



Министерство науки и высшего образования Российской Федерации
Федеральное государственное бюджетное образовательное учреждение
высшего образования
«Московский государственный технический университет
имени Н.Э. Баумана
(национальный исследовательский университет)»
(МГТУ им. Н.Э. Баумана)

ФАКУЛЬТЕТ «Информатика и системы управления»

КАФЕДРА «Программное обеспечение ЭВМ и информационные технологии»

РАСЧЕТНО-ПОЯСНИТЕЛЬНАЯ ЗАПИСКА К ВЫПУСКНОЙ КВАЛИФИКАЦИОННОЙ РАБОТЕ НА ТЕМУ:

*«Метод параллельного выполнения запросов к СУБД
PostgreSQL в пределах одного соединения»*

Студент группы ИУ7-75Б

(Подпись, дата)

О.С. Платонова
(И.О.Фамилия)

Руководитель ВКР

(Подпись, дата)

М.В. Филиппов
(И.О.Фамилия)

Консультант ВКР

(Подпись, дата)

Ю.М. Гаврилова
(И.О.Фамилия)

Нормоконтролер

(Подпись, дата)

(И.О.Фамилия)

2021 г.

РЕФЕРАТ

Расчетно-пояснительная записка 15 страниц, 5 рисунков, 15 источников.

БАЗА ДАННЫХ, POSTGRESQL, МНОГОПОТОЧНЫЕ СУБД

СОДЕРЖАНИЕ

Введение	5
1 Аналитический раздел	7
1.1 Анализ СУБД	7
1.2 Архитектура PostgreSQL	8
1.3 Соединение в PostgreSQL.....	9
1.4 Многопоточность	9
1.5 Пул соединений	11
1.5.1 Frontend pool.....	11
1.5.2 Server pool.....	12
Список использованных источников	14

ОПРЕДЕЛЕНИЯ, ОБОЗНАЧЕНИЯ И СОКРАЩЕНИЯ

База данных (БД) — собрание данных, организованных в соответствии с концептуальной структурой, описывающей характеристики этих данных и взаимоотношения между ними, причем такое собрание данных, которое поддерживает одну или более областей применения. [1]

Система управления базой данных (СУБД) — совокупность программных и лингвистических средств общего или специального назначения, обеспечивающих управление созданием и использованием баз данных. [2]

Массивно-параллельная архитектура (massive parallel processing, MPP) — класс архитектур параллельных вычислительных систем. Главная особенность такой архитектуры состоит в том, что память физически разделена. [3]

Введение

В XXI веке человечество владеет невообразимым объемом данных. Знания, передававшиеся из поколения в поколение в течение многих тысячелетий, продолжают увеличиваться каждый день. Так, ежегодный прирост информации составляет 30%. [4]

С появлением письменности, будь то шумерские таблички или берестяные грамоты, перед человечеством возникает вопрос хранения и обработки данных. Причем с развитием цивилизации, и, как следствие, увеличением документооборота, проблема хранения информации требует систематического решения. Например, в конце XX века данные крупной компании могли занимать несколько этажей, что требовало дополнительных кадров для работы с ними.

Первым этапом решения этого вопроса стало внедрение компьютеров. Многие операции с данными были упрощены, а быстрый рост информационных технологий привел к увеличению скорости работы над данными. Однако хранение информации в виде файлов на одном компьютере стало неэффективным. Во-первых, поиск файла в файловой системе был долгим. Во-вторых, хранение информации в одном файле затрудняло поиск необходимых данных.

Решение проблемы разрозненного хранения данных впервые было представлено на симпозиуме в 1963 году в Санта-Монике. Хотя речь шла о внедрении баз данных в военные приложения, этот момент считается точкой отсчета истории базы данных. Их применение в работе компаний привело к увеличению скорости работы. А автоматизация основных процессов базы данных, таких как создание, просмотр, удаление данных привело к созданию системы управления базы данных.

В 2021 году ни одна сфера жизни не обходится без компьютеризации. Организации используют базы и СУБД для перевода данных в электронный вид. Необходимость перевода заключается не столько в потребности сократить временные и материальные (сокращение кадров) расходы, сколько в

поддержании конкурентоспособности. Переход компании в электронный вид дает возможность приобретения принципиально новых качеств, позволяющих иметь существенные преимущества над другими.

Из-за высокой популярности СУБД возникает вопрос об оптимизации ее работы. Так как один из самых распространенных способов увеличения производительности — параллельное выполнение, следует рассмотреть оптимизацию многопоточной программы. Поскольку операция соединения с базой данных является одной из самых дорогостоящих, следует минимизировать количество соединений.

Целью данной работы является разработка и реализация метода параллельного выполнения запросов к СУБД PostgreSQL в пределах одного соединения. Для достижения поставленной цели необходимо решить следующие задачи:

1. анализ предметной области и существующих методов реализации многопоточного доступа в MPP системах;
2. разработка метода параллельного выполнения запросов к СУБД PostgreSQL в пределах одного соединения;
3. реализация программного модуля для СУБД PostgreSQL;
4. проведение сравнительного анализа стандартного метода обработки запросов к СУБД PostgreSQL с реализуемым методом.

1 Аналитический раздел

В данном разделе будет выполнен анализ существующих СУБД, представлены разработанные методы и алгоритмы решения поставленной задачи. Также будет выполнен сравнительный анализ решений с указанием достоинств и недостатков.

1.1 Анализ СУБД

На рисунке 1.1 представлен рейтинг популярности СУБД, составленный компанией «DB-Engines» по состоянию на конец 2021 года. [5]

Rank			DBMS	Database Model	Score		
Nov 2021	Oct 2021	Nov 2020			Nov 2021	Oct 2021	Nov 2020
1.	1.	1.	Oracle +	Relational, Multi-model ⓘ	1272.73	+2.38	-72.27
2.	2.	2.	MySQL +	Relational, Multi-model ⓘ	1211.52	-8.25	-30.12
3.	3.	3.	Microsoft SQL Server +	Relational, Multi-model ⓘ	954.29	-16.32	-83.35
4.	4.	4.	PostgreSQL + ⓘ	Relational, Multi-model ⓘ	597.27	+10.30	+42.22
5.	5.	5.	MongoDB +	Document, Multi-model ⓘ	487.35	-6.21	+33.52
6.	6.	↑ 7.	Redis +	Key-value, Multi-model ⓘ	171.50	+0.15	+16.08
7.	7.	↓ 6.	IBM Db2	Relational, Multi-model ⓘ	167.52	+1.56	+5.90
8.	8.	8.	Elasticsearch	Search engine, Multi-model ⓘ	159.09	+0.84	+7.54
9.	9.	9.	SQLite +	Relational	129.80	+0.43	+6.48
10.	10.	10.	Cassandra +	Wide column	120.88	+1.61	+2.13

Рисунок 1.1 — Рейтинг популярности СУБД.

Согласно рейтингу, лидирующие позиции занимают реляционные модели баз данных. Данная работа будет основываться на объектно-реляционной СУБД PostgreSQL 12-ой версии, занимающей 4-ую строчку. Выбор аргументирован следующими преимуществами:

- доступность исходного кода;
- кроссплатформенность;
- быстроедействие;
- наследуемость.

1.2 Архитектура PostgreSQL

Одной из наиболее сильных сторон PostgreSQL является архитектура, основанная на модели «клиент-сервер». Выделяют 3 основные подсистемы: клиентская часть, серверная часть и хранилище данных (рисунок 1.2.). [6]

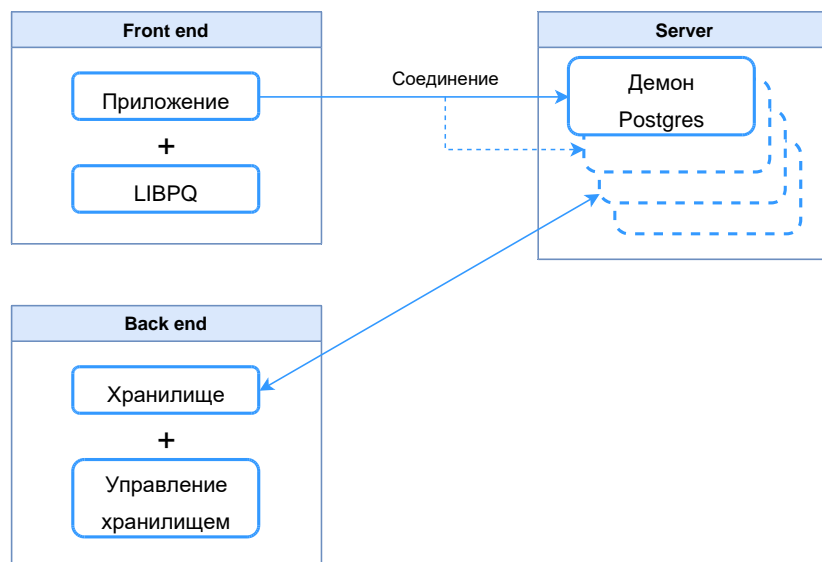


Рисунок 1.2 — «Клиент-серверная» архитектура PostgreSQL.

Клиентская часть состоит из пользовательского приложения и библиотеки LIBPQ, реализующей интерфейс взаимодействия с сервером. [7]

Серверная часть, включающая в себя процесс-демон Postgres и серверные процессы, выполняет обработку запросов. Соединение, установленное клиентом, принимается процессом Postgres. Процесс-демон в дальнейшем с помощью системного вызова *fork()* создаст новый серверный процесс для обслуживания соединения данного клиента. После установки соединения и обработки запроса, результаты будут возвращены обратно клиенту через установленное соединение. [8]

Серверные процессы взаимодействуют между собой через семафоры и разделяемую память, чтобы обеспечить целостность данных при одновременном обращении к ним. [7]

Третья часть сформирована из хранилища данных и средств его управления. Допускается обращение нескольких серверных процессов к информации хранилища одновременно.

Представленная концепция взаимодействия позволяет исключать произвольный доступ клиентов к данным, тем самым поддерживая их целостность и повышая уровень безопасности. Также она дает возможность создания распределенных систем.

1.3 Соединение в PostgreSQL

Как было рассмотрено выше, пользователь устанавливает соединение и посылает его серверному процессу Postgres. Каждое соединение представляется объектом PGconn, который можно получить от функций PQconnectdb, PQconnectdbParams или PQsetdbLogin.

PostgreSQL содержит инструменты для реализации многопоточности. Один из них — библиотека libpq, которая по умолчанию поддерживает повторные вызовы. Однако при реализации многопоточности существует ограничение: «два потока не должны пытаться одновременно работать с одним объектом PGconn. В частности, не допускается параллельное выполнение команд из разных потоков через один объект соединения.». [7]

1.4 Многопоточность

Поскольку на сегодняшний день конкуренция происходит в области затрат и скорости, работа многих приложений основана на многопоточности, которая рассматривается как один из способов увеличения производительности.

С ростом объема БД наблюдается преимущество по времени многопоточной реализации. Так, при работе с базой данных, состоящей из 100000 записей, время выполнения запросов примерно в 1000 раз выше у однопоточной программы. Также однопоточная модель показывает нестабильную работу на больших данных (ошибка OutOfMemory). На рисунке

1.3 приведены результаты сравнения работы однопоточной и многопоточной программ на больших данных. [9]

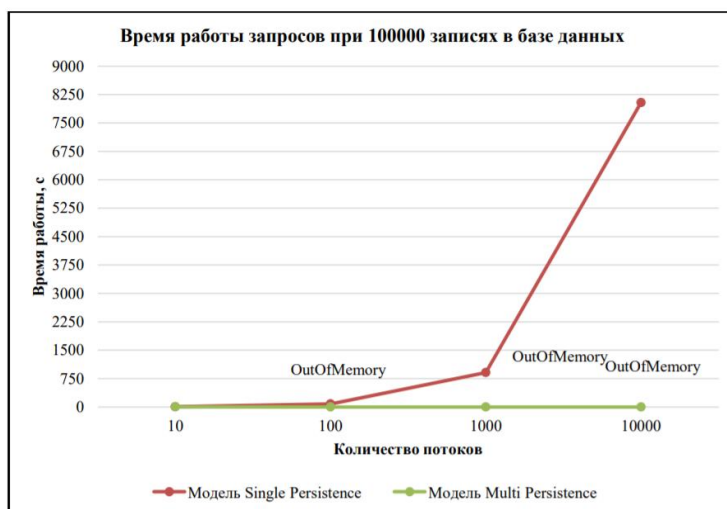


Рисунок 1.3 — Сравнение однопоточной и многопоточной программ, выполняющих запросы чтения БД.

Однако многопоточная реализация имеет свои недостатки. Многопоточные приложения, использующие PostgreSQL, вынуждены открывать новое соединение в каждом потоке.

Поскольку операция подключения — одна из самых дорогостоящих (процесса подключения к БД занимает от 2 до 3 МБ [10]), рост количества потоков может привести к замедлению работы программы: повышенная нагрузка на системные ресурсы и значительное снижение производительности, особенно на многоядерных системах. Увеличение конкуренции при обращении множества процессов к ресурсам PostgreSQL также способно замедлять работу программы.

Открытие соединений на разных потоках может привести к проблеме превышения количества подключений на сервере, что приведет к долгому ожиданию дальнейших запросов или их отклонению.

1.5 Пул соединений

Объектный пул представляет из себя набор инициализированных и готовых к использованию объектов. При необходимости системы обратиться к объекту, вместо его создания будет взят экземпляр из пула. Особенно заметно повышение производительности, когда стоимость и скорость инициализации экземпляра высоки, а количество одновременно используемых объектов в любой момент времени является низким. [11]

Из-за того, что PostgreSQL не имеет встроенного пула подключений [7], большинство клиентских программ вынуждены реализовывать свой собственный.

1.5.1 Frontend pool

При инициализации пула выполняется установка необходимого количества соединений. Предельный размер определяется пользователем в зависимости от контекста задачи. После успешной инициализации из пула может быть извлечено свободное соединение, для выполнения необходимых запросов к БД. После выполнения запросов соединение должно быть возвращено в пул. Если соединение было закрыто, его следует удалить из пула, и вместо него создать новое. [12]

Преимущество данного метода наглядно демонстрирует следующий пример. Устанавливается соединение для 10 клиентов, каждый из которых выполняет 10.000 запросов в БД. Если в среднем выполняется 486 транзакций в секунду, то реализация пула соединений (размером 25) позволяет увеличить это значение примерно на 60% – до 566 транзакций в секунду. [13]

Пул соединений имеет несколько недостатков, один из которых заключается в ограничении максимального количества одновременных подключений к БД. В зависимости от реализации, пользователь может задать размер пула, а также количество соединений, которое может быть добавлено в пул. К другому существенному недостатку следует отнести сложность

реализации, а также встраиваемость кода (особенно в крупных компаниях). Также следует обратить внимание на расчет следующих параметров: минимальное количество соединений, максимальное количество пулов соединений, максимальное время простоя, время ожидания соединения, количество попыток после тайм-аута. От корректной конфигурации пула зависит то, насколько увеличится пропускная способность транзакции.

1.5.2 Server pool

Хотя PostgreSQL не имеет встроенного пула подключений, он был реализован в коммерческой системе Postgres Pro Enterprise. Это объектно-реляционная СУБД, разработанная Postgres Professional в рамках проекта Postgres Pro на основе PostgreSQL [14]. В отличие от внешнего, встроенный пул не требует дополнительного обслуживания и не налагает на клиента никаких ограничений.

Работа встроенного пула аналогична работе внешнего. Число обслуживающих процессов, которые могут использоваться для отдельно взятой БД ограничивается размером пула. При достижении этого значения, процесс-демон Postgres перестает запускать новые процессы, а передает последующее подключение запущенному процессу. Так как один процесс может работать только с одной БД, возникает необходимость поддержки отдельного пула соединения для каждой БД. При появлении подключения к новой БД, добавляется новой пул. В параметрах могут быть указаны пользователи и БД, для которых не требуется реализация пула. Пулы функционируют только на уровне транзакций, т.е. процесс может переключиться на обслуживание нового соединения только после завершения транзакции. [14]

На рисунках 1.4 – 1.5 представлен цикл соединения с БД без пула и с его использованием соответственно. [15]

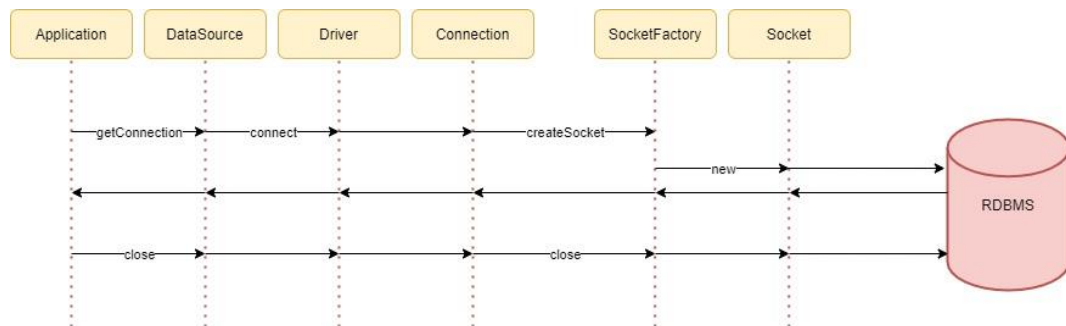


Рисунок 1.4 — Цикл соединения с БД.

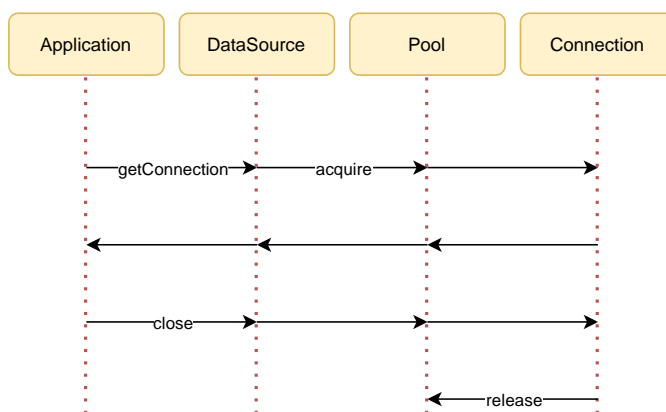


Рисунок 1.5 — Цикл соединения с БД с использованием пула.

Список использованных источников

1. ГОСТ 34.320-96 Информационные технологии. Система стандартов по базам данных. Концепции и терминология для концептуальной схемы и информационной базы.
2. ГОСТ Р ИСО/МЭК ТО 10032-2007: Эталонная модель управления данными.
3. Основные классы современных параллельных компьютеров. [Электронный ресурс]. Режим доступа: <https://parallel.ru/computers/classes.html>.
4. Lyman P., Varian H.R. How much information Архивная копия от 19 февраля 2018 на Wayback Machine. Release of the University of California. Oct.27, 2003. Knowledge Base of Relational and NoSQL Database Management Systems. [Электронный ресурс]. Режим доступа: <https://db-engines.com/en/ranking>.
6. Концепция архитектуры PostgreSQL. [Электронный ресурс]. Режим доступа: <http://www.dataved.ru/2014/09/postgresql.html>
7. PostgreSQL: Документация: 12.8. [Электронный ресурс]. Режим доступа: <https://postgrespro.ru/docs/postgresql/12>.
8. Пан К. С., Цымблер М. Л. Разработка параллельной СУБД на основе последовательной СУБД PostgreSQL с открытым исходным кодом // Вестник ЮУрГУ. Сер. Математическое моделирование и программирование. – 2012. – №18. – С. 277. Режим доступа: <https://cyberleninka.ru/article/n/razrabotka-parallelnoy-subd-na-osnove-posledovatelnoy-subd-postgresql-s-otkryтым-ishodnym-kodom>.
9. Воронова, Н. М. Алгоритмы оценки производительности модуля работы с данными / Н. М. Воронова, А. С. Кованова, Н. С. Корж // Инновации. Наука. Образование. – 2021. – № 38. – С. 647-658.
10. Shetty N. Everything you need to know about Connection Pooling in Postgres. – 2019. [Электронный ресурс]. Режим доступа: <https://www.ashnik.com/everything-you-need-to-know-about-connection-pooling-in-postgres/>.

- 11.Object Pool. [Электронный ресурс]. Режим доступа:
<https://www.oodesign.com/object-pool-pattern.html>.
- 12.Шиндов Д. А. Разработка пула соединений для работы с СУБД MYSQL на языке программирования C++ // ББК 1 А28. – 2021. – С. 112-113.
- 13.Aboagye M. Improve database performance with connection pooling. – 2020.
[Электронный ресурс]. Режим доступа:
<https://stackoverflow.blog/2020/10/14/improve-database-performance-with-connection-pooling/>.
- 14.Документация к Postgres Pro Enterprise 12.8.1. [Электронный ресурс]. Режим доступа: <https://postgrespro.ru/docs/enterprise/12>.
- 15.Shaikh S. S., Pachghare V. K. A Comparative Study of Database Connection Pooling Strategy // International Research Journal of Engineering and Technology. – 2017.