

Т.А. Паромова, Р.К. Кудерметов, Н.В. Луценко
Запорожский национальный технический университет, г. Запорожье
кафедра компьютерных систем и сетей
tparomova@mail.ru, krk@zntu.edu.ua

АНАЛИЗ ВЛИЯНИЯ ВИДА ФРАГМЕНТАЦИИ РАСПРЕДЕЛЕННОЙ БАЗЫ ДАННЫХ НА ВРЕМЯ ВЫПОЛНЕНИЯ ЗАПРОСОВ

Аннотация

Паромова Т.А., Кудерметов Р.К., Луценко Н.В. Анализ влияния вида фрагментации распределенной базы данных на время выполнения запросов. Обсуждаются результаты исследования зависимости времени выполнения запросов от размеров выборки при горизонтальной и вертикальной фрагментациях распределенной базы данных.

Ключевые слова: распределенная база данных, горизонтальная фрагментация, вертикальная фрагментация, производительность

Введение.

Системы управления базами данных (СУБД) являются одним из основных инструментов обработки больших объемов данных. В настоящее время с развитием сетевых технологий все больше используются распределенные базы данных. Это связано, прежде всего, с изменениями в организации бизнеса. Многие кооперации среднего и большого бизнеса, как правило, состоят из распределенных в пространстве подразделений.

При этом часто каждое подразделение может оперировать как с собственными наборами данных, так и с кооперативными. Разработка распределенных баз данных, отражающих организационную структуру бизнеса, преследует цель обеспечить общий и безопасный доступ к общим данным и, вместе с тем, иметь возможность локального хранения собственных данных [1].

Одной из важнейших характеристик распределенных СУБД является их производительность.

Поэтому оценка производительности СУБД является важным этапом проектирования распределенных базы данных, при этом такая оценка важна как для больших объемов баз данных, так и для средних.

Постановка задачи.

При тестировании и экспериментальных оценках быстродействия распределенных баз данных обычно определяется быстродействие без учета структурных свойств базы данных [2, 3], в частности, видов фрагментации. Кроме того, быстродействие распределенной базы данных вполне естественно зависит от характеристик программного обеспечения, аппаратной среды хранения, средств коммуникации, объемов выборки и распределения фрагментов базы данных в среде распределенных средств хранения данных.

В данной работе были поставлены следующие задачи:

1. Исследовать зависимости времени выполнения запросов к базе данных от вида фрагментации и количества распределенных фрагментов базы данных;
2. Получить аналитические выражения для оценки времени выполнения запросов при горизонтальной и вертикальной фрагментациях распределенной базы данных.

Результаты исследований.

Исследования проводились на локальной сети типа FastEthernet и организованных в кластер рабочих станциях. В качестве промежуточного программного обеспечения,

объединяющего рабочие станции в кластер, использовалось промежуточное программное обеспечение (middleware) *mpich* 1.2.4 [4], с помощью которого осуществлялся доступ к базе данных и взаимодействие клиентов с серверами. Экспериментальная распределенная база данных имела 20 таблиц, заполненных произвольными значениями, пропорционально прогнозируемому объему данных. Общий объем распределенной базы составлял 1 Гб.

Исследуемая база данных была распределена между 8 узлами кластера (серверами базы). При этом ставилась задача определения времени выполнения запросов к данным, которые разбиты между этими серверами.

На первом этапе исследовалась зависимость времени обработки запроса от размеров фрагментов распределенной базы данных, находящихся на двух серверах сети, при этом в качестве серверной СУБД использовалась *MS SQL Server*, а в качестве клиента – *СУБД Access*. Эксперименты проводились для выборки строк из горизонтально распределенной базы и столбцов – из вертикально распределенной базы.

В первом случае измерение времени выполнения запросов проводились для вертикально распределенной базы данных по числовому полю между двумя серверами сети. Реализовывалась выборка из 10 полей по индексу, при этом менялось соотношение количества выбираемых столбцов с первого и второго серверов от 0 до 10.

Измерение времени проводились для выборки из 1000 записей одного столбца для каждого варианта распределения базы.

Результаты выполнения запросов при вертикальной фрагментации базы данных представлены на рис. 1, где g_{i_j} – соотношение фрагментов базы данных первый сервер/второй сервер.

В результате аппроксимации полученных измерений времени получены следующие зависимости:

$$T = \frac{size1 \cdot N \cdot t_1}{1000} + \frac{size2 \cdot N \cdot t_2}{1000}, \quad (1)$$

$$size1 = \sum_{i=1}^{S_1} size_of_column_i, \quad (2)$$

$$size2 = \sum_{i=1}^{S_2} size_of_column_i, \quad (3)$$

где *size1* - объем одной записи выбираемых данных с первого сервера; S_1 – количество столбцов, выбираемых с первого сервера; *size2* - объем одной записи выбираемых данных со второго сервера; S_2 – количество столбцов, выбираемых со второго сервера; *N* – количество выбираемых строк; t_1 – время выборки 1000 записей с первого сервера; t_2 – время выборки 1000 записей со второго сервера.

На основе экспериментальных данных можно сделать вывод о зависимости времени обработки запроса от расположения выбираемых столбцов таблицы между серверами сети. Проведенный анализ соответствия экспериментальных и расчетных данных показал, что точность расчетов составляет 4%.

При исследовании горизонтального распределения этой же базы между двумя серверами сети производилась выборка 250 строк по индексу, при этом менялось соотношение количества выбираемых строк с первого и второго серверов.

Среднее время выборки определялось для 1000 записей в каждом случае.

Результаты исследований представлены на рис. 2, где g_{i_j} – соотношение фрагментов базы данных первый сервер/второй сервер.

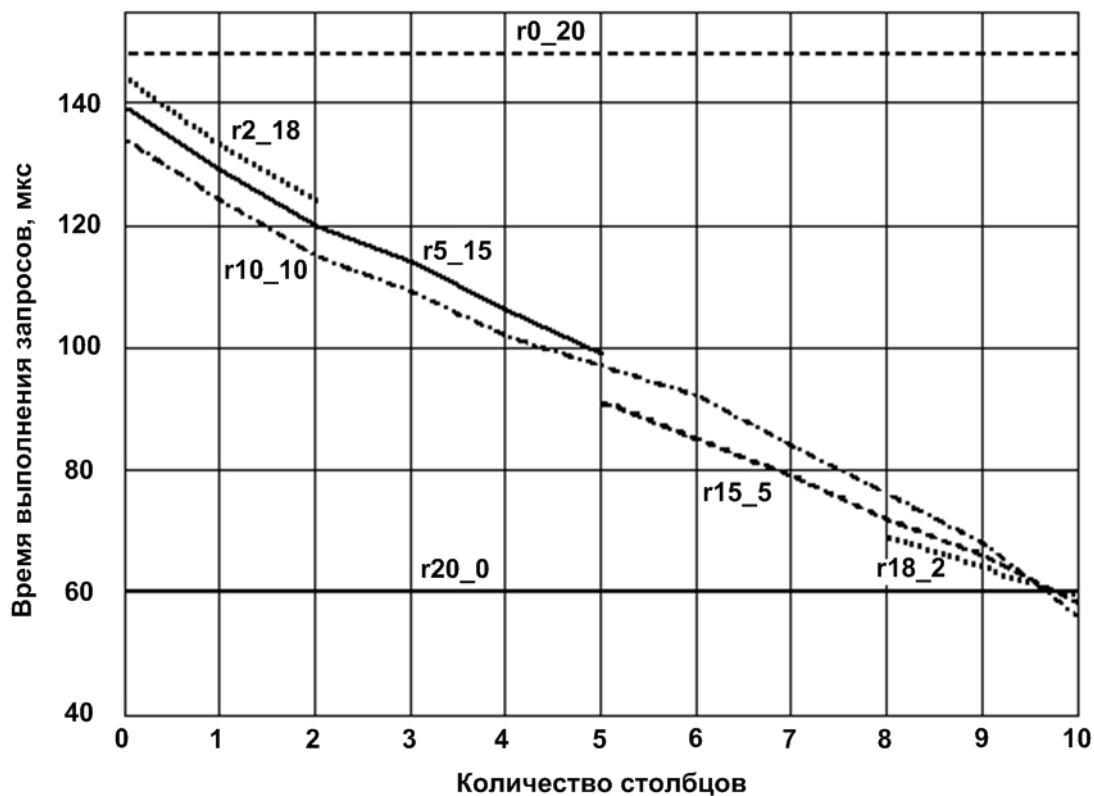


Рисунок 1 - Результаты выполнения запросов при вертикальной фрагментации запросов

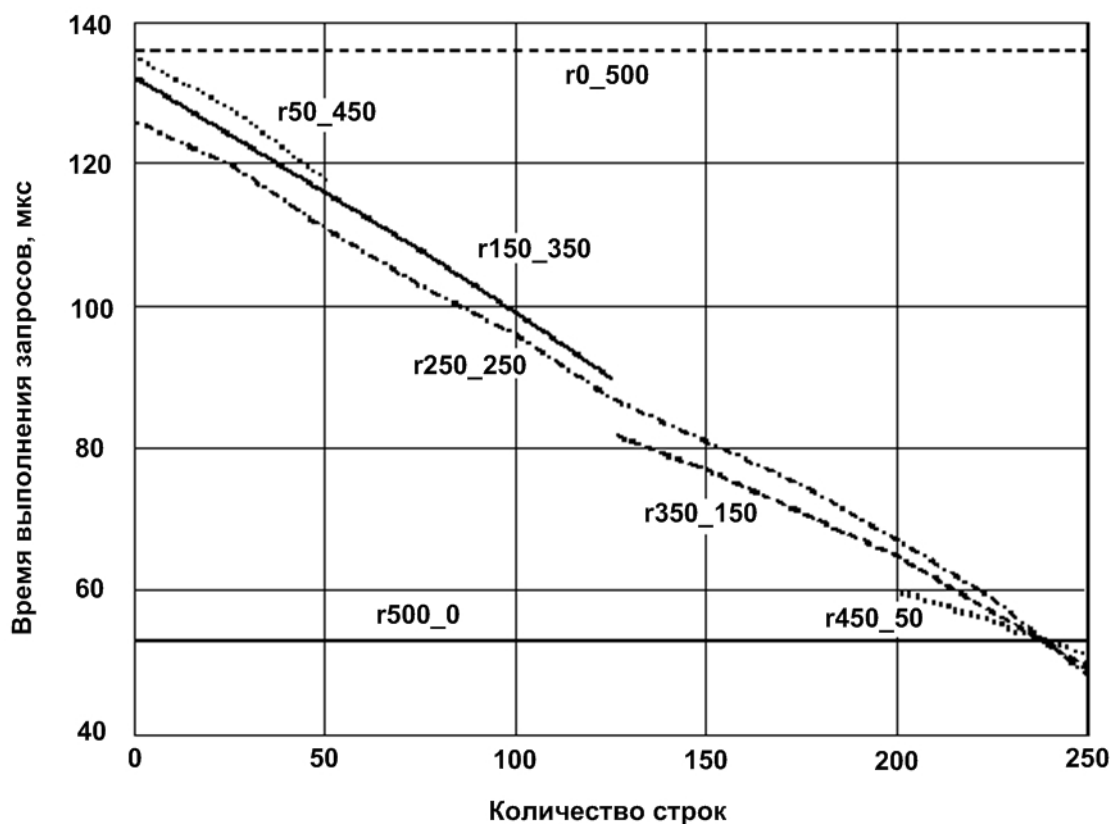


Рисунок 2 - Результаты выполнения запросов при горизонтальной фрагментации запросов

Полученные экспериментальные данные измерения времени можно аппроксимировать с помощью выражений:

$$T = \frac{size \cdot N_1 \cdot t_1}{1000} + \frac{size \cdot N_2 \cdot t_2}{1000}, \quad (4)$$

$$size = \sum_{i=1}^{col_num} size_of_column_i, \quad (5)$$

где $size$ – объем одной записи таблицы; col_num – количество столбцов фрагментируемой таблицы; N_1 – количество выбираемых строк с первого сервера; N_2 – количество выбираемых строк со второго сервера.

Максимальное отклонение значений времени запросов с использованием выражения (4) и экспериментальных данных составляет 4%.

Из выражения (1) можно получить аналитические оценки для базы данных, вертикально фрагментированной между произвольным числом серверов.

Так, например для 8 серверов оценочным выражением будет:

$$T = \frac{size_1 \cdot N \cdot t_1}{1000} + \dots + \frac{size_8 \cdot N \cdot t_8}{1000} = \frac{N}{1000} \sum_{j=1}^8 size_j \cdot t_j, \quad (6)$$

где $size_j = \sum_{i=1}^{S_j} size_of_column_i$ – объем одной записи выбираемых данных из j -го сервера;

S_j – количество столбцов, выбираемых с j -го сервера; t_j – время выборки 1000 записей с j -го сервера.

Аналогично из выражения (4) можно получить оценочное выражения для произвольного числа серверов при горизонтальной фрагментации распределенной базы данных. В частности, для 8 серверов получим:

$$T = \frac{size \cdot N_1 \cdot t_1}{1000} + \dots + \frac{size \cdot N_8 \cdot t_8}{1000} = \frac{size}{1000} \sum_{j=1}^8 N_j \cdot t_j, \quad (7)$$

где N_j – количество выбираемых строк с j -го сервера.

Сравнение времени выполнения запросов позволяет сделать вывод, что горизонтальная фрагментация распределенной базы данных предпочтительнее, т.к. обеспечивает меньшие затраты времени на выполнение запросов (в среднем 10%).

Именно этот способ распределения фрагментов базы данных чаще всего используется на практике.

Выводы.

В результате проведенных исследований получены выражения для аналитических оценок времени выполнения запросов к распределенной базе данных при её горизонтальной и вертикальной фрагментациях.

Погрешности оценок не превышают 4%.

В дальнейшем необходимо проверить справедливость полученных оценок для баз данных с большим числом записей, а также исследовать характер зависимости времени выборки от размеров и типов извлекаемых записей.

На практике большая часть функционирующих баз данных является частью OLTP – систем, в которые постоянно вносится большое количество записей.

Также представляет интерес исследовать затраты времени для различных видов запросов в зависимости от количества операций соединения, так как именно эта операция при выборке данных из нескольких таблиц требует наибольшего времени.

Литература

1. Конноли Т. Базы данных. Проектирование, реализация и сопровождение. Теория и практика / Т. Конноли, К. Бегг, А. Страчан. - Киев: Вильямс, 2001. – 1111 с.
2. Елашкин М. Производительность СУБД и тесты ТРС / М. Елашкин //ByteРоссия, 2004. - №3(67). –с.17. [Электронный ресурс] – Режим доступа: [http:// www.bytemag.ru](http://www.bytemag.ru)
3. Ковригин Д. Сравнительное тестирование трех СУБД на реальной информационной системе / Д. Ковригин, А.Николаев // [Электронный ресурс] - Режим доступа: <http://www.computerinform.ru>
4. Шпаковский Г.И. Программирование для многопроцессорных систем в стандарте MPI / Г.И. Шпаковский, Н.В. Серикова. - Минск: БГУ, 2002. - 323с.

Abstract

Paromova T.O., Kudermetov R.K., Lucenko N.V. The impact analysis of fragmentation kind of distributed database to a queries time. The research results of a queries time dependency from the size of sampling out of distributed database at a horizontal and vertical fragmentation are discussed. Expressions for an analytical estimation of queries time dependency from kind of a database fragmentation are received.

Keywords: *distributed database, horizontal partitioning, vertical partitioning, productivity*

Анотація

Паромова Т.О., Кудерметов Р.К., Луценко Н.В. Аналіз впливу виду фрагментації розподіленої бази даних на час виконання запитів. У статті обговорюються результати дослідження залежності часу виконання запитів від розмірів вибірки для горизонтальної та вертикальної фрагментації розподіленої бази даних.

Ключові слова: *розподілена база даних, горизонтальна фрагментація, вертикальна фрагментація, продуктивність*

Здано в редакцію:
30.04.2010р.

Рекомендовано до друку:
к.т.н, доц. Зеленьова І.Я.