rebuilding of a failed disk with minimum interference with the on-going disk accesses. This is because rebuilding in this case involves copying data from just the failed disk's mirror. In the other RAID levels, rebuilding involves reading the entire contents of all the other disks.

- 10.13** What is scrubbing, in the context of RAID systems, and why is scrubbing important?

Answer: Successfully written sectors which are subsequently damaged, but where the damage has not been detected, are referred to as latent sector errors. In RAID systems, latent errors can lead to data loss even on a single disk failure, if the latent error exists on one of the other disks. Disk scrubbing is a background process that reads disk sectors during idle periods, with the goal of detecting latent sector errors. If a sector error is found, the sector can either be rewritten if the media has not been damaged, or remapped to a spare sector in the disk. The data in the sector can be recovered from the other disks in the RAID array.

- 10.14** In the variable-length record representation, a null bitmap is used to indicate if an attribute has the null value.

- a. For variable length fields, if the value is null, what would be stored in the offset and length fields?
- b. In some applications, tuples have a very large number of attributes, most of which are null. Can you modify the record representation such that the only overhead for a null attribute is the single bit in the null bitmap.

Answer:

- a. It does not matter on what we store in the offset and length fields since we are using a null bitmap to identify null entries. But it would make sense to set the offset and length to zero to avoid having arbitrary values.

- b. We should be able to locate the null bitmap and the offset and length values of non-null attributes using the null bitmap. This can be done by storing the null bitmap at the beginning and then for non-null attributes, store the value (for fixed size attributes), or offset and length values (for variable sized attributes) in the same order as in the bitmap, followed by the values for non-null variable sized attributes. This representation is space efficient but needs extra work to retrieve the attributes.

10.15 Explain why the allocation of records to blocks affects database-system performance significantly.

Answer: If we allocate related records to blocks, we can often retrieve most, or all, of the requested records by a query with one disk access. Disk accesses tend to be the bottlenecks in databases; since this allocation strategy reduces the number of disk accesses for a given operation, it significantly improves performance.

10.16 If possible, determine the buffer-management strategy used by the operating system running on your local computer system and what mechanisms it provides to control replacement of pages. Discuss how the control on replacement that it provides would be useful for the implementation of database systems.

Answer: The typical OS uses variants of LRU, which are cheaper to implement than LRU, for buffer replacement. LRU and its variants are often a bad strategy for databases. As explained in Section 10.8.2 of the text, MRU is the best strategy for nested loop join. In general no single strategy handles all scenarios well, and the database system should be able to manage its own buffer cache for which the replacement policy takes into account all the performance related issues.

Many operating systems provide mechanisms to lock pages in memory, which can be used to ensure buffer pages stay in memory. However, operating systems today generally do not allow any other control on page replacement.

10.17 List two advantages and two disadvantages of each of the following strategies for storing a relational database:

- a. Store each relation in one file.
- b. Store multiple relations (perhaps even the entire database) in one file.

Answer:

- a. Advantages of storing a relation as a file include using the file system provided by the OS, thus simplifying the DBMS, but incurs the disadvantage of restricting the ability of the DBMS to increase performance by using more sophisticated storage structures.

- b. By using one file for the entire database, these complex structures can be implemented through the DBMS, but this increases the size and complexity of the DBMS.

10.18 In the sequential file organization, why is an overflow *block* used even if there is, at the moment, only one overflow record?

Answer: An overflow block is used in sequential file organization because a block is the smallest space which can be read from a disk. Therefore, using any smaller region would not be useful from a performance standpoint. The space saved by allocating disk storage in record units would be overshadowed by the performance cost of allowing blocks to contain records of multiple files.

10.19 Give a normalized version of the *Index_metadata* relation, and explain why using the normalized version would result in worse performance.

Answer: The *Index_metadata* relation can be normalized as follows

Index_metadata(*index_name*, *relation_name*, *index_type*)
Index_Attrib_metadata (*index_name*, *position*, *attribute_name*)

The normalized version will require extra disk accesses to read *Index_Attrib_metadata* everytime an index has to be accessed. Thus, it will lead to worse performance.

10.20 If you have data that should not be lost on disk failure, and the data are write intensive, how would you store the data?

Answer: A **RAID** array can handle the failure of a single drive (two drives in the case of RAID 6) without data loss, and is relatively inexpensive. There are several RAID alternatives, each with different performance and cost implications. For write intensive data with mostly sequential writes, RAID 1 and RAID 5 will both perform well, but with less storage overhead for RAID 5. If writes are random, RAID 1 is preferred, since a random block write requires multiple reads and writes in RAID 5.

10.21 In earlier generation disks the number of sectors per track was the same across all tracks. Current generation disks have more sectors per track on outer tracks, and fewer sectors per track on inner tracks (since they are shorter in length). What is the effect of such a change on each of the three main indicators of disk speed?

Answer: The main performance effect of storing more sectors on the outer tracks and fewer sectors on the inner tracks is that the disk's **data-transfer** rate will be greater on the outer tracks than the inner tracks. This is because the disk spins at a constant rate, so more sectors pass underneath the drive head in a given amount of time when the arm is positioned on an outer track than when on an inner track.

In fact, some high-performance systems are optimized by storing data only in outer tracks, so that disk arm movement is minimized while maximizing data-transfer rates.

- 10.22** Standard buffer managers assume each block is of the same size and costs the same to read. Consider a buffer manager that, instead of LRU, uses the rate of reference to objects, that is, how often an object has been accessed in the last n seconds. Suppose we want to store in the buffer objects of varying sizes, and varying read costs (such as Web pages, whose read cost depends on the site from which they are fetched). Suggest how a buffer manager may choose which block to evict from the buffer.

Answer: A good solution would make use of a *priority queue* to evict pages, where the priority (p) is ordered by the *expected cost* of re-reading a page given it's past access frequency (f) in the last n seconds, it's re-read cost (c), and its size s :

$$p = f * c / s$$

The buffer manager should choose to evict pages with the lowest value of p , until there is enough free space to read in a newly referenced object.

CHAPTER 11



Indexing and Hashing

This chapter covers indexing techniques ranging from the most basic one to highly specialized ones. Due to the extensive use of indices in database systems, this chapter constitutes an important part of a database course.

A class that has already had a course on data-structures would likely be familiar with hashing and perhaps even B⁺-trees. However, this chapter is necessary reading even for those students since data structures courses typically cover indexing in main memory. Although the concepts carry over to database access methods, the details (e.g., block-sized nodes), will be new to such students.

The sections on B-trees (Sections 11.4.5) and bitmap indexing (Section 11.9) may be omitted if desired.

Exercises

- 11.15 When is it preferable to use a dense index rather than a sparse index? Explain your answer.

Answer: It is preferable to use a dense index instead of a sparse index when the file is not sorted on the indexed field (such as when the index is a secondary index) or when the index file is small compared to the size of memory.

- 11.16 What is the difference between a clustering index and a secondary index?

Answer: The clustering index is on the field which specifies the sequential order of the file. There can be only one clustering index while there can be many secondary indices.

- 11.17 For each B⁺-tree of Practice Exercise 11.3, show the steps involved in the following queries:

- Find records with a search-key value of 11.
- Find records with a search-key value between 7 and 17, inclusive.

Answer: With the structure provided by the solution to Practice Exercise 11.3a:

- a. Find records with a value of 11
 - i. Search the first level index; follow the first pointer.
 - ii. Search next level; follow the third pointer.
 - iii. Search leaf node; follow first pointer to records with key value 11.
- b. Find records with value between 7 and 17 (inclusive)
 - i. Search top index; follow first pointer.
 - ii. Search next level; follow second pointer.
 - iii. Search third level; follow second pointer to records with key value 7, and after accessing them, return to leaf node.
 - iv. Follow fourth pointer to next leaf block in the chain.
 - v. Follow first pointer to records with key value 11, then return.
 - vi. Follow second pointer to records with with key value 17.

With the structure provided by the solution to Practice Exercise 12.3b:

- a. Find records with a value of 11
 - i. Search top level; follow second pointer.
 - ii. Search next level; follow second pointer to records with key value 11.
- b. Find records with value between 7 and 17 (inclusive)
 - i. Search top level; follow second pointer.
 - ii. Search next level; follow first pointer to records with key value 7, then return.
 - iii. Follow second pointer to records with key value 11, then return.
 - iv. Follow third pointer to records with key value 17.

With the structure provided by the solution to Practice Exercise 12.3c:

- a. Find records with a value of 11
 - i. Search top level; follow second pointer.
 - ii. Search next level; follow first pointer to records with key value 11.
- b. Find records with value between 7 and 17 (inclusive)
 - i. Search top level; follow first pointer.
 - ii. Search next level; follow fourth pointer to records with key value 7, then return.
 - iii. Follow eighth pointer to next leaf block in chain.
 - iv. Follow first pointer to records with key value 11, then return.
 - v. Follow second pointer to records with key value 17.

- 11.18** The solution presented in Section 11.3.4 to deal with nonunique search keys added an extra attribute to the search key. What effect could this change have on the height of the B⁺-tree?

Answer: The resultant B-tree's extended search key is unique. This results in more number of nodes. A single node (which points to multiple records with the same key) in the original tree may correspond to multiple nodes in the result tree. Depending on how they are organized the height of the tree may increase; it might be more than that of the original tree.

- 11.19** Explain the distinction between closed and open hashing. Discuss the relative merits of each technique in database applications.

Answer: Open hashing may place keys with the same hash function value in different buckets. Closed hashing always places such keys together in the same bucket. Thus in this case, different buckets can be of different sizes, though the implementation may be by linking together fixed size buckets using overflow chains. Deletion is difficult with open hashing as *all* the buckets may have to be inspected before we can ascertain that a key value has been deleted, whereas in closed hashing only that bucket whose address is obtained by hashing the key value need be inspected. Deletions are more common in databases and hence closed hashing is more appropriate for them. For a small, static set of data lookups may be more efficient using open hashing. The symbol table of a compiler would be a good example.

- 11.20** What are the causes of bucket overflow in a hash file organization? What can be done to reduce the occurrence of bucket overflows?

Answer: The causes of bucket overflow are :-

- Our estimate of the number of records that the relation will have was too low, and hence the number of buckets allotted was not sufficient.
- Skew in the distribution of records to buckets. This may happen either because there are many records with the same search key value, or because the hash function chosen did not have the desirable properties of uniformity and randomness.

To reduce the occurrence of overflows, we can :-

- Choose the hash function more carefully, and make better estimates of the relation size.
- If the estimated size of the relation is n_r and number of records per block is f_r , allocate $(n_r/f_r) * (1+d)$ buckets instead of (n_r/f_r) buckets. Here d is a fudge factor, typically around 0.2. Some space is wasted: About 20 percent of the space in the buckets will be empty. But the benefit is that some of the skew is handled and the probability of overflow is reduced.

- 11.21** Why is a hash structure not the best choice for a search key on which range queries are likely?

Answer: A range query cannot be answered efficiently using a hash index, we will have to read all the buckets. This is because key values in the range do not occupy consecutive locations in the buckets, they are distributed uniformly and randomly throughout all the buckets.

11.22 Suppose there is a relation $r(A, B, C)$, with a B^+ -tree index with search key (A, B) .

- What is the worst-case cost of finding records satisfying $10 < A < 50$ using this index, in terms of the number of records retrieved n_1 and the height h of the tree?
- What is the worst-case cost of finding records satisfying $10 < A < 50 \wedge 5 < B < 10$ using this index, in terms of the number of records n_2 that satisfy this selection, as well as n_1 and h defined above?
- Under what conditions on n_1 and n_2 would the index be an efficient way of finding records satisfying $10 < A < 50 \wedge 5 < B < 10$?

Answer:

- What is the worst case cost of finding records satisfying $10 < A < 50$ using this index, in terms of the number of records retrieved n_1 and the height h of the tree?
This query does not correspond to a range query on the search key as the condition on the first attribute if the search key is a comparison condition. It looks up records which have the value of A between 10 and 50. However, each record is likely to be in a different block, because of the ordering of records in the file, leading to many I/O operation. In the worst case, for each record, it needs to traverse the whole tree (cost is h), so the total cost is $n_1 * h$.
- What is the worst case cost of finding records satisfying $10 < A < 50 \wedge 5 < B < 10$ using this index, in terms of the number of records n_2 that satisfy this selection, as well as n_1 and h defined above.
This query can be answered by using an ordered index on the search key (A, B) . For each value of A this is between 10 and 50, the system located records with B value between 5 and 10. However, each record could be likely to be in a different disk block. This amounts to executing the query based on the condition on A , this costs $n_1 * h$. Then these records are checked to see if the condition on B is satisfied. So, the total cost in the worst case is $n_1 * h$.
- Under what conditions on n_1 and n_2 would the index be an efficient way of finding records satisfying $10 < A < 50 \wedge 5 < B < 10$.
 n_1 records satisfy the first condition and n_2 records satisfy the second condition. When both the conditions are queried, n_1 records are output in the first stage. So, in the case where $n_1 = n_2$, no extra records are output in the first stage. Otherwise, the records which

don't satisfy the second condition are also output with an additional cost of h each (worst case).

- 11.23** Suppose you have to create a B⁺-tree index on a large number of names, where the maximum size of a name may be quite large (say 40 characters) and the average name is itself large (say 10 characters). Explain how prefix compression can be used to maximize the average fanout of nonleaf nodes.

Answer: There arises 2 problems in the given scenario. The first problem is names can be of variable length. The second problem is names can be long (maximum is 40 characters), leading to a low fanout and a correspondingly increased tree height. With variable-length search keys, different nodes can have different fanouts even if they are full. The fanout of nodes can be increased by using a technique called prefix compression. With prefix compression, the entire search key value is not stored at internal nodes. Only a prefix of each search key which is sufficient to distinguish between the key values in the subtrees that it separates. The full name can be stored in the leaf nodes, this way we don't lose any information and also maximize the average fanout of internal nodes.

- 11.24** Suppose a relation is stored in a B⁺-tree file organization. Suppose secondary indices stored record identifiers that are pointers to records on disk.

- What would be the effect on the secondary indices if a node split happens in the file organization?
- What would be the cost of updating all affected records in a secondary index?
- How does using the search key of the file organization as a logical record identifier solve this problem?
- What is the extra cost due to the use of such logical record identifiers?

Answer:

- When a leaf page is split in a B⁺-tree file organization, a number of records are moved to a new page. In such cases, all secondary indices that store pointers to the relocated records would have to be updated, even though the values in the records may not have changed.
- Each leaf page may contain a fairly large number of records, and each of them may be in different locations on each secondary index. Thus, a leaf-page split may require tens or even hundreds of I/O operations to update all affected secondary indices, making it a very expensive operation.
- One solution is to store the values of the primary-index search key attributes in secondary indices, in place of pointers to the indexed

records. Relocation of records because of leaf-page splits then does not require any update on any secondary index.

- d. Locating a record using the secondary index now requires two steps: First we use the secondary index to find the primary index search-key values, and then we use the primary index to find the corresponding records. This approach reduces the cost of index update due to file reorganization, although it increases the cost of accessing data using a secondary index.

- 11.25** Show how to compute existence bitmaps from other bitmaps. Make sure that your technique works even in the presence of null values, by using a bitmap for the value *null*.

Answer: The existence bitmap for a relation can be calculated by taking the union (logical-or) of all the bitmaps on that attribute, including the bitmap for value *null*.

- 11.26** How does data encryption affect index schemes? In particular, how might it affect schemes that attempt to store data in sorted order?

Answer: Note that indices must operate on the encrypted data or someone could gain access to the index to interpret the data. Otherwise, the index would have to be restricted so that only certain users could access it. To keep the data in sorted order, the index scheme would have to decrypt the data at each level in a tree. Note that hash systems would not be affected.

- 11.27** Our description of static hashing assumes that a large contiguous stretch of disk blocks can be allocated to a static hash table. Suppose you can allocate only C contiguous blocks. Suggest how to implement the hash table, if it can be much larger than C blocks. Access to a block should still be efficient.

Answer: A separate list/table as shown below can be created.

Starting address of first set of C blocks

C

Starting address of next set of C blocks

$2C$

and so on

Desired block address = Starting address (from the table depending on the block number) + blocksize * (blocknumber % C)

For each set of C blocks, a single entry is added to the table. In this case, locating a block requires 2 steps: First we use the block number to find the actual block address, and then we can access the desired block.

CHAPTER 12



Query Processing

This chapter describes the process by which queries are executed efficiently by a database system. The chapter starts off with measures of cost, then proceeds to algorithms for evaluation of relational algebra operators and expressions. This chapter applies concepts from Chapter 2, 6, 10, and 11.

Query processing algorithms can be covered without tedious and distracting details of size estimation. Although size estimation is covered later, in Chapter 13, the presentation there has been simplified by omitting some details. Instructors can choose to cover query processing but omit query optimization, without loss of continuity with later chapters.

Exercises

- 12.10** Suppose you need to sort a relation of 40 gigabytes, with 4 kilobyte blocks, using a memory size of 40 megabytes. Suppose the cost of a seek is 5 milliseconds, while the disk transfer rate is 40 megabytes per second.
- Find the cost of sorting the relation, in seconds, with $b_b = 1$ and with $b_b = 100$.
 - In each case, how many merge passes are required?
 - Suppose a flash storage device is used instead of a disk, and it has a seek time of 1 microsecond, and a transfer rate of 40 megabytes per second. Recompute the cost of sorting the relation, in seconds, with $b_b = 1$ and with $b_b = 100$, in this setting.

Answer:

- The number of blocks in the main memory buffer available for sorting (M) is $\frac{40 \times 10^6}{4 \times 10^3} = 10^4$. The number of blocks containing records of the given relation (b_r) is $\frac{40 \times 10^9}{4 \times 10^3} = 10^7$. Then the cost of sorting the relation is: $(\text{Number of disk seeks} \times \text{Disk seek cost}) + (\text{Number of block transfers} \times \text{Block transfer time})$. Here Disk seek cost is 5×10^{-3} seconds

and Block transfer time is 10^{-4} seconds ($\frac{4 \times 10^3}{40 \times 10^6}$). The number of block transfers is independent of b_b and is equal to 25×10^6 .

• **Case 1:** $b_b = 1$

Using the equation in Section 12.4, the number of disk seeks is 5002×10^3 . Therefore the cost of sorting the relation is: $(5002 \times 10^3) \times (5 \times 10^{-3}) + (25 \times 10^6) \times (10^{-4}) = 25 \times 10^3 + 2500 = 27500$ seconds.

• **Case 2:** $b_b = 100$

The number of disk seeks is: 52×10^3 . Therefore the cost of sorting the relation is: $(52 \times 10^3) \times (5 \times 10^{-3}) + (25 \times 10^6) \times (10^{-4}) = 260 + 2500 = 2760$ seconds.

b. **Disk storage** The number of merge passes required is given by $\lceil \log_{M-1}(\frac{b_r}{M}) \rceil$. This is independent of b_b . Substituting the values above, we get $\lceil \log_{10^4-1}(\frac{10^7}{10^4}) \rceil$ which evaluates to 1.

c. **Flash storage:**

• **Case 1:** $b_b = 1$

The number of disk seeks is: 5002×10^3 . Therefore the cost of sorting the relation is: $(5002 \times 10^3) \times (1 \times 10^{-6}) + (25 \times 10^6) \times (10^{-4}) = 5.002 + 2500 = 2506$ seconds.

• **Case 2:** $b_b = 100$

The number of disk seeks is: 52×10^3 . Therefore the cost of sorting the relation is: $(52 \times 10^3) \times (1 \times 10^{-6}) + (25 \times 10^6) \times (10^{-4}) = 0.052 + 2500$, which is approx = 2500 seconds.

12.11 Consider the following extended relational-algebra operators. Describe how to implement each operation using sorting, and using hashing.

- Semijoin** (\bowtie_{θ}): $r \bowtie_{\theta} s$ is defined as $\Pi_R(r \bowtie_{\theta} s)$, where R is the set of attributes in the schema of r ; that it selects those tuples r_i in r for which there is a tuple s_j in s such that r_i and s_j satisfy predicate θ .
- Anti-semijoin** ($\bar{\bowtie}_{\theta}$): $r \bar{\bowtie}_{\theta} s$ is defined as $r - \Pi_R(r \bowtie_{\theta} s)$; that it selects those tuples r_i in r for which there is no tuple s_j in s such that r_i and s_j satisfy predicate θ .

Answer: As in the case of join algorithms, semijoin and anti-semijoin can be done efficiently if the join conditions are equijoin conditions. We describe below how to efficiently handle the case of equijoin conditions using sorting and hashing. With arbitrary join conditions, sorting and hashing cannot be used; (block) nested loops join needs to be used instead.

a. **Semijoin:**

- **Semijoin using Sorting:** Sort both r and s on the join attributes in θ . Perform a scan of both r and s similar to the merge al-

algorithm and add tuples of r to the result whenever the join attributes of the current tuples of r and s match.

- **Semijoin using Hashing:** Create a hash index in s on the join attributes in θ . Iterate over r , and for each distinct value of the join attributes, perform a hash lookup in s . If the hash lookup returns a value, add the current tuple of r to the result.

Note that if r and s are large, they can be partitioned on the join attributes first, and the above procedure applied on each partition. If r is small but s is large, a hash index can be built on r , and probed using s ; and if an s tuple matches an r tuple, the r tuple can be output and deleted from the hash index.

b. **Anti-semijoin:**

- **Anti-Semijoin using Sorting:** Sort both r and s on the join attributes in θ . Perform a scan of both r and s similar to the merge algorithm and add tuples of r to the result if no tuple of s satisfies the join predicate for the corresponding tuple of r .
- **Anti-Semijoin using Hashing:** Create a hash index in s on the join attributes in θ . Iterate over r , and for each distinct value of the join attributes, perform a hash lookup in s . If the hash lookup returns a null value, add the current tuple of r to the result.

As for semijoin, partitioning can be used if r and s are large. An index on r can be used instead of an index on s , but then when an s tuple matches an r tuple, the r tuple is deleted from the index. After processing all s tuples, all remaining r tuples in the index are output as the result of the anti-semijoin operation.

- 12.12** Why is it not desirable to force users to make an explicit choice of a query-processing strategy? Are there cases in which it *is* desirable for users to be aware of the costs of competing query-processing strategies? Explain your answer.

Answer: In general it is not desirable to force users to choose a query processing strategy because naive users might select an inefficient strategy. The reason users would make poor choices about processing queries is that they would not know how a relation is stored, nor about its indices. It is unreasonable to force users to be aware of these details since ease of use is a major object of database query languages. If users are aware of the costs of different strategies they could write queries efficiently, thus helping performance. This could happen if experts were using the system.

- 12.13** Design a variant of the hybrid merge-join algorithm for the case where both relations are not physically sorted, but both have a sorted secondary index on the join attributes.

Answer: We merge the leaf entries of the first sorted secondary index with the leaf entries of the second sorted secondary index. The result file contains pairs of addresses, the first address in each pair pointing to a

tuple in the first relation, and the second address pointing to a tuple in the second relation.

This result file is first sorted on the first relation's addresses. The relation is then scanned in physical storage order, and addresses in the result file are replaced by the actual tuple values. Then the result file is sorted on the second relation's addresses, allowing a scan of the second relation in physical storage order to complete the join.

- 12.14** Estimate the number of block transfers and seeks required by your solution to Exercise 12.13 for $r_1 \bowtie r_2$, where r_1 and r_2 are as defined in Practice Exercise 12.3.

Answer: r_1 occupies 800 blocks, and r_2 occupies 1500 blocks. Let there be n pointers per index leaf block (we assume that both the indices have leaf blocks and pointers of equal sizes). Let us assume M pages of memory, $M < 800$. r_1 's index will need $B_1 = \lceil \frac{20000}{n} \rceil$ leaf blocks, and r_2 's index will need $B_2 = \lceil \frac{45000}{n} \rceil$ leaf blocks. Therefore the merge join will need $B_3 = B_1 + B_2$ accesses, without output. The number of output tuples is estimated as $n_o = \lceil \frac{20000 * 45000}{\max(V(C, r_1), V(C, r_2))} \rceil$. Each output tuple will need two pointers, so the number of blocks of join output will be $B_{o1} = \lceil \frac{n_o}{n/2} \rceil$. Hence the join needs $B_j = B_3 + B_{o1}$ disk block accesses.

Now we have to replace the pointers by actual tuples. For the first sorting, $B_{s1} = B_{o1}(2\lceil \log_{M-1}(B_{o1}/M) \rceil + 2)$ disk accesses are needed, including the writing of output to disk. The number of blocks of r_1 which have to be accessed in order to replace the pointers with tuple values is $\min(800, n_o)$. Let n_1 pairs of the form (r_1 tuple, pointer to r_2) fit in one disk block. Therefore the intermediate result after replacing the r_1 pointers will occupy $B_{o2} = \lceil (n_o/n_1) \rceil$ blocks. Hence the first pass of replacing the r_1 -pointers will cost $B_f = B_{s1} + B_{o1} + \min(800, n_o) + B_{o2}$ disk accesses.

The second pass for replacing the r_2 -pointers has a similar analysis. Let n_2 tuples of the final join fit in one block. Then the second pass of replacing the r_2 -pointers will cost $B_s = B_{s2} + B_{o2} + \min(1500, n_o)$ disk accesses, where $B_{s2} = B_{o2}(2\lceil \log_{M-1}(B_{o2}/M) \rceil + 2)$.

Hence the total number of disk accesses for the join is $B_j + B_f + B_s$, and the number of pages of output is $\lceil n_o/n_2 \rceil$.

- 12.15** The hash-join algorithm as described in Section 12.5.5 computes the natural join of two relations. Describe how to extend the hash-join algorithm to compute the natural left outer join, the natural right outer join and the natural full outer join. (Hint: Keep extra information with each tuple in the hash index, to detect whether any tuple in the probe relation matches the tuple in the hash index.) Try out your algorithm on the *takes* and *student* relations.

Answer: For the probe relation tuple t_r under consideration, if no matching tuple is found in the build relation's hash partition, it is padded with nulls and included in the result. This will give us the natural left outer join $t_r \bowtie_{\text{LO}} t_s$. To get the natural right outer join $t_r \bowtie_{\text{RO}} t_s$, we can keep

a boolean flag with each tuple in the current build relation partition s_i residing in memory, and set it whenever any probe relation tuple matches with it. When we are finished with s_i , all the tuples in it which do not have their flag set, are padded with nulls and included in the result. To get the natural full outer join, we do both the above operations together.

To try out our algorithm, we use the sample *student* and *takes* relations of Figures A.9 and A.10. Let us assume that there is enough memory to hold three tuples of the build relation plus a hash index for those three tuples. We use *takes* as the build relation. We use the simple hashing function which returns the *student.ID* mod 10. Taking the partition corresponding to value 7, we get $r_1 = \{("Snow")\}$, and $s_1 = \phi$. The tuple in the probe relation partition will have no matching tuple, so $(("70557", "Snow", "Physics", "0", null))$ is outputted. Proceeding in a similar way, we process all the partitions and complete the join.

12.16 Pipelining is used to avoid writing intermediate results to disk. Suppose you need to sort relation r using sort-merge and merge-join the result with an already sorted relation s .

- Describe how the output of the sort of r can be pipelined to the merge join without being written back to disk.
- The same idea is applicable even if both inputs to the merge join are the outputs of sort-merge operations. However, the available memory has to be shared between the two merge operations (the merge-join algorithm itself needs very little memory). What is the effect of having to share memory on the cost of each sort-merge operation?

Answer:

- Using pipelining, output from the sorting operation on r is written to a buffer B . When B is full, the merge-join processes tuples from B , joining them with tuples from s until B is empty. At this point, the sorting operation is resumed and B is refilled. This process continues until the merge-join is complete.
- If the sort-merge operations are run in parallel and memory is shared equally between the two, each operation will have only $M/2$ frames for its memory buffer. This may increase the number of runs required to merge the data.

12.17 Write pseudocode for an iterator that implements a version of the sort-merge algorithm where the result of the final merge is pipelined to its consumers. Your pseudocode must define the standard iterator functions *open()*, *next()*, and *close()*. Show what state information the iterator must maintain between calls.

Answer: Let M denote the number of blocks in the main memory buffer available for sorting. For simplicity we assume that there are less than M

runs created in the run creation phase. The pseudocode for the iterator functions open, next and close are as follows:

```

SortMergeJoin::open()
begin
    repeat
        read M blocks of the relation;
        sort the in-memory part of the relation;
        write the sorted data to a run file  $R_i$ 
    until the end of the relation
    read one block of each of the  $N$  run files  $R_i$ , into a
        buffer block in memory
     $done_r := false$ ;
end

SortMergeJoin::close()
begin
    clear all the  $N$  runs from main memory and disk;
end

boolean SortMergeJoin::next()
begin
    if the buffer block of any run  $R_i$  is empty and not end-of-file( $R_i$ )
        begin
            read the next block of  $R_i$  (if any) into the buffer block;
        end
    if all buffer blocks are empty
        return false;
    choose the first tuple (in sort order) among the buffer blocks;
    write the tuple to the output buffer;
    delete the tuple from the buffer block and increment its pointer;
    return true;
end

```

- 12.18** Suppose you have to compute ${}_A\mathcal{G}_{sum(C)}(r)$ as well as ${}_{A,B}\mathcal{G}_{sum(C)}(r)$. Describe how to compute these together using a single sorting of r .

Answer: Run the sorting operation on r , grouping by (A, B) , as required for the second result. When evaluating the sum aggregate, keep running totals for both the (A, B) grouping as well as for just the A grouping.

CHAPTER 13



Query Optimization

This chapter describes how queries are optimized. It starts off with statistics used for query optimization, and outlines how to use these statistics to estimate selectivities and query result sizes used for cost estimation. Equivalence rules are covered next, followed by a description of a query optimization algorithm modeled after the classic System R optimization algorithm, and coverage of nested sub-query optimization. The chapter ends with a description of materialized views, their role in optimization and a description of incremental view-maintenance algorithms.

It should be emphasized that the estimates of query sizes and selectivities are approximate, even if the assumptions made, such as uniformity, hold. Further, the cost estimates for various algorithms presented in Chapter 12 assume only a minimal amount of memory, and are thus worst case estimates with respect to buffer space availability. As a result, cost estimates are never very accurate. However, practical experience has shown that such estimates tend to be reasonably accurate, and plans optimal with respect to estimated cost are rarely much worse than a truly optimal plan.

We do *not* expect students to memorize the size-estimates, and we stress only the process of arriving at the estimates, not the exact values. Precision in terms of estimated cost is not a worthwhile goal, so estimates off by a few I/O operations can be considered acceptable.

The theory in this chapter is ideally backed up by lab assignments where students study the query execution plans generated by one or more database systems. Most database products have an “explain plan” feature that lets the user find the evaluation plan used on a query. It is worthwhile asking students to explore the plans generated for different queries, with and without indices. Assignment may be designed in which students measure the performance speedup provided by indices. A more challenging assignment is to design tests to see how clever the query optimizer is, and to guess from these experiments which of the optimization techniques covered in the chapter are used in the system.

Exercises

- 13.15 Suppose that a B⁺-tree index on (*dept_name*, *building*) is available on relation *department*. What would be the best way to handle the following selection?

$$\sigma_{(building < \text{"Watson"}) \wedge (budget < 55000) \wedge (dept_name = \text{"Music"})}(department)$$

Answer: Using the index on (*dept_name*, *building*), we locate the first tuple having (*building* “Watson” and *dept_name* “Music”). We then follow the pointers retrieving successive tuples as long as *building* is less than “Watson”. From the tuples retrieved, the ones not satisfying the condition (*budget* < 55000) are rejected.

- 13.16 Show how to derive the following equivalences by a sequence of transformations using the equivalence rules in Section 13.2.1.
- $\sigma_{\theta_1 \wedge \theta_2 \wedge \theta_3}(E) = \sigma_{\theta_1}(\sigma_{\theta_2}(\sigma_{\theta_3}(E)))$
 - $\sigma_{\theta_1 \wedge \theta_2}(E_1 \bowtie_{\theta_3} E_2) = \sigma_{\theta_1}(E_1 \bowtie_{\theta_3} (\sigma_{\theta_2}(E_2)))$, where θ_2 involves only attributes from E_2

Answer:

- Using rule 1, $\sigma_{\theta_1 \wedge \theta_2 \wedge \theta_3}(E)$ becomes $\sigma_{\theta_1}(\sigma_{\theta_2 \wedge \theta_3}(E))$. On applying rule 1 again, we get $\sigma_{\theta_1}(\sigma_{\theta_2}(\sigma_{\theta_3}(E)))$.
 - $\sigma_{\theta_1 \wedge \theta_2}(E_1 \bowtie_{\theta_3} E_2)$ on applying rule 1 becomes $\sigma_{\theta_1}(\sigma_{\theta_2}(E_1 \bowtie_{\theta_3} E_2))$. This on applying rule 7.a becomes $\sigma_{\theta_1}(E_1 \bowtie_{\theta_3} (\sigma_{\theta_2}(E_2)))$.
- 13.17 Consider the two expressions $\sigma_{\theta}(E_1 \bowtie E_2)$ and $\sigma_{\theta}(E_1 \bowtie E_2)$.
- Show using an example that the two expressions are not equivalent in general.
 - Give a simple condition on the predicate θ , which if satisfied will ensure that the two expressions are equivalent.

Answer:

- Consider relations *dept*(*id*, *deptname*) and *emp*(*id*, *name*, *dept_id*) with sample data as follows:
Sample data for *dept*:

| | |
|-----|----------------|
| 501 | Finance |
| 502 | Administration |
| 503 | Marketing |
| 504 | Sales |

Sample data for *emp*:

| | |
|---|------------|
| 1 | John 501 |
| 2 | Martin 503 |
| 3 | Sarah 504 |

Now consider the expressions

$$\sigma_{deptname < 'Z'}(dept \bowtie emp)$$

and $\sigma_{deptname < 'Z'}(dept \Join emp)$.

The result of the first expression is:

| | | | |
|-----|-----------|---|--------|
| 501 | Finance | 1 | John |
| 503 | Marketing | 2 | Martin |
| 504 | Sales | 3 | Sarah |

whereas the result of the second expression is:

| | | | |
|-----|----------------|------|--------|
| 501 | Finance | 1 | John |
| 502 | Administration | null | null |
| 503 | Marketing | 2 | Martin |
| 504 | Sales | 3 | Sarah |

- b. Considering the same example, if θ included the condition “*name* = ‘Einstein’”, the two expressions would be equivalent, that is they would always have the same result, since any tuple that is in in $dept \Join emp$ but not in $dept \bowtie emp$ would not satisfy the condition since its *name* attribute would be null.

- 13.18** A set of equivalence rules is said to be *complete* if, whenever two expressions are equivalent, one can be derived from the other by a sequence of uses of the equivalence rules. Is the set of equivalence rules that we considered in Section 13.2.1 complete? Hint: Consider the equivalence $\sigma_{3=5}(r) = \{ \}$.

Answer: Two relational expressions are defined to be *equivalent* when on all input relations, they give the same output. The set of equivalence rules considered in Section 13.2.1 is not complete. The expressions $\sigma_{3=5}(r)$ and $\{ \}$ are equivalent, but this cannot be shown by using just these rules.

- 13.19** Explain how to use a histogram to estimate the size of a selection of the form $\sigma_{A \leq v}(r)$.

Answer: Suppose the histogram H storing the distribution of values in r is divided into ranges r_1, \dots, r_n . For each range r_i with low value $r_{i:low}$ and high value $r_{i:high}$, if $r_{i:high}$ is less than v , we add the number of tuples given by

$$H(r_i)$$

to the estimated total. If $v < r_{i:high}$ and $v \geq r_{i:low}$, we assume that values within r_i are uniformly distributed and we add

$$H(r_i) * \frac{v - r_{i:low}}{r_{i:high} - r_{i:low}}$$

to the estimated total.

- 13.20** Suppose two relations r and s have histograms on attributes $r.A$ and $s.A$, respectively, but with different ranges. Suggest how to use the histograms to estimate the size of $r \bowtie s$. Hint: Split the ranges of each histogram further.

Answer: Find the largest unit u that evenly divides the range size of both histograms. Divide each histogram into ranges of size u , assuming that values within a range are uniformly distributed. Then compute the estimated join size using the technique for histograms with equal range sizes.

- 13.21** Consider the query

```
select A, B
from r
where r.B < some (select B
                  from s
                  where s.A = r.A)
```

Show how to decorrelate the above query using the multiset version of the semijoin operation, defined in Exercise 13.14.

Answer: The solution can be written in relational algebra as follows: $\Pi_{A,B}(r \bowtie \theta s)$ where $\theta = (r.B < s.B \wedge s.A = r.A)$.

The SQL query corresponding to this can be written if the database provides a semi join operator, and this varies across implementations.

- 13.22** Describe how to incrementally maintain the results of the following operations, on both insertions and deletions:

- Union and set difference.
- Left outer join.

Answer:

- Given materialized view $v = r \cup s$, when a tuple is inserted in r , we check if it is present in v , and if not we add it to v . When a tuple is deleted from r , we check if it is there in s , and if not, we delete it from v . Inserts and deletes in s are handled in symmetric fashion. For set difference, given view $v = r - s$, when a tuple is inserted in r , we check if it is present in s , and if not we add it to v . When a tuple is deleted from r , we delete it from v if present. When a tuple is inserted in s , we delete it from v if present. When a tuple is deleted from s , we check if it is present in r , and if so we add it to v .
- Given materialized view $v = r \bowtie s$, when a set of tuples i_r is inserted in r , we add the tuples $i_r \bowtie s$ to the view. When i_r is deleted from r , we delete $i_r \bowtie s$ from the view. When a set of tuples i_s is inserted in s , we compute $r \bowtie i_s$. We find all the tuples of r which previously did not match any tuple from s (i.e. those padded with *null* in $r \bowtie s$) but which match i_s . We remove all those *null* padded entries from

the view and add the tuples $r \bowtie s$ to the view. When i_s is deleted from s , we delete the tuples $r \bowtie i_s$ from the view. Also we find all the tuples in r which match i_s but which do not match any other tuples in s . We add all those to the view, after padding them with *null* values.

- 13.23** Give an example of an expression defining a materialized view and two situations (sets of statistics for the input relations and the differentials) such that incremental view maintenance is better than recomputation in one situation, and recomputation is better in the other situation.

Answer: Let r and s be two relations. Consider a materialized view on these defined by $(r \bowtie s)$. Suppose 70% of the tuples in r are deleted. Then recomputation is better than incremental view maintenance. This is because in incremental view maintenance, the 70% of the deleted tuples have to be joined with s while in recomputation, just the remaining 30% of the tuples in r have to be joined with s .

However, if the number of tuples in r is increased by a small percentage, for example 2%, then incremental view maintenance is likely to be better than recomputation.

- 13.24** Suppose you want to get answers to $r \bowtie s$ sorted on an attribute of r , and want only the top K answers for some relatively small K . Give a good way of evaluating the query:

- When the join is on a foreign key of r referencing s , where the foreign key attribute is declared to be not null.
- When the join is not on a foreign key.

Answer:

- Sort r and collect the top K tuples. These tuples are guaranteed to be contained in $r \bowtie s$ since the join is on a foreign key of r referencing s .
- Execute $r \bowtie s$ using a standard join algorithm until the first K results have been found. After K tuples have been computed in the result set, continue executing the join but immediately discard any tuples from r that have attribute values less than all of the tuples in the result set. If a newly joined tuple t has an attribute value greater than at least one of the tuples in the result set, replace the lowest-valued tuple in the result set with t .

- 13.25** Consider a relation $r(A, B, C)$, with an index on attribute A . Give an example of a query that can be answered by using the index only, without looking at the tuples in the relation. (Query plans that use only the index, without accessing the actual relation, are called *index-only* plans.)

Answer: Any query that only involves the attribute A of r can be executed by only using the index. For example, the query


```
select sum(A)
from r
```

only needs to use the values of A , and thus does not need to look at r .

- 13.26** Suppose you have an update query U . Give a simple sufficient condition on U that will ensure that the Halloween problem cannot occur, regardless of the execution plan chosen, or the indices that exist.

Answer: The attributes referred in the WHERE clause of the update query U should not be a part of the SET clauses of U . This will ensure that the Halloween problem cannot occur.

CHAPTER 14



Transactions

This chapter provides an overview of transaction processing. It first motivates the problems of atomicity, consistency, isolation and durability, and introduces the notion of ACID transactions. It then presents some naive schemes, and their drawbacks, thereby motivating the techniques described in Chapters 15 and 16. The rest of the chapter describes the notion of schedules and the concept of serializability.

We strongly recommend covering this chapter in a first course on databases, since it introduces concepts that every database student should be aware of. Details on how to implement the transaction properties are covered in Chapters 15 and 16.

In the initial presentation to the ACID requirements, the isolation requirement on concurrent transactions does not insist on serializability. Following Haerder and Reuter [1983], isolation just requires that the events within a transaction must be hidden from other transactions running concurrently, in order to allow rollback. However, later in the chapter, and in most of the book (except in Chapter 26), we use the stronger condition of serializability as a requirement on concurrent transactions.

Exercises

14.12 List the ACID properties. Explain the usefulness of each.

Answer: The ACID properties, and the need for each of them are:

- **Consistency:** Execution of a transaction in isolation (that is, with no other transaction executing concurrently) preserves the consistency of the database. This is typically the responsibility of the application programmer who codes the transactions.
- **Atomicity:** Either all operations of the transaction are reflected properly in the database, or none are. Clearly lack of atomicity will lead to inconsistency in the database.

- **Isolation:** When multiple transactions execute concurrently, it should be the case that, for every pair of transactions T_i and T_j , it appears to T_i that either T_j finished execution before T_i started, or T_j started execution after T_i finished. Thus, each transaction is unaware of other transactions executing concurrently with it. The user view of a transaction system requires the isolation property, and the property that concurrent schedules take the system from one consistent state to another. These requirements are satisfied by ensuring that only serializable schedules of individually consistency preserving transactions are allowed.
- **Durability:** After a transaction completes successfully, the changes it has made to the database persist, even if there are system failures.

14.13 During its execution, a transaction passes through several states, until it finally commits or aborts. List all possible sequences of states through which a transaction may pass. Explain why each state transition may occur.

Answer: The possible sequences of states are:-

- active* \rightarrow *partially committed* \rightarrow *committed*. This is the normal sequence a successful transaction will follow. After executing all its statements it enters the *partially committed* state. After enough recovery information has been written to disk, the transaction finally enters the *committed* state.
- active* \rightarrow *partially committed* \rightarrow *aborted*. After executing the last statement of the transaction, it enters the *partially committed* state. But before enough recovery information is written to disk, a hardware failure may occur destroying the memory contents. In this case the changes which it made to the database are undone, and the transaction enters the *aborted* state.
- active* \rightarrow *failed* \rightarrow *aborted*. After the transaction starts, if it is discovered at some point that normal execution cannot continue (either due to internal program errors or external errors), it enters the *failed* state. It is then rolled back, after which it enters the *aborted* state.

14.14 Explain the distinction between the terms *serial schedule* and *serializable schedule*.

Answer: A schedule in which all the instructions belonging to one single transaction appear together is called a *serial schedule*. A *serializable schedule* has a weaker restriction that it should be *equivalent* to some serial schedule. There are two definitions of schedule equivalence – conflict equivalence and view equivalence. Both of these are described in the chapter.

14.15 Consider the following two transactions:

```

T13: read(A);
      read(B);
      if A = 0 then B := B + 1;
      write(B).
T14: read(B);
      read(A);
      if B = 0 then A := A + 1;
      write(A).

```

Let the consistency requirement be $A = 0 \vee B = 0$, with $A = B = 0$ the initial values.

- Show that every serial execution involving these two transactions preserves the consistency of the database.
- Show a concurrent execution of T_{13} and T_{14} that produces a nonserializable schedule.
- Is there a concurrent execution of T_{13} and T_{14} that produces a serializable schedule?

Answer:

- There are two possible executions: $T_{13} T_{14}$ and $T_{14} T_{13}$.

Case 1:

| | A | B |
|----------------|---|---|
| initially | 0 | 0 |
| after T_{13} | 0 | 1 |
| after T_{14} | 0 | 1 |

Consistency met: $A = 0 \vee B = 0 \equiv T \vee F = T$

Case 2:

| | A | B |
|----------------|---|---|
| initially | 0 | 0 |
| after T_{14} | 1 | 0 |
| after T_{13} | 1 | 0 |

Consistency met: $A = 0 \vee B = 0 \equiv F \vee T = T$

- Any interleaving of T_{13} and T_{14} results in a non-serializable schedule.

| T_1 | T_2 |
|---------------------------------------|---------------------------------------|
| read (A) | read (B) |
| read (B) | read (A) |
| if A = 0 then B = B + 1 | if B = 0 then A = A + 1 |
| write (B) | write (A) |

| T_{13} | T_{14} |
|---------------------------------------|---------------------------------------|
| read (A) | read (B) |
| read (B) | read (A) |
| if A = 0 then B = B + 1 | if B = 0 then A = A + 1 |
| write (B) | write (A) |

- c. There is no parallel execution resulting in a serializable schedule. From part a. we know that a serializable schedule results in $A = 0 \vee B = 0$. Suppose we start with T_{13} **read**(A). Then when the schedule ends, no matter when we run the steps of T_2 , $B = 1$. Now suppose we start executing T_{14} prior to completion of T_{13} . Then T_2 **read**(B) will give B a value of 0. So when T_2 completes, $A = 1$. Thus $B = 1 \wedge A = 1 \rightarrow \neg (A = 0 \vee B = 0)$. Similarly for starting with T_{14} **read**(B).

14.16 Give an example of a serializable schedule with two transactions such that the order in which the transactions commit is different from the serialization order.

Answer:

| T_1 | T_2 |
|-------------------|------------------|
| read (A) | read (B) |
| unlock (A) | write (B) |
| | read (A) |
| | write (A) |
| | commit |
| commit | |

As we can see, the above schedule is serializable with an equivalent serial schedule T_1, T_2 . In the above schedule T_2 commits before T_1 . Note that the unlock instruction is added to show how this schedule can occur even

with strict two-phase locking, where exclusive locks are held to commit, but shared locks can be released early in two-phase manner.

- 14.17** What is a recoverable schedule? Why is recoverability of schedules desirable? Are there any circumstances under which it would be desirable to allow nonrecoverable schedules? Explain your answer.

Answer: A recoverable schedule is one where, for each pair of transactions T_i and T_j such that T_j reads data items previously written by T_i , the commit operation of T_i appears before the commit operation of T_j . Recoverable schedules are desirable because failure of a transaction might otherwise bring the system into an irreversibly inconsistent state. Non-recoverable schedules may sometimes be needed when updates must be made visible early due to time constraints, even if they have not yet been committed, which may be required for very long duration transactions.

- 14.18** Why do database systems support concurrent execution of transactions, in spite of the extra programming effort needed to ensure that concurrent execution does not cause any problems?

Answer: Transaction-processing systems usually allow multiple transactions to run concurrently. It is far easier to insist that transactions run serially. However there are two good reasons for allowing concurrency:

- Improved throughput and resource utilization. A transaction may involve I/O activity, CPU activity. The CPU and the disk in a computer system can operate in parallel. This can be exploited to run multiple transactions in parallel. For example, while a read or write on behalf of one transaction is in progress on one disk, another transaction can be running in the CPU. This increases the throughput of the system.
- Reduced waiting time. If transactions run serially, a short transaction may have to wait for a preceding long transaction to complete. If the transactions are operating on different parts of the database, it is better to let them run concurrently, sharing the CPU cycles and disk accesses among them. It reduces the unpredictable delays and the average response time.

- 14.19** Explain why the read-committed isolation level ensures that schedules are cascade-free.

Answer:

The read-committed isolation level ensures that a transaction reads only the committed data. A transaction T_i can not read a data item X which has been modified by a yet uncommitted concurrent transaction T_j . This makes T_i independent of the success or failure of T_j . Hence, the schedules which follow read committed isolation level become cascade-free.

- 14.20** For each of the following isolation levels, give an example of a schedule that respects the specified level of isolation, but is not serializable:

- a. Read uncommitted

- b. Read committed
- c. Repeatable read

Answer:

- a. Read Uncommitted:

| T_1 | T_2 |
|-----------------------------------|-----------------------------------|
| read(A) write(A) | |
| read(A) | read(A) write(A) |

In the above schedule, T_2 reads the value of A written by T_1 even before T_1 commits. This schedule is not serializable since T_1 also reads a value written by T_2 , resulting in a cycle in the precedence graph.

- b. Read Committed:

| T_1 | T_2 |
|---|--|
| lock-S(A) read(A) unlock(A) | |
| lock-S(A) read(A) unlock-S(A) commit | lock-X(A) write(A) unlock(A) commit |

In the above schedule, the first time T_1 reads A , it sees a value of A before it was written by T_2 , while the second **read(A)** by T_1 sees the value written by T_2 (which has already committed). The first read results in T_1 preceding T_2 , while the second read results in T_2 preceding T_1 , and thus the schedule is not serializable.

- c. Repeatable Read :

Consider the following schedule, where T_1 reads all tuples in r satisfying predicate P ; to satisfy repeatable read, it must also share-lock these tuples in a two-phase manner.

| T_1 | T_2 |
|-------------------------------------|---|
| pred_read(r, P) | |
| read(A) commit | insert(t) write(A) commit |

Suppose that the tuple t inserted by T_2 satisfies P ; then the insert by T_2 causes T_2 to be serialized after T_1 , since T_1 does not see t . However, the final **read**(A) operation of T_1 forces T_2 to precede T_1 , causing a cycle in the precedence graph.

14.21 Suppose that in addition to the operations **read** and **write**, we allow an operation **pred_read**(r, P), which reads all tuples in relation r that satisfy predicate P .

- Give an example of a schedule using the **pred_read** operation that exhibits the phantom phenomenon, and is non-serializable as a result.
- Give an example of a schedule where one transaction uses the **pred_read** operation on relation r and another concurrent transaction deletes a tuple from r , but the schedule does not exhibit a phantom conflict. (To do so, you have to give the schema of relation r , and show the attribute values of the deleted tuple.)

Answer:

- The repeatable read schedule in the preceding question is an example of a schedule exhibiting the phantom phenomenon and is non-serializable.

- Consider the schedule

| T_1 | T_2 |
|---------------------------------|-----------------------|
| pred_read ($r, r.A=5$) | delete (t) |
| read (B) | write (B) |
| commit | commit |

Suppose that tuple t deleted by T_2 is from relation r , but does not satisfy predicate P , for example because its A value is 3. Then, there is no phantom conflict between T_1 and T_2 , and T_2 can be serialized before T_1 .

CHAPTER 15



Concurrency Control

This chapter describes how to control concurrent execution in a database, in order to ensure the isolation properties of transactions. A variety of protocols are described for this purpose. If time is short, some of the protocols may be omitted. We recommend covering, at the least, two-phase locking (Sections 15.1.1), through 15.1.3, deadlock detection and recovery (Section 15.2, omitting Section 15.2.1), and the phantom phenomenon (Section 15.8.3). The most widely used techniques would thereby be covered.

It is worthwhile pointing out how the graph-based locking protocols generalize simple protocols, such as ordered acquisition of locks, which students may have studied in an operating system course. Although the timestamp protocols by themselves are not widely used, multiversion two-phase locking (Section 15.6.2) is of increasing importance since it allows long read-only transactions to run concurrently with updates.

The phantom phenomenon is often misunderstood by students as showing that two-phase locking is incorrect. It is worth stressing that transactions that scan a relation must read some data to find out what tuples are in the relation; as long as this data is itself locked in a two-phase manner, the phantom phenomenon will not arise.

Exercises

15.20 What benefit does strict two-phase locking provide? What disadvantages result?

Answer: Because it produces only cascadeless schedules, recovery is very easy. But the set of schedules obtainable is a subset of those obtainable from plain two phase locking, thus concurrency is reduced.

15.21 Most implementations of database systems use strict two-phase locking. Suggest three reasons for the popularity of this protocol.

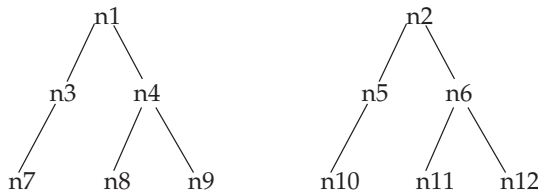
Answer: It is relatively simple to implement, imposes low rollback overhead because of cascadeless schedules, and usually allows an acceptable level of concurrency.

15.22 Consider a variant of the tree protocol called the *forest* protocol. The database is organized as a forest of rooted trees. Each transaction T_i must follow the following rules:

- The first lock in each tree may be on any data item.
- The second, and all subsequent, locks in a tree may be requested only if the parent of the requested node is currently locked.
- Data items may be unlocked at any time.
- A data item may not be relocked by T_i after it has been unlocked by T_i .

Show that the forest protocol does *not* ensure serializability.

Answer: Take a system with 2 trees:



We have 2 transactions, T_1 and T_2 . Consider the following legal schedule:

| T_1 | T_2 |
|---|---|
| lock(n1) lock(n3) write(n3) unlock(n3) lock(n5) read(n5) unlock(n5) unlock(n1) | lock(n2) lock(n5) write(n5) unlock(n5) lock(n3) read(n3) unlock(n3) unlock(n2) |

This schedule is not serializable.

- 15.23 Under what conditions is it less expensive to avoid deadlock than to allow deadlocks to occur and then to detect them?

Answer: Deadlock avoidance is preferable if the consequences of abort are serious (as in interactive transactions), and if there is high contention and a resulting high probability of deadlock.

- 15.24 If deadlock is avoided by deadlock-avoidance schemes, is starvation still possible? Explain your answer.

Answer: A transaction may become the victim of deadlock-prevention rollback arbitrarily many times, thus creating a potential starvation situation.

- 15.25 In multiple-granularity locking, what is the difference between implicit and explicit locking?

Answer: When a transaction *explicitly* locks a node in shared or exclusive mode, it *implicitly* locks all the descendents of that node in the same mode. The transaction need not explicitly lock the descendent nodes. There is no difference in the functionalities of these locks, the only difference is in the way they are acquired, and their presence tested.

- 15.26 Although SIX mode is useful in multiple-granularity locking, an exclusive and intention-shared (XIS) mode is of no use. Why is it useless?

Answer: An exclusive lock is incompatible with any other lock kind. Once a node is locked in exclusive mode, none of the descendents can be simultaneously accessed by any other transaction in any mode. Therefore an exclusive and intend-shared declaration has no meaning.

- 15.27 The multiple-granularity protocol rules specify that a transaction T_i can lock a node Q in S or IS mode only if T_i currently has the parent of Q locked in either IX or IS mode. Given that SIX and S locks are stronger than IX or IS locks, why does the protocol not allow locking a node in S or IS mode if the parent is locked in either SIX or S mode?

Answer:

If T_i has locked the parent node P in S or SIX mode then it means it has implicit S locks on all the descendent nodes of the parent node including Q . So, there is no need for locking Q in S or IS mode and the protocol does not allow doing that.

- 15.28 When a transaction is rolled back under timestamp ordering, it is assigned a new timestamp. Why can it not simply keep its old timestamp?

Answer: A transaction is rolled back because a newer transaction has read or written the data which it was supposed to write. If the rolled back transaction is re-introduced with the same timestamp, the same reason for rollback is still valid, and the transaction will have to be rolled back again. This will continue indefinitely.

15.29 Show that there are schedules that are possible under the two-phase locking protocol, but are not possible under the timestamp protocol, and vice versa.

Answer: A schedule which is allowed in the two-phase locking protocol but not in the timestamp protocol is:

| step | T_0 | T_1 | Precedence marks |
|------|------------------|------------------|-----------------------|
| 1 | lock-S(A) | | |
| 2 | read(A) | | |
| 3 | | lock-X(B) | |
| 4 | | write(B) | |
| 5 | | unlock(B) | |
| 6 | lock-S(B) | | |
| 7 | read(B) | | $T_1 \rightarrow T_0$ |
| 8 | unlock(A) | | |
| 9 | unlock(B) | | |

This schedule is not allowed in the timestamp protocol because at step 7, the W-timestamp of B is 1.

A schedule which is allowed in the timestamp protocol but not in the two-phase locking protocol is:

| step | T_0 | T_1 | T_2 |
|------|-----------------|-----------------|-----------------|
| 1 | write(A) | | |
| 2 | | write(A) | |
| 3 | | | write(A) |
| 4 | write(B) | | |
| 5 | | write(B) | |

This schedule cannot have lock instructions added to make it legal under two-phase locking protocol because T_1 must unlock (A) between steps 2 and 3, and must lock (B) between steps 4 and 5.

15.30 Under a modified version of the timestamp protocol, we require that a commit bit be tested to see whether a read request must wait. Explain how the commit bit can prevent cascading abort. Why is this test not necessary for write requests?

Answer: Using the commit bit, a read request is made to wait if the transaction which wrote the data item has not yet committed. Therefore, if the writing transaction fails before commit, we can abort that transaction alone. The waiting read will then access the earlier version in case of a

multiversion system, or the restored value of the data item after abort in case of a single-version system. For writes, this commit bit checking is unnecessary. That is because either the write is a “blind” write and thus independent of the old value of the data item or there was a prior read, in which case the test was already applied.

- 15.31** As discussed in Exercise 15.19, snapshot isolation can be implemented using a form of timestamp validation. However, unlike the multiversion timestamp-ordering scheme, which guarantees serializability, snapshot isolation does not guarantee serializability. Explain what is the key difference between the protocols that results in this difference.

Answer:

The timestamp validation step for the snapshot isolation level checks for the presence of common written data items between the transactions. However, write skew can occur, where a transaction T_1 updates an item A whose old version is read by T_2 , while T_2 updates an item B whose old version is read by T_1 , resulting in a non-serializable execution. There is no validation of reads against writes in the snapshot isolation protocol.

The multiversion timestamp-ordering protocol on the other hand avoids the write skew problem by rolling back a transaction that writes a data item which has been already read by a transaction with a higher timestamp.

- 15.32** Outline the key similarities and differences between the timestamp based implementation of the first-committer-wins version of snapshot isolation, described in Exercise 15.19, and the optimistic-concurrency-control-without-read-validation scheme, described in Section 15.9.3.

Answer: Both the schemes do not ensure serializability. The version number check in the optimistic-concurrency-control-without-read-validation implements the first committer-wins rule used in the snapshot isolation. Unlike the snapshot isolation, the reads performed by a transaction in optimistic-concurrency-control-without-read-validation may not correspond to the snapshot of the database. Different reads by the same transaction may return data values corresponding to different snapshots of the database.

- 15.33** Explain the phantom phenomenon. Why may this phenomenon lead to an incorrect concurrent execution despite the use of the two-phase locking protocol?

Answer: The phantom phenomenon arises when, due to an insertion or deletion, two transactions logically conflict despite not locking any data items in common. The insertion case is described in the book. Deletion can also lead to this phenomenon. Suppose T_i deletes a tuple from a relation while T_j scans the relation. If T_i deletes the tuple and then T_j reads the relation, T_i should be serialized before T_j . Yet there is no tuple that both T_i and T_j conflict on.

An interpretation of 2PL as just locking the accessed tuples in a relation is incorrect. There is also an index or a relation data that has information

about the tuples in the relation. This information is read by any transaction that scans the relation, and modified by transactions that update, or insert into, or delete from the relation. Hence locking must also be performed on the index or relation data, and this will avoid the phantom phenomenon.

- 15.34** Explain the reason for the use of degree-two consistency. What disadvantages does this approach have?

Answer: The degree-two consistency avoids cascading aborts and offers increased concurrency but the disadvantage is that it does not guarantee serializability and the programmer needs to ensure it.

- 15.35** Give example schedules to show that with key-value locking, if any of lookup, insert, or delete do not lock the next-key value, the phantom phenomenon could go undetected.

Answer: In the next-key locking technique, every index lookup or insert or delete must not only the keys found within the range (or the single key, in case of a point lookup) but also the next-key value- that is, the key value greater than the last key value that was within the range. Thus, if a transaction attempts to insert a value that was within the range of the index lookup of another transaction, the two transactions would conflict on the key value next to the inserted key value. The next-key value should be locked to ensure that conflicts with subsequent range lookups of other queries are detected, thereby detecting phantom phenomenon.

- 15.36** Many transactions update a common item (e.g., the cash balance at a branch), and private items (e.g., individual account balances). Explain how you can increase concurrency (and throughput) by ordering the operations of the transaction.

Answer: The private items can be updated by the individual transactions independently. They can acquire the exclusive locks for the private items (as no other transaction needs it) and update the data items. But the exclusive lock for the common item is shared among all the transactions. The common item should be locked before the transaction decides to update it. And when it holds the lock for the common item, all other transactions should wait till its released. But in order that the common item is updated correctly, the transaction should follow a certain pattern. A transaction can update its private item as and when it requires, but before updating the private item again, the common item should be updated. So, essentially the private and the common items should be accessed alternately, otherwise the private item's update will not be reflected in the common item.

- a. No possibility of deadlock and no starvation. The lock for the common item should be granted based on the time of requests.
- b. The schedule is serializable.

- 15.37** Consider the following locking protocol: All items are numbered, and once an item is unlocked, only higher-numbered items may be locked.

Locks may be released at any time. Only X-locks are used. Show by an example that this protocol does not guarantee serializability.

Answer: We have 2 transactions, T_1 and T_2 . Consider the following legal schedule:

| T_1 | T_2 |
|------------------|------------------|
| lock(A) | |
| write(A) | |
| unlock(A) | |
| | lock(A) |
| | read(A) |
| | lock(B) |
| | write(B) |
| | unlock(B) |
| lock(B) | |
| read(B) | |
| unlock(B) | |

Explanation: In the given example schedule, let's assume A is a higher numbered item than B.

- T_i executes write(A) before T_j executes read(A). So, there's an edge $T_i \rightarrow T_j$.
- T_j executes write(B) before T_i executes read(B). So, there's an edge $T_j \rightarrow T_i$.

There's a cycle in the graph which means the given schedule is not conflict serializable.

CHAPTER 16



Recovery System

This chapter covers failure models and a variety of failure recovery techniques. Recovery in a real-life database systems supporting concurrent transactions is rather complicated. To help the student understand concepts better, the chapter presents recovery models in increasing degree of complexity.

The coverage of recovery in the 6th edition has changed significantly from the previous editions. In particular, the number of different algorithms presented has been reduced from earlier, avoiding the confusion between algorithms that perform undo first, then redo, versus algorithms that perform redo first, then undo; we have standardized on the latter in this edition.

Coverage of recovery algorithms now follows the sequence below:

- In Section 16.3 we describe a number of key concepts in recovery, before presenting the basic recovery algorithm in Section 16.4; the algorithm in this section supports concurrency with strict two-phase locking.
- Sections 16.5 and 16.6 discuss how to extend the recovery algorithm to deal with issues in buffer management and failure of non-volatile storage.
- Section 16.7 discusses how to handle early lock release by using logical undo operations; the algorithm described here was referred to as the “Advanced Recovery Algorithm” in the 5th edition. The coverage here has been extended, with more examples.
- Section 16.8 presents the ARIES recovery algorithm, variants of which are widely used in practice. ARIES includes some support for early lock release, along with several optimizations that speed up recovery greatly.

Exercises

- 16.12** Explain the difference between the three storage types—volatile, non-volatile, and stable—in terms of I/O cost.

Answer: Volatile storage is storage which fails when there is a power failure. Cache, main memory, and registers are examples of volatile storage.

Non-volatile storage is storage which retains its content despite power failures. An example is magnetic disk. Stable storage is storage which theoretically survives any kind of failure (short of a complete disaster!). This type of storage can only be approximated by replicating data.

In terms of I/O cost, volatile memory is the fastest and non-volatile storage is typically several times slower. Stable storage is slower than non-volatile storage because of the cost of data replication.

16.13 Stable storage cannot be implemented.

- a. Explain why it cannot be.
- b. Explain how database systems deal with this problem.

Answer:

- a. Stable storage cannot really be implemented because all storage devices are made of hardware, and all hardware is vulnerable to mechanical or electronic device failures.
- b. Database systems approximate stable storage by writing data to multiple storage devices simultaneously. Even if one of the devices crashes, the data will still be available on a different device. Thus data loss becomes extremely unlikely.

16.14 Explain how the database may become inconsistent if some log records pertaining to a block are not output to stable storage before the block is output to disk.

Answer: Consider a banking scheme and a transaction which transfers \$50 from account A to account B . The transaction has the following steps:

- a. **read**(A, a_1)
- b. $a_1 := a_1 - 50$
- c. **write**(A, a_1)
- d. **read**(B, b_1)
- e. $b_1 := b_1 + 50$
- f. **write**(B, b_1)

Suppose the system crashes after the transaction commits, but before its log records are flushed to stable storage. Further assume that at the time of the crash the update of A in the third step alone had actually been propagated to disk whereas the buffer page containing B was not yet written to disk. When the system comes up it is in an inconsistent state, but recovery is not possible because there are no log records corresponding to this transaction in stable storage.

16.15 Outline the drawbacks of the no-steal and force buffer management policies.

Answer: Drawback of the no-steal policy: The no-steal policy does not work with transactions that perform a large number of updates, since the buffer may get filled with updated pages that cannot be evicted to disk, and the transaction cannot then proceed.

Drawback of the force policy: The force policy might slow down the commit of transactions as it forces all modified blocks to be flushed to disk before commit. Also, the number of output operations is more in the case of force policy. This is because frequently updated blocks are output multiple times, once for each transaction. Further, the policy results in more seeks, since the blocks written out are not likely to be consecutive on disk.

- 16.16** Physiological redo logging can reduce logging overheads significantly, especially with a slotted page record organization. Explain why.

Answer: If a slotted page record organization is used, the deletion of a record from a page may result in many other records in the page being shifted. With physical redo logging, all bytes of the page affected by the shifting of records must be logged. On the other hand, if physiological redo logging is used, only the deletion operation is logged. This results in a much smaller log record. Redo of the deletion operation would delete the record and shift other records as required.

- 16.17** Explain why logical undo logging is used widely, whereas logical redo logging (other than physiological redo logging) is rarely used.

Answer: The class of operations which release locks early are called logical operations. Once such lower level locks are released, such operations cannot be undone by using the old values of updated data items. Instead, they must be undone by executing a compensating operation called a logical undo operation. In order to allow logical undo of operations, special log records are necessary to store the necessary logical undo information. Thus logical undo logging is used widely.

Redo operations are performed exclusively using physical log records. This is because the state of the database after a system failure may reflect some updates of an operation and not of other operations, depending on what buffer blocks had been written to disk before the failure. The database state on disk might not be in an *operation consistent* state, i.e., it might have partial effects of operations. Logical undo or redo operations cannot be performed on an inconsistent data structure.

- 16.18** Consider the log in Figure 16.5. Suppose there is a crash just before the $\langle T_0 \text{ abort} \rangle$ log record is written out. Explain what would happen during recovery.

Answer: Recovery would happen as follows:

Redo phase:

- a. Undo-List = T_0, T_1

- b. Start from the checkpoint entry and perform the redo operation.
- c. $C = 600$
- d. T_1 is removed from the Undo-list as there is a commit record.
- e. T_2 is added to the Undo list on encountering the $\langle T_2 \text{ start} \rangle$ record.
- f. $A = 400$
- g. $B = 2000$

Undo phase:

- a. Undo-List = T_0, T_2
- b. Scan the log backwards from the end.
- c. $A = 500$; output the redo-only record $\langle T_2, A, 500 \rangle$
- d. output $\langle T_2 \text{ abort} \rangle$
- e. $B = 2000$; output the redo-only record $\langle T_0, B, 2000 \rangle$
- f. output $\langle T_0 \text{ abort} \rangle$

At the end of the recovery process, the state of the system is as follows:

$A = 500$
 $B = 2000$
 $C = 600$

The log records added during recovery are:

$\langle T_2, A, 500 \rangle$
 $\langle T_2 \text{ abort} \rangle$
 $\langle T_0, B, 2000 \rangle$
 $\langle T_0 \text{ abort} \rangle$

Observe that B is set to 2000 by two log records, one created during normal rollback of T_0 , and the other created during recovery, when the abort of T_0 is completed. Clearly the second one is redundant, although not incorrect. Optimizations described in the ARIES algorithm (and equivalent optimizations described in Section 16.7 for the case of logical operations) can help avoid carrying out redundant operations, which create such redundant log records.

- 16.19** Suppose there is a transaction that has been running for a very long time, but has performed very few updates.

- a. What effect would the transaction have on recovery time with the recovery algorithm of Section 16.4, and with the ARIES recovery algorithm.
- b. What effect would the transaction have on deletion of old log records?

Answer:

- a. If a transaction has been running for a very long time, with few updates, it means that during recovery, the undo phase will have to scan the log backwards till the beginning of this transaction. This will increase the recovery time in the case of the recovery algorithm of Section 16.4. However, in the case of ARIES the effect is not that bad as ARIES considers the LastLSN and PrevLSN values of transactions in the undo list during its backward scan, allowing it to skip intermediate records belonging to completed transactions.
- b. A long running transaction implies that no log records which are written after it started, can be deleted till it either commits or aborts. This might lead to a very large log file being generated, though most of the transactions in the log file have completed. This transaction becomes a bottleneck for deletion of old log records.

16.20 Consider the log in Figure 16.6. Suppose there is a crash during recovery, just before the operation abort log record is written for operation O_1 . Explain what would happen when the system recovers again.

Answer: **Errata note:** The question above in the book has the text “.. just before after the operation abort”; the word “after” should be deleted from the text.

There is no checkpoint in the log, so recovery starts from the beginning of the log, and replays each action that is found in the log.

The redo phase would add the following log records:

```
< T0, B, 2050 >
< T0, C, 600 >
< T1, C, 400 >
< T0, C, 500 >
```

At the end of the redo phase, the undo list contains transactions T_0 and T_1 , since their start records are found, but not their end of abort records.

During the undo phase, scanning backwards in the log, the following events happen:

```
< T0, C, 400 >
< T1, C, 600 > /* The operation end of T1.O2 is found, */
/* and logical undo adds +200 to the current value of C. */
/* Other log records of T1 are skipped till T1.O2 operation */
/* begin is found. Log records of other txns would be */
/* processed, but there are none here. */
< T1, O2, operation-abort > /* On finding T1.O2 operation begin */
```

```

< T1, abort > /* On finding T1 start */
/* Next, the operation end of T0.O1 is found, and */
/* logical undo adds +100 to the current value of C. */
< T0, C, 700 >
/* Other operations of T1 till T0.O1 begin are skipped */
/* And when T0.O1 operation begin is found: */
< T0, O1, operation-abort >
< T0, B, 2000 >
< T0, abort >

```

Finally the values of data items B and C would be 2000, and 700, which, which were their original values before T_0 or T_1 started.

- 16.21** Compare log-based recovery with the shadow-copy scheme in terms of their overheads, for the case when data is being added to newly allocated disk pages (in other words, there is no old value to be restored in case the transaction aborts).

Answer: In general, with logging each byte that is written is actually written twice, once as part of the page to which it is written, and once as part of the log record for the update. In contrast, shadow-copy schemes avoids log writes, at the expense of increased IO operations, and lower concurrency.

In the case of data added to newly allocated pages:

- There is no old value, so there is no need to manage a shadow copy for recovery, which benefits the shadow-copy scheme.
- Log-based recovery system would unnecessarily write old value, unless optimizations to skip the old value part are implemented for such newly allocated pages. Without such an optimization, logging has a higher overhead.
- Newly allocated pages are often formatted, for example by zeroing out all bytes, and then setting bytes corresponding to data structures in a slotted-page architecture. These operations also have to be logged, adding to the overhead of logging.
- However, the benefit of shadow-copy over logging is really significant only if the page is filled with a significant amount of data before it is written, since a lot of logging is avoided in this case.

- 16.22** In the ARIES recovery algorithm:

- If at the beginning of the analysis pass, a page is not in the checkpoint dirty page table, will we need to apply any redo records to it? Why?
- What is RecLSN, and how is it used to minimize unnecessary redos?

Answer:

- a. If a page is not in the checkpoint dirty page table at the beginning of the analysis pass, redo records prior to the checkpoint record need not be applied to it as it means that the page has been flushed to disk and been removed from the DirtyPageTable before the checkpoint. However, the page may have been updated after the checkpoint, which means it will appear in the dirty page table at the end of the analysis pass.

For pages that appear in the checkpoint dirty page table, redo records prior to the checkpoint may also need to be applied.

- b. The RecLSN is an entry in the DirtyPageTable, which reflects the LSN at the end of the log when the page was added to DirtyPageTable. During the redo pass of the ARIES algorithm, if the LSN of the update log record encountered, is less than the RecLSN of the page in DirtyPageTable, then that record is not redone but skipped. Further, the redo pass starts at RedoLSN, which is the earliest of the RecLSNs among the entries in the checkpoint DirtyPageTable, since earlier log records would certainly not need to be redone. (If there are no dirty pages in the checkpoint, the RedoLSN is set to the LSN of the checkpoint log record.)

16.23 Explain the difference between a system crash and a “disaster.”

Answer: In a system crash, the CPU goes down, and disk may also crash. But stable-storage at the site is assumed to survive system crashes. In a “disaster”, *everything* at a site is destroyed. Stable storage needs to be distributed to survive disasters.

16.24 For each of the following requirements, identify the best choice of degree of durability in a remote backup system:

- Data loss must be avoided but some loss of availability may be tolerated.
- Transaction commit must be accomplished quickly, even at the cost of loss of some committed transactions in a disaster.
- A high degree of availability and durability is required, but a longer running time for the transaction commit protocol is acceptable.

Answer:

- Two very safe is suitable here because it guarantees durability of updates by committed transactions, though it can proceed only if both primary and backup sites are up. Availability is low, but it is mentioned that this is acceptable.
- One safe committing is fast as it does not have to wait for the logs to reach the backup site. Since data loss can be tolerated, this is the best option.

- c. With two safe committing, the probability of data loss is quite low, and also commits can proceed as long as at least the primary site is up. Thus availability is high. Commits take more time than in the one safe protocol, but that is mentioned as acceptable.

16.25 The Oracle database system uses undo log records to provide a snapshot view of the database, under snapshot-isolation. The snapshot view seen by transaction T_i reflects updates of all transactions that had committed when T_i started, and the updates of T_i ; updates of all other transactions are not visible to T_i .

Describe a scheme for buffer handling whereby transactions are given a snapshot view of pages in the buffer. Include details of how to use the log to generate the snapshot view. You can assume that operations as well as their undo actions affect only one page.

Answer: First, determine if a transaction is currently modifying the buffer. If not, then return the current contents of the buffer. Otherwise, examine the records in the undo log pertaining to this buffer. Make a copy of the buffer, then for each relevant operation in the undo log, apply the operation to the buffer copy starting with the most recent operation and working backwards until the point at which the modifying transaction began. Finally, return the buffer copy as the snapshot buffer.

CHAPTER 17



Database-System Architectures

The chapter is suitable for an introductory course. We recommend covering it, at least as self-study material, since students are quite likely to use the non-centralized (particularly client-server) database architectures when they enter the real world. The material in this chapter could potentially be supplemented by the two-phase commit protocol (2PC), (Section 19.4.1 from Chapter 19) to give students an overview of the most important details of non-centralized database architectures.

Exercises

- 17.7 Why is it relatively easy to port a database from a single processor machine to a multiprocessor machine if individual queries need not be parallelized?

Answer: Porting is relatively easy to a shared memory multiprocessor machine. Databases designed for single-processor machines already provide multitasking, allowing multiple processes to run on the same processor in a time-shared manner, giving a view to the user of multiple processes running in parallel. Thus, coarse-granularity parallel machines logically appear to be identical to single-processor machines, making the porting relatively easy.

The only difference is that a multiprocessor machine has multiple processor caches, and depending on the specific machine architecture, it is possible that a memory write executed on one processor may not be visible to a read on another processor for some time. To ensure that writes are visible at other processors, the database system needs to execute a special instruction (often called a *fence* instruction) which flushes a data item from all processor caches, ensuring the next read sees the latest value. Typically this instruction is executed just before releasing a latch or lock.

Porting a database to a shared disk or shared nothing multiprocessor architecture is clearly a little harder.

- 17.8 Transaction-server architectures are popular for client-server relational databases, where transactions are short. On the other hand, data-server

architectures are popular for client-server object-oriented database systems, where transactions are expected to be relatively long. Give two reasons why data servers may be popular for object-oriented databases but not for relational databases.

Answer: Data servers are good if data transfer is small with respect to computation, which is often the case in applications of OODBs such as computer aided design. In contrast, in typical relational database applications such as transaction processing, a transaction performs little computation but may touch several pages, which will result in a lot of data transfer with little benefit in a data server architecture. Another reason is that structures such as indices are heavily used in relational databases, and will become spots of contention in a data server architecture, requiring frequent data transfer. There are no such points of frequent contention in typical OODB applications which involve complex data that is rarely updated concurrently by multiple processes.

- 17.9** What is lock de-escalation, and under what conditions is it required? Why is it not required if the unit of data shipping is an item?

Answer: In a client-server system with page shipping, when a client requests an item, the server typically grants a lock not on the requested item, but on the *page* having the item, thus implicitly granting locks on all the items in the page. The other items in the page are said to be *prefetched*. If some other client subsequently requests one of the prefetched items, the server may ask the owner of the page lock to transfer back the lock on this item. If the page lock owner doesn't need this item, it de-escalates the page lock that it holds, to item locks on all the items that it is actually accessing, and then returns the locks on the unwanted items. The server can then grant the latter lock request.

If the unit of data shipping is an item, there are no coarser granularity locks; even if prefetching is used, it is typically implemented by granting individual locks on each of the prefetched items. Thus when the server asks for a return of a lock, there is no question of de-escalation, the requested lock is just returned if the client has no use for it.

- 17.10** Suppose you were in charge of the database operations of a company whose main job is to process transactions. Suppose the company is growing rapidly each year, and has outgrown its current computer system. When you are choosing a new parallel computer, what measure is most relevant—speedup, batch scaleup, or transaction scaleup? Why?

Answer: With increasing scale of operations, we expect that the number of transactions submitted per unit time increases. On the other hand, we wouldn't expect most of the individual transactions to grow longer, nor would we require that a given transaction should execute more quickly now than it did before. Hence transaction scale-up is the most relevant measure in this scenario.

- 17.11 Database systems are typically implemented as a set of processes (or threads) sharing a shared memory area.
- How is access to the shared memory area controlled?
 - Is two-phase locking appropriate for serializing access to the data structures in shared memory? Explain your answer.

Answer:

- A locking system is necessary to control access to the shared data structures. Since many database transactions only involve reading from a data structure, **reader-writer** locks should be used, allowing multiple processes to concurrently read from a data structure by using the lock in **shared** mode. A process that wishes to modify the data structure needs to obtain an **exclusive** lock on the data structure which prohibits concurrent access from other reading or writing processes.
Locks used for controlling access to shared memory data structures are typically not held in a two-phase manner (as described below), and are often called *latches* to distinguish them from locks held in a two-phase manner.
- Shared memory areas are usually hot spots of contention, since every transaction needs to access the shared memory frequently. If such memory areas are locked in a two-phase manner, concurrency would be greatly reduced, reducing performance correspondingly. Instead, locks (latches) on shared memory data structures are released after performing operations on the data structure.
Serializability at a lower level, such as the exact layout of data on a page, or the structure of a B-tree or hash-index, is not important as long as the differences caused by non-serial execution are not visible at a higher level (typically, at the relational abstraction). Locks on tuples (or higher granularity) are retained in a two-phase manner to ensure serializability at the higher level.

- 17.12 Is it wise to allow a user process to access the shared memory area of a database system? Explain your answer.

Answer: No, the shared memory area may contain data that the user's process is not authorized to see, and thus allowing direct access to shared memory is a security risk.

- 17.13 What are the factors that can work against linear scaleup in a transaction processing system? Which of the factors are likely to be the most important in each of the following architectures: shared memory, shared disk, and shared nothing?

Answer: Increasing contention for shared resources prevents linear scale-up with increasing parallelism. In a shared memory system, contention for memory (which implies bus contention) will result in falling scale-up with

increasing parallelism. In a shared disk system, it is contention for disk and bus access which affects scale-up. In a shared-nothing system, inter-process communication overheads will be the main impeding factor. Since there is no shared memory, acquiring locks, and other activities requiring message passing between processes will take more time with increased parallelism.

- 17.14** Memory systems can be divided into multiple modules, each of which can be serving a separate request at a given time. What impact would such a memory architecture have on the number of processors that can be supported in a shared-memory system?

Answer: If all memory requests have to go through a single memory module, the memory module would become the bottleneck as the number of processors increases. After some point, adding processors will not result in any performance improvement. However, if the memory system is itself able to run multiple requests in parallel, a larger number of processors can be supported in a shared-memory system, providing better speedup and/or scaleup.

- 17.15** Consider a bank that has a collection of sites, each running a database system. Suppose the only way the databases interact is by electronic transfer of money between themselves, using persistent messaging. Would such a system qualify as a distributed database? Why?

Answer: In a distributed database, it should be possible to run queries across sites, and to run transactions across sites using protocols such as two-phase commit. Each site provides an environment for execution of both global transactions initiated at remote sites and local transactions. The system described in the question does not have these properties, and hence it cannot qualify as a distributed database.

CHAPTER 18



Parallel Databases

This chapter is suitable for an advanced course, but can also be used for independent study projects by students of a first course. The chapter covers several aspects of the design of parallel database systems — partitioning of data, parallelization of individual relational operations, and parallelization of relational expressions. The chapter also briefly covers some systems issues, such as cache coherency and failure resiliency.

The most important applications of parallel databases today are for warehousing and analyzing large amounts of data. Therefore partitioning of data and parallel query processing are covered in significant detail. Query optimization is also of importance, for the same reason. However, parallel query optimization is still not a fully solved problem; exhaustive search, as is used for sequential query optimization, is too expensive in a parallel system, forcing the use of heuristics.

The description of parallel query processing algorithms is based on the shared-nothing model. Students may be asked to study how the algorithms can be improved if shared-memory machines are used instead.

Exercises

- 18.9 For each of the three partitioning techniques, namely round-robin, hash partitioning, and range partitioning, give an example of a query for which that partitioning technique would provide the fastest response.

Answer:

Round robin partitioning:

When relations are large and queries read entire relations, round-robin gives good speed-up and fast response time.

Hash partitioning

For point queries on the partitioning attributes, this gives the fastest response, as each disk can process a different query simultaneously. If the hash partitioning is uniform, entire relation scans can be performed efficiently.

Range partitioning For range queries on the partitioning attributes, which access a few tuples, range partitioning gives the fastest response.

18.10 What factors could result in skew when a relation is partitioned on one of its attributes by:

- a. Hash partitioning?
- b. Range partitioning?

In each case, what can be done to reduce the skew?

Answer:

- a. Hash-partitioning:
Too many records with the same value for the hashing attribute, or a poorly chosen hash function without the properties of randomness and uniformity, can result in a skewed partition. To improve the situation, we should experiment with better hashing functions for that relation.
- b. Range-partitioning:
Non-uniform distribution of values for the partitioning attribute (including duplicate values for the partitioning attribute) which are not taken into account by a bad partitioning vector is the main reason for skewed partitions. Sorting the relation on the partitioning attribute and then dividing it into n ranges with equal number of tuples per range will give a good partitioning vector with very low skew.

18.11 Give an example of a join that is not a simple equi-join for which partitioned parallelism can be used. What attributes should be used for partitioning?

Answer: We give two examples of such joins.

- a. $r \bowtie_{(r.A=s.B) \wedge (r.A < s.C)} s$
Here we have an equi-join condition which can be executed first, and the extra conditions can be checked independently on each tuple in the join result. Partitioned parallelism is useful to execute the equi-join,
- b. $r \bowtie_{(r.A \geq (\lfloor s.B/20 \rfloor) * 20) \wedge (r.A < ((\lfloor s.B/20 \rfloor) + 1) * 20)} s$
This is a query in which an r tuple and an s tuple join with each other if they fall into the same range of values. Hence partitioned parallelism applies naturally to this scenario, even though the join is not an equi-join.

For both the queries, r should be partitioned on attribute A and s on attribute B . For the second query, the partitioning of s should actually be done on $(\lfloor s.B/20 \rfloor) * 20$.

18.12 Describe a good way to parallelize each of the following:

- a. The difference operation
- b. Aggregation by the **count** operation
- c. Aggregation by the **count distinct** operation
- d. Aggregation by the **avg** operation
- e. Left outer join, if the join condition involves only equality
- f. Left outer join, if the join condition involves comparisons other than equality
- g. Full outer join, if the join condition involves comparisons other than equality

Answer:

- a. We can parallelize the difference operation by partitioning the relations on all the attributes, and then computing differences locally at each processor. As in aggregation, the cost of transferring tuples during partitioning can be reduced by partially computing differences at each processor, before partitioning.
- b. Let us refer to the group-by attribute as attribute A , and the attribute on which the aggregation function operates, as attribute B . **count** is performed just like **sum** (mentioned in the book) except that, a count of the number of values of attribute B for each value of attribute A is transferred to the correct destination processor, instead of a sum. After partitioning, the partial counts from all the processors are added up locally at each processor to get the final result.
- c. For this, partial counts cannot be computed locally before partitioning. Each processor instead transfers all unique B values for each A value to the correct destination processor. After partitioning, each processor locally counts the number of unique tuples for each value of A , and then outputs the final result.
- d. This can again be implemented like **sum**, except that for each value of A , a **sum** of the B values as well as a **count** of the number of tuples in the group, is transferred during partitioning. Then each processor outputs its local result, by dividing the total sum by total number of tuples for each A value assigned to its partition.
- e. This can be performed just like partitioned natural join. After partitioning, each processor computes the left outer join locally using any of the strategies of Chapter 12.
- f. The left outer join can be computed using an extension of the Fragment-and-Replicate scheme to compute non equi-joins. Consider $r \bowtie s$. The relations are partitioned, and $r \bowtie s$ is computed at

each site. We also collect tuples from r that did not match any tuples from s ; call the set of these dangling tuples at site i as d_i . After the above step is done at each site, for each fragment of r , we take the intersection of the d_i 's from every processor in which the fragment of r was replicated. The intersections give the real set of dangling tuples; these tuples are padded with nulls and added to the result. The intersections themselves, followed by addition of padded tuples to the result, can be done in parallel by partitioning.

- g. The algorithm is basically the same as above, except that when combining results, the processing of dangling tuples must be done for both relations.

18.13 Describe the benefits and drawbacks of pipelined parallelism.

Answer:

- **Benefits:** No need to write intermediate relations to disk only to read them back immediately.
- **Drawbacks:**
 - a. Cannot take advantage of high degrees of parallelism, as typical queries do not have large number of operations.
 - b. Not possible to pipeline operators which need to look at all the input before producing any output.
 - c. Since each operation executes on a single processor, the most expensive ones take a long time to finish. Thus speed-up will be low despite the use of parallelism.

18.14 Suppose you wish to handle a workload consisting of a large number of small transactions by using shared-nothing parallelism.

- a. Is intraquery parallelism required in such a situation? If not, why, and what form of parallelism is appropriate?
- b. What form of skew would be of significance with such a workload?
- c. Suppose most transactions accessed one *account* record, which includes an *account_type* attribute, and an associated *account_type_master* record, which provides information about the account type. How would you partition and/or replicate data to speed up transactions? You may assume that the *account_type_master* relation is rarely updated.

Answer:

- a. Intraquery parallelism is probably not appropriate for this situation. Since each individual transaction is small, the overhead of parallelizing each query may exceed the potential benefits. Interquery parallelism would be a better choice, allowing many transactions to run in parallel.

- b. Partition skew can be a performance issue in this type of system, especially with the use of shared-nothing parallelism. A load imbalance amongst the processors of the distributed system can significantly reduce the speedup gained by parallel execution. For example, if all transactions happen to involve only the data in a single partition, the processors not associated with that partition will not be used at all.
 - c. Since *account_type_master* is rarely updated, it can be replicated in entirety across all nodes. If the *account* relation is updated frequently and accesses are well-distributed, it should be partitioned across nodes.
- 18.15** The attribute on which a relation is partitioned can have a significant impact on the cost of a query.
- a. Given a workload of SQL queries on a single relation, what attributes would be candidates for partitioning?
 - b. How would you choose between the alternative partitioning techniques, based on the workload?
 - c. Is it possible to partition a relation on more than one attribute? Explain your answer.

Answer:

- a. The candidate attributes would be
 - i. Attributes on which one or more queries has a selection condition. The corresponding selection condition can then be evaluated at a single processor, instead of being evaluated at all processors.
 - ii. Attributes involved in join conditions. If such an attribute is used for partitioning, it is possible to perform the join without repartitioning the relation. This effect is particularly beneficial for very large relations, for which repartitioning can be very expensive.
 - iii. Attributes involved in group-by clauses; similar to joins, it is possible to perform aggregation without repartitioning the corresponding relation.
- b. A cost-based approach works best in choosing between alternatives. In this approach, candidate partitioning choices are generated, and for each candidate the cost of executing all the queries/updates in a workload is estimated. The choice leading to the least cost is picked. One issue is that the number of candidate choices is generally very large. Algorithms and heuristics designed to limit the number of candidates for which costs need to be estimated are widely used in practice.

Another issue is that the workload may have a very large number of queries/updates. Techniques to reduce this number include the following (a) combining repeated occurrences of a query that only differ in constants, replacing them by one parametrized query along with a count of number of occurrences and (b) dropping queries which are very cheap anyway, or not likely to be affected by the partitioning choice.

- c. It is possible to partition a relation on more than one attribute, in two ways. One is to involve multiple attributes in a single composite partitioning key. The other way is to keep more than one copy of the same relation, partitioned in different ways. The latter approach is increases update costs, but can speed up some queries significantly.

CHAPTER 19



Distributed Databases

Distributed databases in general, and *heterogeneous* distributed databases in particular, are of increasing practical importance, as organizations attempt to integrate databases across physical and organizational boundaries. Such interconnection of databases to create a distributed or multidatabase is in fact proving crucial to competitiveness for many companies. This chapter reconsiders the issues addressed earlier in the text, such as query processing, recovery and concurrency control, from the standpoint of distributed databases.

This is a long chapter, and is appropriate only for an advanced course. Single topics may be chosen for inclusion in an introductory course. Good choices include distributed data storage, heterogeneity and two-phase commit.

Exercises

19.16 Discuss the relative advantages of centralized and distributed databases.

Answer:

- A distributed database allows a user convenient and transparent access to data which is not stored at the site, while allowing each site control over its own local data. A distributed database can be made more reliable than a centralized system because if one site fails, the database can continue functioning, but if the centralized system fails, the database can no longer continue with its normal operation. Also, a distributed database allows parallel execution of queries and possibly splitting one query into many parts to increase throughput.
- A centralized system is easier to design and implement. A centralized system is cheaper to operate because messages do not have to be sent.

19.17 Explain how the following differ: fragmentation transparency, replication transparency, and location transparency.

Answer:

- a. With fragmentation transparency, the user of the system is unaware of any fragmentation the system has implemented. A user may for-

multate queries against global relations and the system will perform the necessary transformation to generate correct output.

- b. With replication transparency, the user is unaware of any replicated data. The system must prevent inconsistent operations on the data. This requires more complex concurrency control algorithms.
- c. Location transparency means the user is unaware of where data are stored. The system must route data requests to the appropriate sites.

19.18 When is it useful to have replication or fragmentation of data? Explain your answer.

Answer: Replication is useful when there are many read-only transactions at different sites wanting access to the same data. They can all execute quickly in parallel, accessing local data. But updates become difficult with replication. Fragmentation is useful if transactions on different sites tend to access different parts of the database.

19.19 Explain the notions of transparency and autonomy. Why are these notions desirable from a human-factors standpoint?

Answer: Autonomy is the amount of control a single site has over the local database. It is important because users at that site want quick and correct access to local data items. This is especially true when one considers that local data will be most frequently accessed in a database. Transparency hides the distributed nature of the database. This is important because users should not be required to know about location, replication, fragmentation or other implementation aspects of the database.

19.20 If we apply a distributed version of the multiple-granularity protocol of Chapter 15 to a distributed database, the site responsible for the root of the DAG may become a bottleneck. Suppose we modify that protocol as follows:

- Only intention-mode locks are allowed on the root.
- All transactions are given all possible intention-mode locks on the root automatically.

Show that these modifications alleviate this problem without allowing any nonserializable schedules.

Answer: Serializability is assured since we have not changed the rules for the multiple granularity protocol. Since transactions are automatically granted all intention locks on the root node, and are not given other kinds of locks on it, there is no need to send any lock requests to the root. Thus the bottleneck is relieved.

19.21 Study and summarize the facilities that the database system you are using provides for dealing with inconsistent states that can be reached with lazy propagation of updates.

Answer:

PostgreSQL does not have built-in support for replication. However, there are several external projects that add replication support to the database engine.

- **Slony-I** adds basic master-slave replication functionality to PostgreSQL with updates to the replicated databases performed lazily. Slony-I uses the trigger mechanism of PostgreSQL to implement replication. A consistent view of the replicated data is provided by preserving transactions across the replicated sites. Thus, the replicated sites will always have a consistent, although potentially older, version of the data.
Slony-I does not support multi-master replication, and therefore does not have to deal with inconsistent states due to updates at multiple sites.
- **Postgres-R** extends PostgreSQL to add synchronous replication, thus avoiding the consistency issues of performing lazy updates. However, Postgres-R is experimental software, and not ready for production use.

Several databases support lazy multi-master replication, which involves the risk of inconsistent data. The first problem lies in detecting the conflict, and the second problem lies in resolving the conflict. Key values and timestamps are the primary means of detecting conflicts:

- Uniqueness conflicts are detected when there are duplicates on a primary key.
- Update conflicts are detected when two sites update the same item independently. Timestamps or version numbers sent with updates are used to detect such conflicts, by recording the timestamp/version number prior to and after the update.
- Delete conflicts occur when a transaction at one site deletes a tuple, which another site concurrently updates (before the delete propagated to that site).

Oracle provides several options for resolving conflicts:

- Latest time stamp: Most recent update wins
- Overwrite: Overwrites current value with new value, without checking for conflict.
- Discard: Ignores the value
- Additive: Difference of the two values is added to the current value (current value = current value + (new value - old value)). Alternatives include retaining the minimum or maximum value.

As an alternative to using the predefined conflict resolution policies in Oracle, administrators can create stored procedures to resolve each type of conflict on a relation.

In Microsoft SQL Server 2008, multi-master replication is called peer-to-peer replication. SQL Server supports conflict detection between all combinations of insert, delete and update, based on the primary key, and a node identifier plus a version number stored with each tuple. In case conflicts are detected, the default is to stop replication, but the default can be overridden using stored procedures to resolve conflicts.

In addition, SQL Server also supports a form of replication called Merge Replication, which is based on a single-publisher/multiple-subscriber model, where each subscriber can update data and send updates back to the publisher. The resolution model simply retains one of the conflicting updates, with the publisher having higher priority, and a priority scheme between subscribers used to choose the winning subscriber in case two subscribers updated the same data.

- 19.22** Discuss the advantages and disadvantages of the two methods that we presented in Section 19.5.2 for generating globally unique timestamps.

Answer: The centralized approach has the problem of a possible bottleneck at the central site and the problem of electing a new central site if it goes down. The distributed approach has the problem that many messages must be exchanged to keep the system fair, or one site can get ahead of all other sites and dominate the database.

- 19.23** Consider the relations:

employee (*name*, *address*, *salary*, *plant_number*)
machine (*machine_number*, *type*, *plant_number*)

Assume that the *employee* relation is fragmented horizontally by *plant_number*, and that each fragment is stored locally at its corresponding plant site. Assume that the *machine* relation is stored in its entirety at the Armonk site. Describe a good strategy for processing each of the following queries.

- a. Find all employees at the plant that contains machine number 1130.
- b. Find all employees at plants that contain machines whose type is “milling machine.”
- c. Find all machines at the Almaden plant.
- d. Find employee \bowtie machine.

Answer:

- a. i. Perform $\Pi_{\text{plant_number}} (\sigma_{\text{machine_number}=1130} (\text{machine}))$ at Armonk.

- ii. Send the query $\Pi_{name} (employee)$ to all site(s) which are in the result of the previous query.
 - iii. Those sites compute the answers.
 - iv. Union the answers at the destination site.
- b. This strategy is the same as in part *a* of this exercise, except the first step should be to perform
- $$\Pi_{plant_number} (\sigma_{type="milling\ machine"} (machine)) \text{ at Armonk.}$$
- c. i. Perform $\sigma_{plant_number = x} (machine)$ at Armonk, where x is the plant number for Almaden.
- ii. Send the answers to the destination site.
- d. Strategy 1:
- i. Group *machine* at Armonk by plant number.
 - ii. Send the groups to the sites with the corresponding plant number.
 - iii. Perform a local join between the local data and the received data.
 - iv. Union the results at the destination site.

Strategy 2:

Send the *machine* relation at Armonk, and all the fragments of the *employee* relation to the destination site. Then perform the join at the destination site.

There is parallelism in the join computation according to the first strategy but not in the second. Nevertheless, in a WAN the amount of data to be shipped is the main cost factor. We expect that each plant will have more than one machine, hence the result of the local join at each site will be a cross-product of the employee tuples and machines at that plant. This cross-product's size is greater than the size of the *employee* fragment at that site. As a result the second strategy will result in less data shipping, and will be more efficient.

19.24 For each of the strategies of Exercise 19.23, state how your choice of a strategy depends on:

- a. The site at which the query was entered.
- b. The site at which the result is desired.

Answer:

- a. Assuming that the cost of shipping the query itself is minimal, the site at which the query was submitted does not affect our strategy for query evaluation.
- b. For the first query, we find out the plant numbers where the machine number 1130 is present, at Armonk. Then the employee tuples at all those plants are shipped to the destination site. We can see that this strategy is more or less independent of the destination site. The same can be said of the second query. For the third query, the selection is performed at Armonk and results shipped to the destination site. This strategy is obviously independent of the destination site. For the fourth query, we have two strategies. The first one performs local joins at all the plant sites and their results are unioned at the destination site. In the second strategy, the *machine* relation at Armonk as well as all the fragments of the *employee* relation are first shipped to the destination, where the join operation is performed. There is no obvious way to optimize these two strategies based on the destination site. In the answer to Exercise 19.23 we saw the reason why the second strategy is expected to result in less data shipping than the first. That reason is independent of destination site, and hence we can in general prefer strategy two to strategy one, regardless of the destination site.

19.25 Is the expression $r_i \bowtie r_j$ necessarily equal to $r_j \bowtie r_i$? Under what conditions does $r_i \bowtie r_j = r_j \bowtie r_i$ hold?

Answer: In general, $r_i \bowtie r_j \neq r_j \bowtie r_i$. This can be easily seen from Exercise 19.11, in which $r \bowtie s \neq s \bowtie r$. $r \bowtie s$ was given in 19.11, while

$s \bowtie r =$

| C | D | E |
|---|---|---|
| 3 | 4 | 5 |
| 3 | 6 | 8 |
| 2 | 3 | 2 |

By definition, $r_i \bowtie r_j = \Pi_{R_i}(r_i \bowtie r_j)$ and $r_j \bowtie r_i = \Pi_{R_j}(r_i \bowtie r_j)$, where R_i and R_j are the schemas of r_i and r_j respectively. For $\Pi_{R_i}(r_i \bowtie r_j)$ to be always equal to $\Pi_{R_j}(r_i \bowtie r_j)$, the schemas R_i and R_j must be the same.

19.26 If a cloud data-storage service is used to store two relations r and s and we need to join r and s , why might it be useful to maintain the join as a materialized view? In your answer, be sure to distinguish among various meanings of “useful”: overall throughput, efficient use of space, and response time to user queries.

Answer: Performing a join on a cloud data-storage system can be very expensive, if either of the relations to be joined is partitioned on attributes

other than the join attributes, since a very large amount of data would need to be transferred to perform the join. However, if $r \bowtie s$ is maintained as a materialized view, it can be updated at a relatively low cost each time either r or s is updated, instead of incurring a very large cost when the query is executed. Thus, queries are benefitted at some cost to updates.

Considering the various notions of usefulness, with the materialized view, overall throughput will be much better if the join query is executed reasonably often relative to updates, but may be worse if the join is rarely used, but updates are frequent.

The materialized view will certainly require extra space, but given that disk capacities are very high relative to IO (seek) operations and transfer rates, the extra space is likely to not be an major overhead since a large number of disks are needed anyway to handle the IO load and data transfer load.

The materialized view will obviously be very useful to evaluate join queries, reducing time greatly by avoiding a large amount of data transfer across machines.

- 19.27** Why do cloud-computing services support traditional database systems best by using a virtual machine instead of running directly on the service provider's actual machine?

Answer: By using a virtual machine, if a physical machine fails, virtual machines running on that physical machine can be restarted quickly on one or more other physical machines, improving availability. (Assuming of course that data remains accessible, either by storing multiple copies of data, or by storing data in an highly available external storage system.)

- 19.28** Describe how LDAP can be used to provide multiple hierarchical views of data, without replicating the base-level data.

Answer: This can be done using referrals. For example an organization may maintain its information about departments either by geography (i.e. all departments in a site of the the organization) or by structure (i.e. information about a department from all sites). These two hierarchies can be maintained by defining two different schemas with department information at a site as the base information. The entries in the two hierarchies will refer to the base information entry using referrals.

CHAPTER 20



Data Warehousing and Mining

This chapter covers data warehousing and data mining. Considering the growing importance of all the topics covered in this chapter, some of the sections of the chapter can be assigned as supplementary reading material, even in an introductory course. The material in the chapter is also suitable for laying the groundwork for an advanced course, or for professionals to keep in touch with recent developments.

Exercises

- 20.6 Draw a diagram that shows how the *classroom* relation of our university example as shown in Appendix A would be stored under a column-oriented storage structure.

Answer: The relation would be stored in three files, one per attribute, as shown below. We assume that the row number can be inferred implicitly from position, by using fixed size space for each attribute. Otherwise, the row number would also have to be stored explicitly.

| <i>building</i> |
|-----------------|
| Packard |
| Painter |
| Taylor |
| Watson |
| Watson |

| <i>room_number</i> |
|--------------------|
| 101 |
| 514 |
| 3128 |
| 100 |
| 120 |

| capacity |
|----------|
| 500 |
| 10 |
| 70 |
| 30 |
| 50 |

20.7 Explain why the nested-loops join algorithm (see Section 12.5.1) would work poorly on database stored in a column-oriented manner. Describe an alternative algorithm that would work better and explain why your solution is better.

Answer: If the nested-loops join algorithm is used as is, it would require tuples for each of the relations to be assembled before they are joined. Assembling tuples can be expensive in a column store, since each attribute may come from a separate area of the disk; the overhead of assembly would be particularly wasteful if many tuples do not satisfy the join condition and would be discarded. In such a situation it would be better to first find which tuples match by accessing only the join columns of the relations. Sort-merge join, hash join, or indexed nested loops join can be used for this task. After the join is performed, only tuples that get output by the join need to be assembled; assembly can be done by sorting the join result on the record identifier of one of the relations and accessing the corresponding attributes, then resorting on record identifiers of the other relation to access its attributes.

20.8 Construct a decision-tree classifier with binary splits at each node, using tuples in relation $r(A, B, C)$ shown below as training data; attribute C denotes the class. Show the final tree, and with each node show the best split for each attribute along with its information gain value.

(1, 2, a), (2, 1, a), (2, 5, b), (3, 3, b), (3, 6, b),
(4, 5, b), (5, 5, c), (6, 3, b), (6, 7, c)

Answer: Figure 20.1 shows one possible decision tree for the data. Using the Gini purity metric, the purity of the initial data set is

$$1 - \sum_{i=1}^k p_i^2 = 1 - ((\frac{2}{9})^2 + (\frac{5}{9})^2 + (\frac{2}{9})^2) = 0.595259$$

The first branch splits on $B \leq 2$, giving a purity score of $1 - 1^2 = 0$ for those attributes with $B \leq 2$ (all are classified as a), and a purity score of

$$1 - ((\frac{2}{7})^2 + (\frac{5}{7})^2) = 0.40816$$

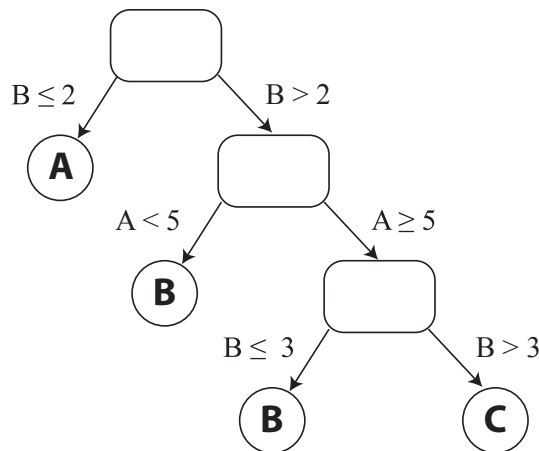


Figure 20.1 Decision tree for data on relation $r(A, B, C)$

for the remaining items. The weighted purity of the entire set is

$$\frac{2}{9} * 0 + \frac{7}{9} * 0.40816 = 0.31746$$

The information gain from this split is $0.595259 - 0.31746 = 0.27513$. Next, we split on $A < 5$. The 4 data items with $A < 5$ all have class b , and thus have purity 0. The remaining 3 items have purity

$$1 - ((\frac{1}{3})^2 + (\frac{2}{3})^2) = 0.44444$$

The weighted purity of these sets is

$$\frac{4}{7} * 0 + \frac{3}{7} * 0.44444 = 0.19048$$

The information gain from the second split is $0.40816 - 0.19048 = 0.21769$. Finally, we split on $B \leq 3$. One data item satisfies this predicate and has class b . The other two items both have class c . The purity of these two sets is 0. The information gain from the final split is $0.21769 - 0 = 0.21769$.

- 20.9** Suppose half of all the transactions in a clothes shop purchase jeans, and one third of all transactions in the shop purchase T-shirts. Suppose also that half of the transactions that purchase jeans also purchase T-shirts. Write down all the (nontrivial) association rules you can deduce from the above information, giving support and confidence of each rule.

Answer: The rules are as follows. The last rule can be deduced from the previous ones.

| Rule | Support | Conf. |
|---|---------|-------|
| $\forall \text{ transactions } T, \text{ true} \Rightarrow \text{buys}(T, \text{jeans})$ | 50% | 50% |
| $\forall \text{ transactions } T, \text{ true} \Rightarrow \text{buys}(T, \text{t-shirts})$ | 33% | 33% |
| $\forall \text{ transactions } T, \text{buys}(T, \text{jeans}) \Rightarrow \text{buys}(T, \text{t-shirts})$ | 25% | 50% |
| $\forall \text{ transactions } T, \text{buys}(T, \text{t-shirts}) \Rightarrow \text{buys}(T, \text{jeans})$ | 25% | 75% |

20.10 Consider the problem of finding large itemsets.

- Describe how to find the support for a given collection of itemsets by using a single scan of the data. Assume that the itemsets and associated information, such as counts, will fit in memory.
- Suppose an itemset has support less than j . Show that no superset of this itemset can have support greater than or equal to j .

Answer:

- Let $\{S_1, S_2, \dots, S_n\}$ be the collection of item-sets for which we want to find the support. Associate a counter $\text{count}(S_i)$ with each item-set S_i .
Initialize each counter to zero. Now examine the transactions one-by-one. Let $S(T)$ be the item-set for a transaction T . For each item-set S_i that is a subset of $S(T)$, increment the corresponding counter $\text{count}(S_i)$.
When all the transactions have been scanned, the values of $\text{count}(S_i)$ for each i will give the support for item-set S_i .
- Let A be an item-set. Consider any item-set B which is a superset of A . Let τ_A and τ_B be the sets of transactions that purchase all items in A and all items in B , respectively. For example, suppose A is $\{a, b, c\}$, and B is $\{a, b, c, d\}$.
A transaction that purchases all items from B must also have purchased all items from A (since $A \subseteq B$). Thus, every transaction in τ_B is also in τ_A . This implies that the number of transactions in τ_B is at most the number of transactions in τ_A . In other words, the support for B is at most the support for A .
Thus, if any item-set has support less than j , all supersets of this item-set have support less than j .

20.11 Create a small example of a set of transactions showing that although many transactions contain two items, that is, the itemset containing the two items has a high support, purchase of one of the items may have a negative correlation with purchase of the other.

Answer: The following set of transactions involve fruit purchases:

| Transaction_ID | Item |
|----------------|--------|
| T-1 | orange |
| T-1 | banana |
| T-1 | apple |
| T-2 | orange |
| T-2 | banana |
| T-3 | orange |
| T-3 | apple |
| T-4 | orange |
| T-4 | banana |
| T-4 | grapes |
| T-5 | banana |
| T-5 | apple |
| T-6 | banana |
| T-6 | grapes |

Consider the association rule

$$orange \Rightarrow banana$$

This rule is satisfied in 3 out of the 6 transactions, so the support value is 50 percent. However, the correlation between purchasing oranges and purchasing bananas in this data set is -0.32 .

- 20.12** The organization of parts, chapters, sections, and subsections in a book is related to clustering. Explain why, and to what form of clustering.

Answer: The organization of a book's content is a form of **hierarchical clustering**. Contents within a single subsection are closely related, whereas different parts of a book cover a more diverse range of topics.

- 20.13** Suggest how predictive mining techniques can be used by a sports team, using your favorite sport as an example.

Answer: Given the large amount of statistics collected during and about sporting events, there are many ways a sports team can make use of predictive data mining:

- Some players may be more effective in certain situations or environments, so data mining can predict when each player should be used.
- Specific strategies may be more effective against certain teams or during certain situations in the game.
- Predictive mining can estimate the outcome of a match beforehand, information which could be useful to a team before entering a tournament.

CHAPTER 21



Information Retrieval

This chapter covers advanced querying techniques for databases and information retrieval. Advanced querying techniques include decision support systems, online analytical processing, including SQL:1999 support for OLAP, and data mining.

Although information retrieval has been considered as a separate field from databases in the research community, there are strong connections. Distributed information retrieval is growing in importance with the explosion of documents on the world wide web and the resultant importance of web search techniques.

Considering the growing importance of all the topics covered in this chapter, some of the sections of the chapter can be assigned as supplementary reading material, even in an introductory course. These could include OLAP, some parts of data mining, and some parts of information retrieval. The material in the chapter is also suitable for laying the groundwork for an advanced course, or for professionals to keep in touch with recent developments.

Exercises

- 21.6 Using a simple definition of term frequency as the number of occurrences of the term in a document, give the TF-IDF scores of each term in the set of documents consisting of this and the next exercise.

Answer:

Term frequency $TF(d, t) = \log(1 + n(d, t)/n(d))$

where $n(d, t)$ denotes the number of occurrences of term t in the document d and $n(d)$ denotes the number of terms in the document.

using - $\log(1 + 1/75)$

a - $\log(1 + 5/75)$

simple - $\log(1 + 1/75)$

definition - $\log(1 + 1/75)$

of - $\log(1 + 6/75)$

term - $\log(1 + 3/75)$

frequency - $\log(1 + 1/75)$

as - $\log(1 + 1/75)$
 the - $\log(1 + 1/75)$
 number - $\log(1 + 7/75)$
 occurrences - $\log(1 + 1/75)$
 in - $\log(1 + 2/75)$
 document - $\log(1 + 1/75)$
 give - $\log(1 + 1/75)$
 TFIDF - $\log(1 + 1/75)$
 scores - $\log(1 + 1/75)$
 each - $\log(1 + 3/75)$
 set - $\log(1 + 1/75)$
 documents - $\log(1 + 3/75)$
 consisting - $\log(1 + 1/75)$
 this - $\log(1 + 1/75)$
 and - $\log(1 + 2/75)$
 next - $\log(1 + 1/75)$
 exercise - $\log(1 + 1/75)$
 create - $\log(1 + 1/75)$
 small - $\log(1 + 2/75)$
 example - $\log(1 + 1/75)$
 4 - $\log(1 + 1/75)$
 with - $\log(1 + 1/75)$
 PageRank - $\log(1 + 3/75)$
 inverted - $\log(1 + 1/75)$
 lists - $\log(1 + 1/75)$
 sorted - $\log(1 + 1/75)$
 by - $\log(1 + 1/75)$
 you - $\log(1 + 1/75)$
 do - $\log(1 + 1/75)$
 not - $\log(1 + 1/75)$
 need - $\log(1 + 1/75)$
 to - $\log(1 + 1/75)$
 compute - $\log(1 + 1/75)$
 just - $\log(1 + 1/75)$
 assume - $\log(1 + 1/75)$
 some - $\log(1 + 1/75)$
 values - $\log(1 + 1/75)$
 page - $\log(1 + 1/75)$

- 21.7 Create a small example of four small documents, each with a PageRank, and create inverted lists for the documents sorted by the PageRank. You do not need to compute PageRank, just assume some values for each page.
- Answer:** Given 4 documents - A, B, C, D where the PageRanks are decreasing in that order, which means A has the highest PageRank and D has the lowest PageRank. We have, pages that are pointed to from more

web pages have higher PageRank. Similarly, pages pointed to by Web pages with a high PageRank will also have a higher PageRank. One way of creating an inverted list is:

- a. B, C, D all point to A . $A \leftarrow B, A \leftarrow C, A \leftarrow D$.
- b. A points to B . $B \leftarrow A$.
- c. B points to C . $C \leftarrow B$.
- d. C points to D . $D \leftarrow C$.

- 21.8** Suppose you wish to perform keyword querying on a set of tuples in a database, where each tuple has only a few attributes, each containing only a few words. Does the concept of term frequency make sense in this context? And that of inverse document frequency? Explain your answer. Also suggest how you can define the similarity of two tuples using TF-IDF concepts.

Answer: Term frequency is the logarithm of the number of occurrences of the term divided by the number of terms in the document. When it comes to small databases with few attributes, each containing only a few words, the concept of term frequency may not make sense. The relevance of a term may not depend on the number of occurrences of the term, and also when the domain is very small the logarithmic increase we used in the term frequency may not be a good indicator.

The inverse document frequency which is the inverse of the number of documents that contain this term may not also be very relevant in this case. For example, the primary key value and some other key may be having the inverse document frequency of 1, but we can't assume their weights to be equal.

The similarity of the two tuples can be measured by the *cosine similarity* metric. But one major difference is only values that belong to the same attribute should be considered. Two different attributes from two tuples may be having the same value, but that doesn't increase the similarity factor.

- 21.9** Web sites that want to get some publicity can join a Web ring, where they create links to other sites in the ring, in exchange for other sites in the ring creating links to their site. What is the effect of such rings on popularity ranking techniques such as PageRank?

Answer: PageRank is a measure of popularity of a page based on the popularity of the pages that link to the page. It may be noted that the pages that are pointed to from many Web pages are more likely to be visited, and thus will have a higher PageRank. Similarly, pages pointed to by Web pages with a high PageRank will also have a higher probability of being visited, and thus will have a higher PageRank. In the given scenario where Web sites join a Web ring and create links to other sites, the PageRank of all the pages increases. The number of links referencing to each page increases, which only increases the PageRank.

- 21.10** The Google search engine provides a feature whereby Web sites can display advertisements supplied by Google. The advertisements supplied are based on the contents of the page. Suggest how Google might choose which advertisements to supply for a page, given the page contents.

Answer: Google might use the concepts in similarity based retrieval. Here, they can give the system a document A and the set of advertisements B, and ask the system to retrieve advertisements that are similar to A. One approach is to find k terms in A with highest values of $TF(A, t) * IDF(t)$, and to use these k terms as a query to find relevance of other documents. The metric '*cosinesimilarity*' can also be used to determine which advertisements to supply for a page, given the page contents.

- 21.11** One way to create a keyword-specific version of PageRank is to modify the random jump such that a jump is only possible to pages containing the keyword. Thus pages that do not contain the keyword but are close (in terms of links) to pages that contain the keyword also get a nonzero rank for that keyword.
- Give equations defining such a keyword-specific version of PageRank.
 - Give a formula for computing the relevance of a page to a query containing multiple keywords.

Answer:

- Give equations defining such a keyword-specific version of PageRank.

$$P[j] = \delta / N_i + (1 - \delta) * \sum_{i=1}^N (T[i, j] * P[i])$$

where δ is a constant between 0 and 1, N is the number of pages, N_i is the number of pages containing the keyword; δ represents the probability of a step in the walk being a jump to a page containing the keyword.

- Give a formula for computing the relevance of a page to a query containing multiple keywords.
The relevance of a document to a query containing multiple keywords is estimated by combining by adding the relevance measures of the page to each keyword. Some weights also might be considered.

$$r(d, Q) = \sum_{t=1}^n TF(d, t) * IDF(t)$$

where Q is the set of keywords of size n , TF is the *term frequency* and IDF is the *inversedocument frequency*. This measure can be further

refined if the user is permitted to specify weights $w(t)$ for terms in the query, in which case the user specified weights are also taken into account by multiplying $TF(d, t)$ by $w(t)$ in the above formula.

21.12 The idea of popularity ranking using hyperlinks can be extended to relational and XML data, using foreign key and IDREF edges in place of hyperlinks. Suggest how such a ranking scheme may be of value in the following applications:

- a. A bibliographic database that has links from articles to authors of the articles and links from each article to every article that it references.
- b. A sales database that has links from each sales record to the items that were sold.

Also suggest why prestige ranking can give less than meaningful results in a movie database that records which actor has acted in which movies.

Answer:

- a. A bibliographic database, which has links from articles to authors of the articles and links from each article to every article that it references.

This helps us in ranking articles according to their popularity. If an article is referenced by many articles, then its more popular, so each article has a rank associated with it. and we could also find the authors for those articles.

- b. A sales database which has links from each sales record to the items that were sold.

This helps us in determining which items are the most popular. If the item is referenced by many sales records, then its more popular. So, ranking the database helps in determining the popularity of the items that were sold.

A movie which has many actors associated with it is deemed to be more popular when prestige ranking is taken into account. Same goes with the actor as well. But that may not be true in real life, the popularity of a movie or that of an actor cannot be determined by ranking the movie lists which map the actors to the movies.

21.13 What is the difference between a false positive and a false drop? If it is essential that no relevant information be missed by an information retrieval query, is it acceptable to have either false positives or false drops? Why?

Answer: False drop - A few relevant documents may not be retrieved. False positive - A few irrelevant documents may be retrieved. It is acceptable to have false positives but not any false drops when it is essential that no relevant information is missed because by permitting false positive; the system can later filter the results away later by looking at the keywords than they actually contain, but by permitting false drops, some

relevant information is missed out. By allowing false positives and not allowing false drops, no relevant information is missed out.

CHAPTER 22



Object-Based Databases

This chapter describes extensions to relational database systems to provide complex data types and object-oriented features. Such extended systems are called object-relational systems. Since the chapter was introduced in the 3th edition most commercial database systems have added some support for object-relational features, and these features have been standardized as part of SQL:1999.

It would be instructive to assign students exercises aimed at finding applications where the object-relational model, in particular complex objects, would be better suited than the traditional relational model.

Exercises

- 22.7 Redesign the database of Practice Exercise 22.2 into first normal form and fourth normal form. List any functional or multivalued dependencies that you assume. Also list all referential-integrity constraints that should be present in the first and fourth normal form schemas.

Answer: To put the schema into first normal form, we flatten all the attributes into a single relation schema.

Employee-details = (*ename, cname, bday, bmonth, byear, stype, xyear, xcity*)

We rename the attributes for the sake of clarity. *cname* is *Children.name*, and *bday, bmonth, byear* are the *Birthday* attributes. *stype* is *Skills.type*, and *xyear* and *xcity* are the *Exams* attributes. The FDs and multivalued dependencies we assume are:-

ename, cname → *bday, bmonth, byear*
ename →→ *cname, bday, bmonth, byear*
ename, stype →→ *xyear, xcity*

The FD captures the fact that a child has a unique birthday, under the assumption that one employee cannot have two children of the same

name. The MVDs capture the fact there is no relationship between the children of an employee and his or her skills-information. The redesigned schema in fourth normal form is:-

Employee = (*ename*)
Child = (*ename*, *cname*, *bday*, *bmonth*, *byear*)
Skill = (*ename*, *stype*, *xyear*, *xcity*)

ename will be the primary key of *Employee*, and (*ename*, *cname*) will be the primary key of *Child*. The *ename* attribute is a foreign key in *Child* and in *Skill*, referring to the *Employee* relation.

22.8 Consider the schema from Practice Exercise 22.2.

- a. Give SQL DDL statements to create a relation *EmpA* which has the same information as *Emp*, but where multiset-valued attributes *ChildrenSet*, *SkillsSet* and *ExamsSet* are replaced by array-valued attributes *ChildrenArray*, *SkillsArray* and *ExamsArray*.
- b. Write a query to convert data from the schema of *Emp* to that of *EmpA*, with the array of children sorted by birthday, the array of skills by the skill type and the array of exams by the year.
- c. Write an SQL statement to update the *Emp* relation by adding a child Jeb, with a birthdate of February 5, 2001, to the employee named George.
- d. Write an SQL statement to perform the same update as above but on the *EmpA* relation. Make sure that the array of children remains sorted by year.

Answer:

- a.


```

create type Exams (year int, city varchar(30))
create type SkillsA (type varchar(30),
  ExamsArray Exams array [20])
create type Children (name varchar(30), birthday date)
create table EmpA (ename varchar(30),
  ChildrenArray Children array [10],
  SkillArray SkillsA array [25])
      
```
- b.


```

select ename,
  array(select name, birthday
    from unnest(E.ChildrenSet) as CS
    order by CS.birthday) as ChildrenArray,
  array(select type,
    array(select year, city
      
```

```

        from unnest(SS.ExamSet)
        order by SS.year) as ExamsArray
    from unnest(E.SkillSet) as SS)
    as SkillsArray
from Emp as E

```

c.

```

update Emp
set ChildrenSet = ChildrenSet union
    multiset[('Jeb', '2/5/2001')]
where ename = 'George'

```

- d. We make use of the infix array concatenation operator, "||", which takes two arrays as arguments and returns the concatenation of the two arrays.

```

update EmpA
set ChildrenArray = array(
    select name, birthday
    from unnest(ChildrenArray ||
        array[('Jeb', '2/5/2001')])
    order by birthday)
where ename = 'George'

```

- 22.9 Consider the schemas for the table *people*, and the tables *students* and *teachers*, which were created under *people*, in Section 22.4. Give a relational schema in third normal form that represents the same information. Recall the constraints on subtables, and give all constraints that must be imposed on the relational schema so that every database instance of the relational schema can also be represented by an instance of the schema with inheritance.

Answer: A corresponding relational schema in third normal form is given below:-

```

People = (name, address)
Students = (name, degree, student-department)
Teachers = (name, salary, teacher-department)

```

name is the primary key for all the three relations, and it is also a foreign key referring to *People*, for both *Students* and *Teachers*.

Instead of placing only the *name* attribute of *People* in *Students* and *Teachers*, both its attributes can be included. In that case, there will be a slight change, namely – (*name, address*) will become the foreign key in *Students* and *Teachers*. The primary keys will remain the same in all tables.

22.10 Explain the distinction between a type x and a reference type $\text{ref}(x)$. Under what circumstances would you choose to use a reference type?

Answer: If the type of an attribute is x , then in each tuple of the table, corresponding to that attribute, there is an actual object of type x . If its type is $\text{ref}(x)$, then in each tuple, corresponding to that attribute, there is a *reference* to some object of type x . We choose a reference type for an attribute, if that attribute's intended purpose is to refer to an independent object.

22.11 Consider the E-R diagram in Figure 22.7, which contains specializations, using subtypes and subtables.

- Give an SQL schema definition of the E-R diagram.
- Give an SQL query to find the names of all people who are not secretaries.
- Give an SQL query to print the names of people who are neither employees nor students.
- Can you create a person who is an employee and a student with the schema you created? Explain how, or explain why it is not possible.

Answer:

```
a. create type Person
    ID varchar(10),
    name varchar(30),
    address varchar(40))
create type Employee
    under Person
    (salary integer)
create type Student
    under Person
    (tot_credits integer)
create type Instructor
    under Employee
    (rank varchar(10))
create type Secretary
    under Employee
    (hours_per_week integer)

create table person of Person
create table employee of Employee
    under person
create table student of Student
    under person
create table instructor of Instructor
    under employee
```

create table *secretary* **of** *Secretary*
under *employee*

b.

```
select *
from person
where ID not in
      (select ID from secretary)
```

c. Give an SQL query to print the names of people who are neither employees nor students.

```
select *
from only person
```

d. It is not possible to create a person who is an employee and a student, since there is no most-specific type in our schema corresponding to such a person. To do so, we would have to create a corresponding type (such as TeachingAssistant) using multiple inheritance, and a corresponding table that is under employee and student. However, SQL does not permit multiple inheritance of tables, so this is not possible in SQL.

22.12 Suppose a JDO database had an object *A*, which references object *B*, which in turn references object *C*. Assume all objects are on disk initially. Suppose a program first dereferences *A*, then dereferences *B* by following the reference from *A*, and then finally dereferences *C*. Show the objects that are represented in memory after each dereference, along with their state (hollow or filled, and values in their reference fields).

Answer: See figures 22.1 through 22.3. Gray boxes indicate persistent objects and white boxes indicate hollow objects.



Figure 22.1 State of the program after *A* is referenced.



Figure 22.2 State of the program after *B* is referenced.



Figure 22.3 State of the program after *C* is referenced.

CHAPTER 23



XML

XML is today widely used in the exchange of data between applications, and for storing data, whether simple or complex, in flat files. All the recent standard formats for storing documents, spreadsheets, presentations, and so on are based on XML, cementing the success of XML. In the context of databases, XML has been quite successful, although not with anything like the success of XML outside of databases; relational databases continue to be used for most mission critical data. However, support for storing XML in databases has improved significantly in the last 5 years.

Our view of XML is decidedly database centric. In this view, XML is a data model that provides a number of features beyond that provided by the relational model, in particular the ability to package related information into a single unit, by using nested structures. Specific application domains for data representation and interchange need their own standards that define the data schema.

Given the extensive nature of XML and related standards, this chapter only attempts to provide an introduction, and does not attempt to provide a complete description. For a course that intends to explore XML in detail, supplementary material may be required. These could include online information on XML and books on XML.

Exercises

- 23.10** Show, by giving a DTD, how to represent the non-1NF *books* relation from Section 22.2, using XML.

Answer:

```
<!DOCTYPE bib [
  <!ELEMENT book (title, author+, publisher, keyword+)>
  <!ELEMENT publisher (pub-name, pub-branch) >
  <!ELEMENT title ( #PCDATA )>
  <!ELEMENT author ( #PCDATA )>
  <!ELEMENT keyword ( #PCDATA )>
  <!ELEMENT pub-name( #PCDATA )>
  <!ELEMENT pub-branch( #PCDATA )>
]>
```

- 23.11** Write the following queries in XQuery, assuming the schema from Practice Exercise 23.2.

- Find the names of all employees who have a child who has a birthday in March.
- Find those employees who took an examination for the skill type “typing” in the city “Dayton”.
- List all skill types in *Emp*.

Answer:

- Find the names of all employees who have a child who has a birthday in March.

```
for $e in /db/emp,
  $m in distinct-values($e/children/birthday/month)
where $m = 'March'
return $e/ename
```

- Find those employees who took an examination for the skill type “typing” in the city “Dayton”.

```
for $e in /db/emp
  $s in $e/skills[type='typing']
  $exam in $s/exams
where $exam/city= 'Dayton'
return $e/ename
```

- List all skill types in *Emp*.

```
for $t in distinct-values (/db/emp/skills/type)
return $t
```

- 23.12** Consider the XML data shown in Figure 23.3. Suppose we wish to find purchase orders that ordered two or more copies of the part with identifier 123. Consider the following attempt to solve this problem:

```
for $p in purchaseorder
where $p/part/id = 123 and $p/part/quantity >= 2
return $p
```

Explain why the query may return some purchase orders that order less than two copies of part 123. Give a correct version of the above query.

Answer: Reason:

The expression $x = y$ evaluates to true if any of the values returned by the first expression is equal to any one of the values returned by the second expression. In this case, if any of the values in the $\text{part/quantity} \geq 2$ then it evaluates to true, so the query returns parts which have $\text{id} = 123$ irrespective of the quantity.

Query:

```
for $b in purchaseorder
where some $p in $b/part satisfies
    $p/id = 123 AND $p/quantity >= 2
return {$b}
```

- 23.13** Give a query in XQuery to flip the nesting of data from Exercise 23.10. That is, at the outermost level of nesting the output must have elements corresponding to authors, and each such element must have nested within it items corresponding to all the books written by the author.

Answer:

```
<bib>
  for $x in distinct-values(/books/author)
  return
    <author> <name> { $x } </name>
    { for $y in /books[author = $x/author],
      return <book>
        <title> { $y/title } </title>
        <publisher> { $y/publisher } </publisher>
        <keyword> { $y/keyword } </keyword>
      </book>
    } </author>
</bib>
```


- 23.14** Give the DTD for an XML representation of the information in Figure 7.29. Create a separate element type to represent each relationship, but use ID and IDREF to implement primary and foreign keys.

Answer:

```
<!DOCTYPE bookstore [
  <!ELEMENT basket (contains+, basket-of)>
  <!ATTLIST basket
    basketid ID #REQUIRED >
  <!ELEMENT customer (name, address, phone)>
  <!ATTLIST customer
    email ID #REQUIRED >
  <!ELEMENT book (year, title, price, written-by, published-by)>
  <!ATTLIST book
    ISBN ID #REQUIRED >
  <!ELEMENT warehouse (address, phone, stocks)>
  <!ATTLIST warehouse
    code ID #REQUIRED >
  <!ELEMENT author (name, address, URL)>
  <!ATTLIST author
    authid ID #REQUIRED >
  <!ELEMENT publisher (address, phone, URL)>
  <!ATTLIST publisher
    name ID #REQUIRED >
  <!ELEMENT basket-of >
  <!ATTLIST basket-of
    owner IDREF #REQUIRED >
  <!ELEMENT contains >
  <!ATTLIST contains
    book IDREF #REQUIRED
    number CDATA #REQUIRED >
  <!ELEMENT stocks >
  <!ATTLIST stocks
    book IDREF #REQUIRED
    number CDATA #REQUIRED >
  <!ELEMENT written-by >
  <!ATTLIST written-by
    authors IDREFS #REQUIRED >
  <!ELEMENT published-by >
  <!ATTLIST published-by
    publisher IDREF #REQUIRED >
  <!ELEMENT name (#PCDATA)>
  <!ELEMENT address (#PCDATA)>
  <!ELEMENT phone (#PCDATA)>
```

```

<!ELEMENT year (#PCDATA )>
<!ELEMENT title (#PCDATA )>
<!ELEMENT price (#PCDATA )>
<!ELEMENT number (#PCDATA )>
<!ELEMENT URL (#PCDATA )>
] >

```

23.15 Give an XML Schema representation of the DTD from Exercise 23.14.

Answer:

```

<xs:schema xmlns:xs="http://www.w3.org/2001/XMLSchema">
  <xs:element name="basket" type="BasketType">
    <xs:element name="customer">
      <xs:complexType>
        <xs:attribute name="email" use="required">
        <xs:sequence>
          <xs:element name="name" type="xs:string">
          <xs:element name="address" type="xs:string">
          <xs:element name="phone" type="xs:decimal">
        </xs:sequence>
      </xs:complexType>
    </xs:element>
    <xs:key name="customerKey">
      <xs:selector xpath="/bookstore/customer">
      <xs:field xpath="email">
    </xs:key>
    <xs:element name="book">
      <xs:complexType>
        <xs:attribute name="ISBN" use="required">
        <xs:sequence>
          <xs:element name="year" type="xs:decimal">
          <xs:element name="title" type="xs:string">
          <xs:element name="price" type="xs:decimal">
          <xs:element name="written-by" type="xs:string">
          <xs:element name="published-by" type="xs:string">
        </xs:sequence>
      </xs:complexType>
    </xs:element>
    <xs:key name="bookKey">
      <xs:selector xpath="/bookstore/book">
      <xs:field xpath="ISBN">
    </xs:key>
    <xs:element name="warehouse">
      <xs:complexType>

```

```

        <xs:attribute name="code" use="required">
        <xs:sequence>
        <xs:element name="address" type="xs:string">
        <xs:element name="phone" type="xs:decimal">
        <xs:element name="stocks" type="xs:decimal">
        </xs:sequence>
    </xs:complexType>
</xs:element>
<xs:key name="warehouseKey">
    <xs:selector xpath="/bookstore/warehouse">
    <xs:field xpath="code">
</xs:key>
<xs:element name="author">
    <xs:complexType>
        <xs:attribute name="authid" use="required">
        <xs:sequence>
        <xs:element name="name" type="xs:string">
        <xs:element name="address" type="xs:string">
        <xs:element name="URL" type="xs:string">
        </xs:sequence>
    </xs:complexType>
</xs:element>
<xs:key name="authorKey">
    <xs:selector xpath="/bookstore/author">
    <xs:field xpath="authid">
</xs:key>
<xs:element name="publisher">
    <xs:complexType>
        <xs:attribute name="name" use="required">
        <xs:sequence>
        <xs:element name="address" type="xs:string">
        <xs:element name="phone" type="xs:decimal">
        </xs:sequence>
    </xs:complexType>
</xs:element>
<xs:key name="publisherKey">
    <xs:selector xpath="/bookstore/publisher">
    <xs:field xpath="name">
</xs:key>
<xs:element name="basket-of">
    <xs:attribute name="owner" use="required">
</xs:element>
<xs:keyref name="basketCustomerFKey" refer="customerKey">
    <xs:selector xpath="/bookstore/customer">
    <xs:field xpath="owner">

```

```

</xs:key>
<xs:element name="contains">
  <xs:attribute name="book" use="required">
    <xs:attribute name="number" use="required" type="xs:decimal">
</xs:element>
<xs:keyref name="bookBasketFKey" refer="bookKey">
  <xs:selector xpath="/bookstore/book">
    <xs:field xpath="book/">
</xs:key>
<xs:element name="stocks">
  <xs:attribute name="book" use="required">
    <xs:attribute name="number" use="required" type="xs:decimal">
</xs:element>
<xs:keyref name="bookwarehouseFKey" refer="bookKey">
  <xs:selector xpath="/bookstore/book">
    <xs:field xpath="book/">
</xs:key>
<xs:element name="written-by">
  <xs:attribute name="authors" use="required">
</xs:element>
<xs:keyref name="authorbookFKey" refer="authorKey">
  <xs:selector xpath="/bookstore/author">
    <xs:field xpath="authors">
</xs:key>
<xs:element name="published-by">
  <xs:attribute name="publisher" use="required">
</xs:element>
<xs:keyref name="bookpublisherFKey" refer="publisherKey">
  <xs:selector xpath="/bookstore/publisher">
    <xs:field xpath="publisher">
</xs:key>
  <xs:complexType name="BasketType">
    <xs:attribute name="basketid" use="required">
    <xs:sequence>
      <xs:element name="contains" minOccurs="1"
        maxOccurs="unbounded">
      <xs:element name="basket-of">
    </xs:sequence>
  </xs:complexType>
<xs:key name="basketKey">
  <xs:selector xpath="/bookstore/basket">
    <xs:field xpath="basketid">
</xs:key>
</xs:schema>

```

23.16 Write queries in XQuery on the bibliography DTD fragment in Figure 23.16, to do the following:

- Find all authors who have authored a book and an article in the same year.
- Display books and articles sorted by year.
- Display books with more than one author.
- Find all books that contain the word “database” in their title and the word “Hank” in an author’s name (whether first or last).

Answer:

- Find all authors who have authored a book and an article in the same year.

```
for $a in distinct-values (/bib/book/author),
    $y in /bib/book[author=$a]/year,
    $art in /bib/article[author=$a and year=$y]
return $a
```

- Display books and articles sorted by year.

```
for $a in ((/bib/book) | (/bib/article))
order by $a/year
return $a
```

- Display books with more than one author.

```
for $b in ((/bib/book[author/count()>1])
return $b
```

- Find all books that contain the word “database” in their title and the word “Hank” in an author’s name (whether first or last).

```
for $a in /bibliography/book/
where $b in $a satisfies
    (contains($b/title,“database”) AND
    (contains($b/author(@last_name),“Hank”)
    OR contains($b/author(@first_name),“Hank”))
return $a
```

23.17 Give a relational mapping of the XML purchase order schema illustrated in Figure 23.3, using the approach described in Section 23.6.2.3. Suggest how to remove redundancy in the relational schema, if item identifiers functionally determine the description and purchase and supplier names functionally determine the purchase and supplier address, respectively.

Answer: Removing redundancy:

Item ids functionally determine the description The table 'item' can be broken down into 2 tables

```
item(item_id, quantity, price, itemlist_id)
and
item_info(item_id, description)
```

Then the item description doesn't get repeated whenever the item appears on the itemlist.

Purchaser_names functionally determine the address The table purchaser can be subdivided as shown below

```
purchaser(name, purchaseorder_id)
and
purchaser_info(name, address)
```

Supplier_names functionally determine the address The table supplier can be subdivided as shown below

```
supplier(name, purchaseorder_id)
and
supplier_info(name, address)
```

The address does not get repeated every time the supplier or purchaser is involved in a purchase.

```
create table purchase_order
(purchaseorder_id char(20),
purchaser_name char(20),
supplier_name char(20),
itemlist_id char(20),
total_cost numeric(16,2),
payment_terms char(20),
shipping_mode char(20),
primary key (purchaseorder_id),
foreign key (purchaser_name) references purchaser,
foreign key (supplier_name) references supplier,
foreign key (itemlist_id) references itemlist)
```

```
create table purchaser
  (purchaser_name char(20),
   address char(20),
   purchaseorder_id char(20),
   primary key (name),
   foreign key (purchaseorder_id) references purchase_order)
```

```
create table supplier
  (supplier_name char(20),
   address char(20),
   purchaseorder_id char(20),
   primary key (name),
   foreign key (purchaseorder_id) references purchase_order)
```

```
create table itemlist
  (itemlist_id char(20),
   item_id char(20),
   purchaseorder_id char(20),
   primary key (itemlist_id),
   foreign key (item_id) references item,
   foreign key (purchaseorder_id) references purchase_order)
```

```
create table item
  (item_id char(20),
   description char(20),
   quantity numeric(16),
   price numeric(16,2),
   itemlist_id char(20),
   primary key (item_id),
   foreign key (itemlist_id) references itemlist)
```

- 23.18** Write queries in SQL/XML to convert university data from the relational schema we have used in earlier chapters to the *university-1* and *university-2* XML schemas.

Answer: **university-1:**

```

select xmlelement(name "department",
  xmlelement(name "dept_name", d.dept_name),
  xmlelement(name "building" d.building),
  xmlelement(name "budget" d.budget),
  xmlelement(name "courses",
    (select xmlagg(xmlelement(name "course",
      xmlelement(name "course_id" c.course_id),
      xmlelement(name "title" c.title),
      xmlelement(name "credits" c.credits)))
    from course c
    where c.dept_name = d.dept_name
    order by c.course_id)))
from department d
select xmlelement(name "instructor",
  xmlelement(name "IID", i.IID),
  xmlelement(name "name", i.name),
  xmlelement(name "dept_name", i.dept_name),
  xmlelement(name "salary", i.salary),
  (select xmlagg(xmlelement(name "course_id", t.course_id))
    from teaches t
    where teaches.IID = i.IID
    order by t.course_id))
from instructor i

```


university-2:

```

select xmlelement(name "instructor",
  xmlelement(name "IID", i.IID),
  xmlelement(name "name", i.name),
  xmlelement(name "dept_name", i.dept_name),
  xmlelement(name "salary", i.salary),
  xmlelement(name "teaches",
    (select xmlagg(xmlelement(name "course",
      xmlelement(name "course_id" c.course_id),
      xmlelement(name "title" c.title),
      xmlelement(name "credits" c.credits)
    from course c, teaches t
    where c.course_id = t.course_id and t.IID = i.IID
    order by c.course_id))))))
from instructor i
select xmlelement(name "department",
  xmlelement(name "dept_name", d.dept_name),
  xmlelement(name "building" d.building),
  xmlelement(name "budget" d.budget),
  (select xmlagg(xmlelement(name "course_id", c.course_id))
    from course c
    where c.dept_name = d.dept_name)
from department d

```

23.19 As in Exercise 23.18, write queries to convert university data to the *university-1* and *university-2* XML schemas, but this time by writing XQuery queries on the default SQL/XML database to XML mapping.

Answer: **university-1:**

```
<university-1> {
  for $x in /university/department/row
  return
    <department>
      <dept_name> {$x/dept_name} </dept_name>
      <building> {$x/building} </building>
      <budget> {$x/budget} </budget>
      for $c in /university/course/row[dept_name=$x/dept_name]
      return <course>
        <course_id> {$c/course_id} </course_id>
        <title> {$c/title} </title>
        <credits> {$c/credits} </credits>
      </course>
    </department>
  for $i in /university/instructor/row
  return
    <instructor>
      <IID> {$i/IID} </IID>
      <name> {$i/name} </name>
      <dept_name> {$i/dept_name} </dept_name>
      <salary> {$i/salary} </salary>
      for $t in /university/teaches/row[IID=$i/IID]
      return {
        <course_id> {$t/course_id} </course_id>
      }
    </instructor>
} </university-1>
```

university-2:

```

<university-2> {
  for $x in /university/instructor/row/
  return
    <instructor>
      <IID> {$i/IID} </IID>
      <name> {$i/name} </name>
      <dept_name> {$i/dept_name} </dept_name>
      <salary> {$i/salary} </salary>
      <teaches> {
        for $t in /university/teaches/row[IID=$i/IID]
        $c in /university/course/row[$t.course_id=$c/course_id]
        return <course>
          <course_id> {$c/course_id} </course_id>
          <title> {$c/title} </title>
          <dept_name> {$c/dept_name} </dept_name>
          <credits> {$c/credits} </credits>
        </course>
      } </teaches>
    </instructor>
  for $x in /university/department/row/
  return
    <department>
      <dept_name> {$x/dept_name} </dept_name>
      <building> {$x/building} </building>
      <budget> {$x/budget} </budget>
      for $c in /university/course/row[dept_name=$x/dept_name]
      return
        <course_id> {$c/course_id} </course_id>
    </department>
} </university-2>

```

- 23.20** One way to shred an XML document is to use XQuery to convert the schema to an SQL/XML mapping of the corresponding relational schema, and then use the SQL/XML mapping in the backward direction to populate the relation.

As an illustration, give an XQuery query to convert data from the *university-1* XML schema to the SQL/XML schema shown in Figure 23.15.

Answer:

```

<department> {
  for $x in /university-1/department
  return
    <row>
      <dept_name> { $x/dept_name } </dept_name>
      <building> { $x/building } </building>
      <budget> { $x/budget } </budget>
    </row>
} </department>

<course> {
  for $x in /university-1/department/course
  return
    <row> { $x } </row>
} </course>

<instructor> {
  for $x in /university-1/instructor
  return
    <row>
      <IID> { $x/IID } </IID>
      <name> { $x/name } </name>
      <dept_name> { $x/dept_name } </dept_name>
      <salary> { $x/salary } </salary>
    </row>
} </instructor>

<teaches> {
  for $i in /university-1/instructor, $cid in $i/course_id
  return
    <row>
      <IID> { $i/IID } </IID>
      <course_id> { $cid } </course_id>
    </row>
} </teaches>

```

- 23.21** Consider the example XML schema from Section 23.3.2, and write XQuery queries to carry out the following tasks:

- a. Check if the key constraint shown in Section 23.3.2 holds.
- b. Check if the keyref constraint shown in Section 23.3.2 holds.

Answer:

- a. Check if the key constraint shown in Section 23.3.2 holds.

```

let $x = /university/department
let $y = distinct-values(/university/department/dept_name)
return {
  if fn:count($y) == fn:count($x)
    then 1
    else 0
}

```

- b. Check if the keyref constraint shown in Section 23.3.2 holds.

```

return
  (every $c in /university/course satisfies
    (some $d in /university/department satisfies
      $c/dept_name = $d/dept_name))

```

23.22 Consider Practice Exercise 23.7, and suppose that authors could also appear as top-level elements. What change would have to be done to the relational schema?

Answer: Author id can be added as an attribute to the author top element. It can be used to refer to the author in the book and article topelements and also keep track of the order.

CHAPTER 24



Advanced Application Development

Exercises

- 24.6 Find out all performance information your favorite database system provides. Look for at least the following: what queries are currently executing or executed recently, what resources each of them consumed (CPU and I/O), what fraction of page requests resulted in buffer misses (for each query, if available), and what locks have a high degree of contention. You may also be able to get information about CPU and I/O utilization from the operating system.

Answer:

- PostgreSQL: The *EXPLAIN* command lets us see what query plan the system creates for any query. The numbers that are currently quoted by *EXPLAIN* are:
 - a. Estimated start-up cost
 - b. Estimated total cost
 - c. Estimated number of rows output by this plan node
 - d. Estimated average width of rows
- SQL: There is a Microsoft tool called *SQL Profiler*. The data is logged, and then the performance can be monitored. The following performance counters can be logged.
 - a. Memory
 - b. Physical Disk
 - c. Process
 - d. Processor
 - e. SQLServer:Access Methods
 - f. SQLServer:Buffer Manager
 - g. SQLServer:Cache Manager

- h. SQLServer:Databases
- i. SQLServer:General Statistics
- j. SQLServer:Latches
- k. SQLServer:Locks
- l. SQLServer:Memory Manager
- m. SQLServer:SQL Statistics
- n. SQLServer:SQL Settable

24.7 a. What are the three broad levels at which a database system can be tuned to improve performance?

b. Give two examples of how tuning can be done for each of the levels.

Answer:

a. We refer to performance tuning of a database system as the modification of some system components in order to improve transaction response times, or overall transaction throughput. Database systems can be tuned at various levels to enhance performance. viz.

- i. Schema and transaction design
- ii. Buffer manager and transaction manager
- iii. Access and storage structures
- iv. Hardware - disks, CPU, busses etc.

b. We describe some examples for performance tuning of some of the major components of the database system.

i. **Tuning the schema**

In this chapter we have seen two examples of schema tuning, viz. vertical partition of a relation (or conversely - join of two relations), and denormalization (or conversely - normalization). These examples reflect the general scenario, and ideas therein can be applied to tune other schemas.

ii. **Tuning the transactions**

One approach used to speed-up query execution is to improve the its plan. Suppose that we need the natural join of two relations - say *account* and *depositor* from our sample bank database. A *sort-merge-join* (Section 12.5.4) on the attribute *account-number* may be quicker than a simple nested-loop join on the relations.

Other ways of tuning transactions are - breaking up long update transactions and combining related sets of queries into a single query. Generic examples for these approaches are given in this chapter.

For client-server systems, wherein the query has to be transmitted from client to server, the query transmission time itself may form a large fraction of the total query cost. Using *stored procedures* can significantly reduce the queries response time.

iii. **Tuning the buffer manager**

The buffer manager can be made to increase or decrease the number of pages in the buffer according to changing page-fault rates. However, it must be noted that a larger number of pages may mean higher costs for latch management and maintenance of other data-structures like free-lists and page map tables.

iv. **Tuning the transaction manager**

The transaction schedule affects system performance. A query that computes statistics for customers at each branch of the bank will need to scan the relations *account* and *depositor*. During these scans, no updates to any customer's balance will be allowed. Thus, the response time for the update transactions is high. Large queries are best executed when there are few updates, such as at night.

Checkpointing also incurs some cost. If recovery time is not critical, it is preferable to examine a long log (during recovery) rather than spend a lot of (checkpointing) time during normal operation. Hence it may be worthwhile to tune the checkpointing interval according to the expected rate of crashes and the required recovery time.

v. **Tuning the access and storage structures**

A query's response time can be improved by creating an appropriate index on the relation. For example, consider a query in which a depositor enquires about her balance in a particular account. This query would result in the scan of the relation *account* if it has no index on *account-number*. Similar indexing considerations also apply to computing joins. i.e an index on *account-number* in the *account* relation saves scanning *account* when a natural join of *account* is taken with *depositor*.

In contrast, performance of update transactions may suffer due to indexing. Let us assume that frequent updates to the balance are required. Also suppose that there is an index on *balance* (presumably for range queries) in *account*. Now, for each update to the value of the balance, the index too will have to be updated. In addition, concurrent updates to the index structure will require additional locking overheads. Note that the response time for each update would not be more if there were no index on *balance*.

The type of index chosen also affects performance. For a range query, an order preserving index (like B-trees) is better than a hashed index.

Clustering of data affects the response time for some queries. For example, assume that the tuples of the *account* relation are clustered on *branch-name*. Then the average execution time for a query that finds the total balance amount deposited at a partic-

ular branch can be improved. Even more benefit accrues from having a clustered index on *branch-name*.

If the database system has more than one disk, *declustering* of data will enable parallel access. Suppose that we have five disks and that in a hypothetical situation where each customer has five accounts and each account has a lot of historical information that needs to be accessed. Storing one account per customer per disk will enable parallel access to all accounts of a particular customer. Thus, the speed of a scan on *depositor* will increase about five-fold.

vi. **Tuning the hardware**

The hardware for the database system typically consists of disks, the processor, and the interconnecting architecture (busses etc.). Each of these components may be a bottleneck and by increasing the number of disks or their block-sizes, or using a faster processor, or by improving the bus architecture, one may obtain an improvement in system performance.

- 24.8** When carrying out performance tuning, should you try to tune your hardware (by adding disks or memory) first, or should you try to tune your transactions (by adding indices or materialized views) first. Explain your answer.

Answer: The 3 levels of tuning - hardware, database system parameters, schema and transactions - interact with one another; so we must consider them together when tuning a system. For example, tuning at the transaction level may result in the hardware bottleneck changing from the disk system to the CPU, or vice versa.

- 24.9** Suppose that your application has transactions that each access and update a single tuple in a very large relation stored in a B⁺-tree file organization. Assume that all internal nodes of the B⁺-tree are in memory, but only a very small fraction of the leaf pages can fit in memory. Explain how to calculate the minimum number of disks required to support a workload of 1000 transactions per second. Also calculate the required number of disks, using values for disk parameters given in Section 10.2.

Answer: Given that all internal nodes of the B⁺-tree are in memory, and only a very small fraction of the leaf pages can fit in memory. We can deduce that each I/O transaction that access and update a single tuple requires just 1 I/O operation. The disk with the parameters given in the chapter would support a little under 100 random-access I/O operations of 4 kilbytes each per second. So, number of disks needed to support a workload of 1000 transactions is $1000/100 = 10$ disks.

The disk parameters given in Section 10.2 are almost the same as the values in the current chapter. So, the number of disks required will be around 10 in this case also.

- 24.10** What is the motivation for splitting a long transaction into a series of small ones? What problems could arise as a result, and how can these problems be averted?

Answer: Long update transactions cause a lot of log information to be written, and hence extend the checkpointing interval and also the recovery time after a crash. A transaction that performs many updates may even cause the system log to overflow before the transaction commits.

To avoid these problems with a long update transaction it may be advisable to break it up into smaller transactions. This can be seen as a *group* transaction being split into many small *mini-batch* transactions. The same effect is obtained by executing both the group transaction and the mini-batch transactions, which are scheduled in the order that their operations appear in the group transaction.

However, executing the mini-batch transactions in place of the group transaction has some costs, such as extra effort when recovering from system failures. Also, even if the group transaction satisfies the *isolation* requirement, the mini-batch may not. Thus the transaction manager can release the locks held by the mini-batch only when the last transaction in the mini-batch completes execution.

- 24.11** Suppose the price of memory falls by half, and the speed of disk access (number of accesses per second) doubles, while all other factors remain the same. What would be the effect of this change on the 5 minute and 1 minute rule?

Answer: There will be no effect of these changes on the 5 minute or the 1 minute rule. The value of n , i.e. the frequency of page access at the break-even point, is proportional to the product of memory price and speed of disk access, other factors remaining constant. So when memory price falls by half and access speed doubles, n remains the same.

- 24.12** List at least 4 features of the TPC benchmarks that help make them realistic and dependable measures.

Answer: Some features that make the TPC benchmarks realistic and dependable are -

- Ensuring full support for ACID properties of transactions,
- Calculating the throughput by observing the *end-to-end* performance,
- Making sizes of relations proportional to the expected rate of transaction arrival, and
- Measuring the dollar cost per unit of throughput.

- 24.13** Why was the TPC-D benchmark replaced by the TPC-H and TPC-R benchmarks?

Answer: Various TPC-D queries can be significantly speeded up by using materialized views and other redundant information, but the overheads of using them should be properly accounted. Hence TPC-R and TPC-H

were introduced as refinements of TPC-D, both of which use same schema and workload. TPC-R models periodic reporting queries, and the database running it is permitted to use materialized views. TPC-H, on the other hand, models ad hoc querying, and prohibits materialized views and other redundant information.

- 24.14** Explain what application characteristics would help you decide which of TPCC, TPC-H, or TPC-R best models the application.

Answer: Depending on the application characteristics, different benchmarks are used to model it.

The TPCC benchmark is widely used for transaction processing. It is appropriate for applications which concentrate on the main activities in an order-entry environment, such as entering and delivering orders, recording payments, checking status of orders, and monitoring levels of stock.

The TPC-H (H represents *ad hoc*) benchmark models the applications which prohibit materialized views and other redundant information, and permits indices only on primary and foreign keys. This benchmark models ad-hoc querying where the queries are not known beforehand.

The TPC-R (R represents for *reporting*) models the applications which has queries, inserts, updates, and deletes. The application is permitted to use materialized views and other redundant information.

CHAPTER 25



Advanced Data Types and New Applications

This chapter covers advanced data types and new applications, including temporal databases, spatial and geographic databases, multimedia databases, and mobility and personal databases. In particular, the data types mentioned above have grown in importance in recent years, and commercial database systems are increasingly providing support for such data types through extensions to the database system variously called cartridges or extenders.

This chapter is suited as a means to lay the groundwork for an advanced course. Some of the material, such as temporal and spatial data types, may be suitable for self-study in a first course.

Exercises

- 25.9 Will functional dependencies be preserved if a relation is converted to a temporal relation by adding a time attribute? How is the problem handled in a temporal database?

Answer: Functional dependencies may be violated when a relation is augmented to include a time attribute. For example, suppose we add a time attribute to the relation *account* in our sample bank database. The dependency *account-number* \rightarrow *balance* may be violated since a customer's balance would keep changing with time.

To remedy this problem temporal database systems have a slightly different notion of functional dependency, called *temporal functional dependency*. For example, the temporal functional dependency:

$$account\text{-}number \xrightarrow{\tau} balance$$

over *Account-schema* means that for each instance *account* of *Account-schema*, all snapshots of *account* satisfy the functional dependency

account-number→*balance*;

i.e at any time instance, each account will have a unique bank balance corresponding to it.

- 25.10** Consider two-dimensional vector data where the data items do not overlap. Is it possible to convert such vector data to raster data? If so, what are the drawbacks of storing raster data obtained by such conversion, instead of the original vector data?

Answer: To convert non-overlapping vector data to raster data, we set the values for exactly those pixels that lie on any one of the data items (regions); the other pixels have a default value.

The disadvantages to this approach are: loss of precision in location information (since raster data loses resolution), a much higher storage requirement, and loss of abstract information (like the shape of a region).

- 25.11** Study the support for spatial data offered by the database system that you use, and implement the following:
- A schema to represent the geographic location of restaurants along with features such as the cuisine served at the restaurant and the level of expensiveness.
 - A query to find moderately priced restaurants that serve Indian food and are within 5 miles of your house (assume any location for your house).
 - A query to find for each restaurant the distance from the nearest restaurant serving the same cuisine and with the same level of expensiveness.

Answer: PostgreSQL includes support for R-tree indices over spatial data, as well as a number of built-in geometric data types (points, boxes, circles, lines, and paths) to represent spatial data, and functions to manipulate this data.

a.

```
create table restaurants (
    name varchar(30),
    location point,
    cuisine varchar(30),
    price int)
```

- b. Assume your house is at coordinates (21.5, 14.2), and that a price value of 2 means “moderately priced”.
 <-> is the PostgreSQL operator representing “distance between”.

```
select name
```

```

from restaurants
where ((point '(21.5, 14.2)) <-> location) < 5.0
      and cuisine = 'Indian'
      and price <= 2

```

c.

```

select r1.name, min(r1.location <-> r2.location)
from restaurants as r1, restaurants as r2
where r1.cuisine = r2.cuisine
      and r1.price = r2.price
group by r1.name

```

- 25.12 What problems can occur in a continuous-media system if data are delivered either too slowly or too fast?

Answer: Continuous media systems typically handle a large amount of data, which have to be delivered at a steady rate. Suppose the system provides the picture frames for a television set. The delivery rate of data from the system should be matched with the frame display rate of the TV set. If the delivery rate is too low, the display would periodically freeze or blank out, since there will be no new data to be displayed for some time. On the other hand, if the delivery rate is too high, the data buffer at the destination TV set will overflow causing loss of data; the lost data will never get displayed.

- 25.13 List three main features of mobile computing over wireless networks that are distinct from traditional distributed systems.

Answer: Some of the main distinguishing features are as follows.

- In distributed systems, disconnection of a host from the network is considered to be a *failure*, whereas allowing such disconnection is a *feature* of mobile systems.
- Distributed systems are usually centrally administered, whereas in mobile computing, each personal computer that participates in the system is administered by the user (owner) of the machine and there is little central administration, if any.
- In conventional distributed systems, each machine has a fixed location and network address(es). This is not true for mobile computers, and in fact, is antithetical to the very purpose of mobile computing.
- Queries made on a mobile computing system may involve the location and velocity of a host computer.
- Each computer in a distributed system is allowed to be arbitrarily large and may consume a lot of (almost) uninterrupted electrical power. Mobile systems typically have small computers that run on low wattage, short-lived batteries.

- 25.14** List three factors that need to be considered in query optimization for mobile computing that are not considered in traditional query optimizers.

Answer: The most important factor influencing the cost of query processing in traditional database systems is that of disk I/O. However, in mobile computing, minimizing the amount of energy required to execute a query is an important task of a query optimizer. To reduce the consumption of energy (battery power), the query optimizer on a mobile computer must minimize the size and number of queries to be transmitted to remote computers as well as the time for which the disk is spinning.

In traditional database systems, the cost model typically does not include connection time and the amount of data transferred. However, mobile computer users are usually charged according to these parameters. Thus, these parameters should also be minimized by a mobile computer's query optimizer.

- 25.15** Give an example to show that the version-vector scheme does not ensure serializability. (Hint: Use the example from Practice Exercise 25.8, with the assumption that documents 1 and 2 are available on both mobile computers A and B, and take into account the possibility that a document may be read without being updated.)

Answer: Consider the example given in the previous exercise. Suppose that both host A and host B are not connected to each other. Further, assume that identical copies of document 1 and document 2 are stored at host A and host B.

Let $\{X = 5\}$ be the initial contents of document 1, and $\{X = 10\}$ be the initial contents of document 2. Without loss of generality, let us assume that all version-vectors are initially zero.

Suppose host A updates the number its copy of document 1 with that in its copy of document 2. Thus, the contents of both the documents (at host A) are now $\{X = 10\}$. The version number $V_{1,A,A}$ is incremented to 1.

While host B is disconnected from host A, it updates the number in its copy of document 2 with that in its copy of document 1. Thus, the contents of both the documents (at host B) are now $\{X = 5\}$. The version number $V_{2,B,B}$ is incremented to 1.

Later, when host A and host B connect, they exchange version-vectors. The version-vector scheme updates the copy of document 1 at host B to $\{X = 10\}$, and the copy of document 2 at host A to $\{X = 5\}$. Thus, both copies of each document are identical, viz. document 1 contains $\{X = 10\}$ and document 2 contains $\{X = 5\}$.

However, note that a serial schedule for the two updates (one at host A and another at host B) would result in both documents having the *same* contents. Hence this example shows that the version-vector scheme does not ensure serializability.

CHAPTER 26



Advanced Transaction Processing

In this chapter, we go beyond the basic transaction processing schemes discussed previously, and cover more advanced transaction-processing concepts, including transaction-processing monitors, workflow systems, main-memory databases, real-time transaction systems, and handling of long-duration transactions by means of nested transactions, multi-level transactions and weak degrees of consistency. We end the chapter by covering weak degrees of consistency used to handle multidatabase systems.

This chapter is suited to an advanced course. The sections on TP monitors and workflows may also be covered in an introductory course as independent-study material.

Coverage of remote backup systems has been moved from this chapter to the chapter on recovery, while coverage of transaction processing in multidatabases has been moved into this chapter from its earlier position in the distributed database chapter.

Exercises

- 26.7 Explain how a TP monitor manages memory and processor resources more effectively than a typical operating system.

Answer: In a typical OS, each client is represented by a process, which occupies a lot of memory. Also process multi-tasking over-head is high. A TP monitor is more of a service provider, rather than an environment for executing client processes. The client processes run at their own sites, and they send requests to the TP monitor whenever they wish to avail of some service. The message is routed to the right server by the TP monitor, and the results of the service are sent back to the client.

The advantage of this scheme is that the same server process can be serving several clients simultaneously, by using multithreading. This saves memory space, and reduces CPU overheads on preserving ACID prop-

erties and on scheduling entire processes. Even without multi-threading, the TP monitor can dynamically change the number of servers running, depending on whatever factors affect good performance. All this is not possible with a typical OS setup.

- 26.8** Compare TP-monitor features with those provided by Web servers supporting servlets (such servers have been nicknamed *TP-lite*).

Answer: Web application servers supporting servlets (as also Web application servers providing similar functionality through other language APIs) have many features of TP monitors. For example, they allow a single process, or a few processes, to serve a large number of requests by exploiting multi-threading. Systems built to handle large numbers of requests usually have routers that divide up incoming traffic between a large number of Web application servers. Web application servers typically lack a few features such as support for transaction coordination (using two-phase commit), which TP-monitors support, although some application servers do have add-on features that support transaction coordination.

- 26.9** Consider the process of admitting new students at your university (or new employees at your organization).

- a. Give a high-level picture of the workflow starting from the student application procedure.
- b. Indicate acceptable termination states and which steps involve human intervention.
- c. Indicate possible errors (including deadline expiry) and how they are dealt with.
- d. Study how much of the workflow has been automated at your university.

Answer:

- a. Students typically apply online to the university. Once they have filled in all details of their form, they submit it, and pay application fees. The fee payment is itself a sub-workflow, involving an external party such as a credit card company or bank, and possibly a third-party payment gateway as an intermediary to the credit card company or bank.

At this point the university has to process the form. The first step of processing by the university is to check that the applications are complete; if anything is incomplete, the student has to be informed and given a chance to correct the form. Typically there is a date for starting the decision process; all forms received up to that date are taken together, and after filtering on some criteria, and inputs from various people involved in the admission process, an admission decision is made. Subsequently forms may be placed in accepted, rejected, or wait-listed state.

This decision is conveyed back to students. The notification is done by email and SMS, and the deliver of email/SMS can itself be considered as an automated workflow since it involves external services, with persistent message delivery queues.

- b. Acceptable termination states include: student admitted and student rejected; intermediate states may include student wait-listed, which will eventually end up in either admitted or reject state. The eventual decision on admission or rejection (or temporary wait-listing) will require human decisions. Applications may be rejected early if some requirements are not met, such as application fees not paid, without human intervention.
- c. At IIT Bombay, as in most institutions today, applications, including fee payment are online and automated. Humans are involved in the decision making process, but there is a fair degree of bulk processing, for example cutoff marks are specified for standardized exams, and all students below the cutoff are automatically rejected. Once the decision is made, the notification by email/SMS is itself a workflow, which is automated.

26.10 Answer the following questions regarding electronic payment systems:

- a. Explain why electronic transactions carried out using credit-card numbers may be insecure.
- b. An alternative is to have an electronic payment gateway maintained by the credit-card company, and the site receiving payment redirects customers to the gateway site to make the payment.
 - i. Explain what benefits such a system offers if the gateway does not authenticate the user.
 - ii. Explain what further benefits are offered if the gateway has a mechanism to authenticate the user.
- c. Some credit-card companies offer a one-time-use credit-card number as a more secure method of electronic payment. Customers connect to the credit-card company's Web site to get the one-time-use number. Explain what benefit such a system offers, as compared to using regular credit-card numbers. Also explain its benefits and drawbacks as compared to electronic payment gateways with authentication.
- d. Does either of the above systems guarantee the same privacy that is available when payments are made in cash? Explain your answer.

Answer:

- a. Credit card numbers can be easily stolen by someone who handles the card (for offline transactions) or may be stolen in bulk from computers where they are stored in (online) shops where a purchase was made. Some who has got the number (including verification

numbers such as the widely used CVV number) can make purchases online using these numbers.

- b.
 - i. Even if the gateway does not authorize the user, as long as the user can check the authenticity of the gateway (using HTTPS protocol certificates), the user knows that the credit card numbers are only handled by a trusted party; the site receiving the payment never gets to even see the card number.
 - ii. If the gateway further authenticates the user, using additional passwords (usually this is done by further redirecting the user to the bank which issued the credit card, so the user need not trust the gateway with passwords), even someone who has access to the credit card number cannot make online purchases using the credit card. Such authentication is mandatory for online payments in some countries.
- c. One-time-use credit card numbers cannot be used again, so even if it is stored at the site receiving the payment, and subsequently compromised, it cannot be used to make any further purchases online. Thus, it has similar benefits to a system using authentication. However, its drawback is that the user has to perform extra actions to get a one-time-use number.
- d. With cash, it is quite possible to make completely anonymous purchases, with complete privacy. None of the above systems guarantees the same privacy as cash. First, the credit card company knows who you bought things from, so there is not question of a fully anonymous purchase. Second, even though your identity can be concealed from the site where you make a purchase, law enforcement officials can get information from both the site and the credit card company, compromising your privacy completely.

26.11 If the entire database fits in main memory, do we still need a database system to manage the data? Explain your answer.

Answer: Even if the entire database fits in main memory, a DBMS is needed to perform tasks like concurrency control, recovery, logging etc, in order to preserve ACID properties of transactions.

26.12 In the group-commit technique, how many transactions should be part of a group? Explain your answer.

Answer: As log-records are written to stable storage in multiples of a block, we should group transaction commits in such a way that the last block containing log-records for the current group is almost full.

26.13 In a database system using write-ahead logging, what is the worst-case number of disk accesses required to read a data item from a specified disk page. Explain why this presents a problem to designers of real-time database systems. Hint: consider the case when the disk buffer is full.

Answer: In the worst case, a read can cause a buffer page to be written to disk (preceded by the corresponding log records), followed by the reading from disk of the page containing the data to be accessed. This takes two or more disk accesses, and the time taken is several orders of magnitude more than the main-memory reference required in the best case. Hence transaction execution-time variance is very high and can be estimated only poorly. It is therefore difficult to plan schedules which need to finish within a deadline.

- 26.14** What is the purpose of compensating transactions? Present two examples of their use.

Answer: A compensating transaction is used to perform a semantic undo of changes made previously by committed transactions. For example, a person might deposit a check in their savings account. Then the database would be updated to reflect the new balance. Since it takes a few days for the check to clear, it might be discovered later that the check bounced, in which case a compensating transaction would be run to subtract the amount of the bounced check from the depositor's account. Another example of when a compensating transaction would be used is in a grading program. If a student's grade on an assignment is to be changed after it is recorded, a compensating program (usually an option of the grading program itself) is run to change the grade and redo averages, etc.

- 26.15** Explain the connections between a workflow and a long-duration transaction.

Answer: A long-duration transaction is still expected to finish within the time of a user interaction, whereas a workflow may last for a very long time. Long duration transactions can be aborted at any point, whereas a workflow cannot simply be rolled back, instead it has to be moved to an acceptable termination state.

