

**Харківський національний університет
радіоелектроніки**

КВАЛІФІКАЦІЙНА РОБОТА

«Методи обробки запитів в системах паралельних баз даних»

Виконав: студент групи СПм-20-1 Бессараб Є.В.

Керівник: доц. каф. ЕОМ Філімончук Т.В.

Аналіз предметної області

2

Метою кваліфікаційної роботи є дослідження методів паралельного виконання SQL - запитів у великомасштабних системах для підвищення продуктивності розподілених СКБД та зменшення часу виконання запиту.

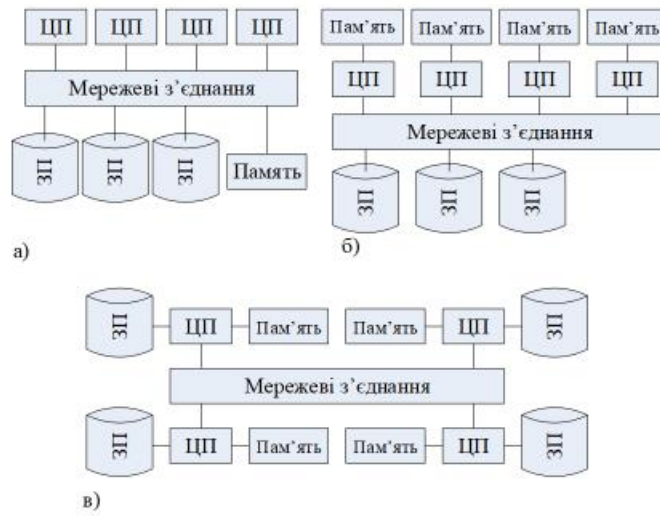
Об'єкт дослідження: методи паралельного виконання SQL-запитів.

Завдання:

- докладний аналіз ефективності існуючих методів розпаралелювання запитів у гетерогенних середовищах;
- дослідження алгоритмів та програмних засобів для паралельного виконання SQL – запитів, які забезпечують сумісність з існуючими СКБД;
- отримання апріорних оцінок часу виконання запитів;
- отримання експериментальної оцінки методів та алгоритмів, які були використані для розпаралелювання запитів у паралельних системах баз даних.

Архітектури паралельних систем

3



Класифікація Стоунбрекера для паралельних систем

4

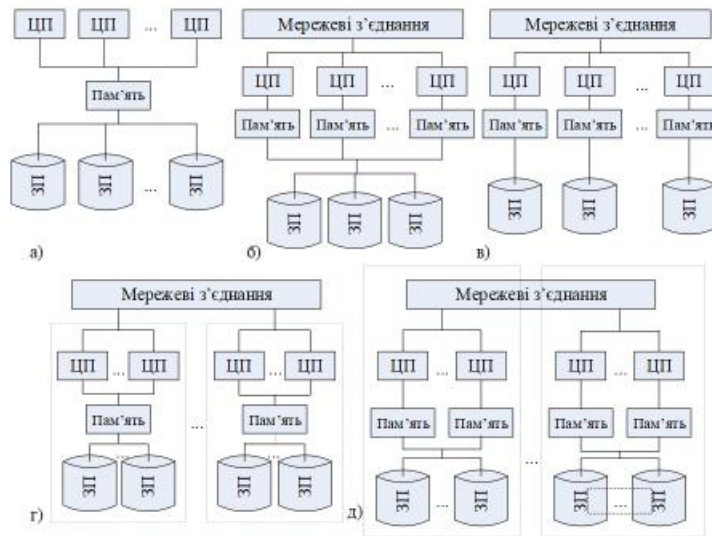
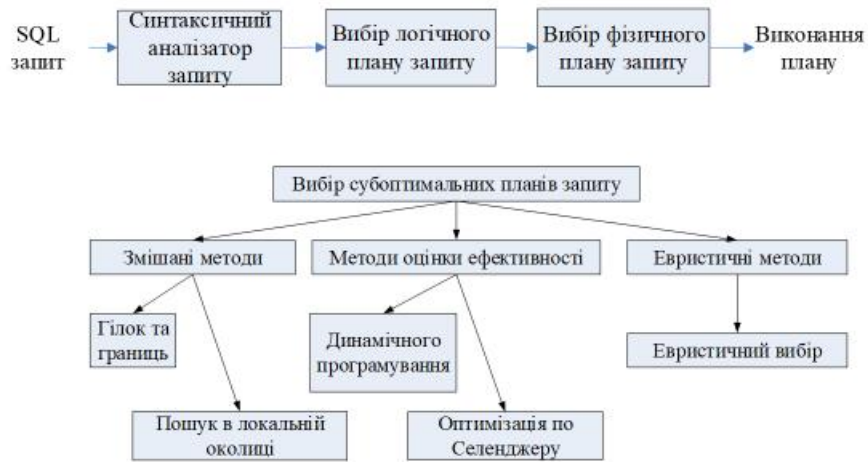


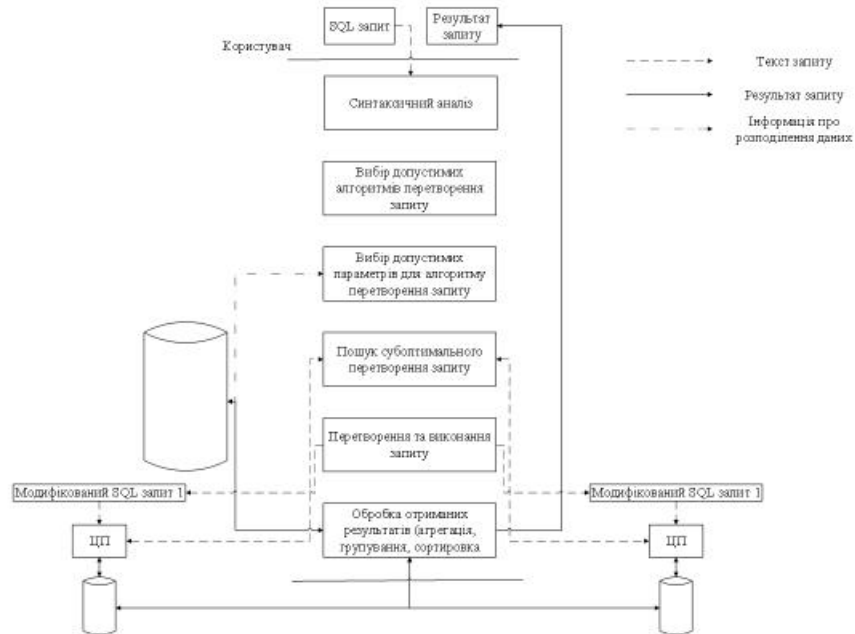
Схема компіляції SQL-запиту. Методи вибору субоптимального запиту

5

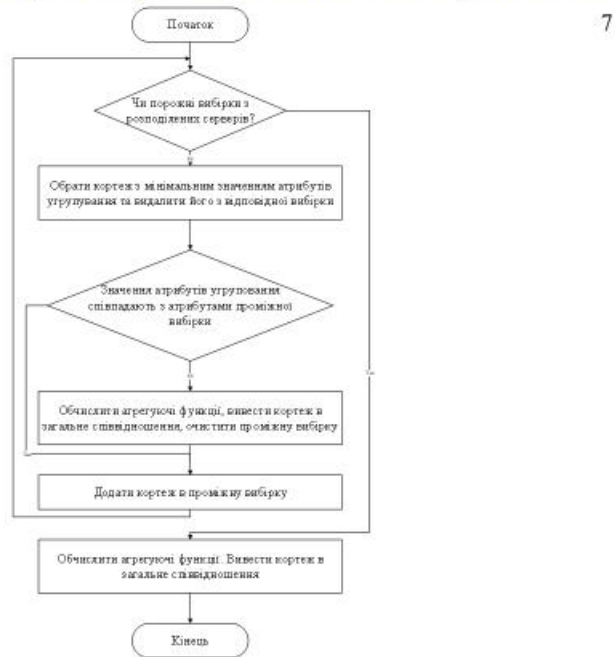


Архітектура системи для паралельного використання запитів

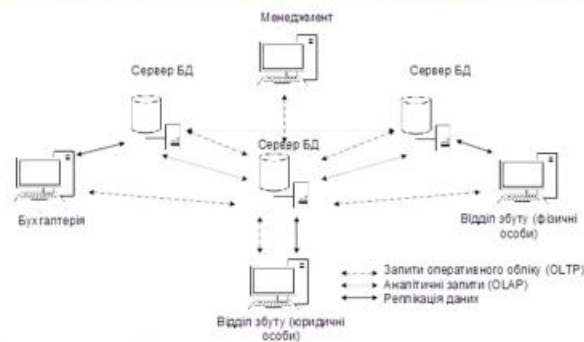
6



Агрегація атрибутів проміжних співвідношень

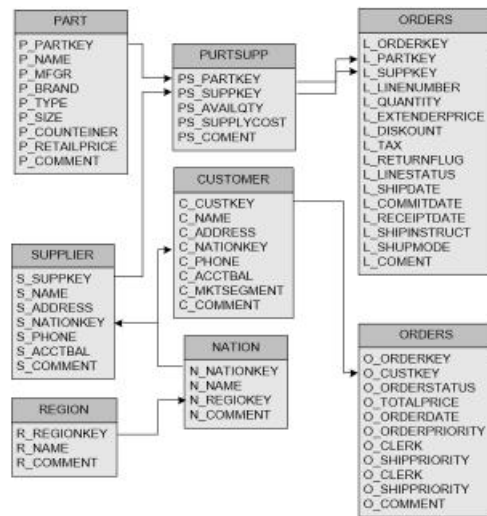


Структура інформаційної системи підприємства



Назва атрибуту	Тип	Коментар
SysID	int	Унікальний ключ опису індексу
TableName	varchar(255)	Ім'я відношення, якому належить атрибут
FieldName	varchar(255)	Назва атрибуту в таблиці
IndexName	varchar(255)	Назва індексу
IndexType	int	Тип індексу: 0 = B-Tree, 1= Cluster
FieldType	varchar(255)	Тип атрибуту
MinValue	varchar(255)	Мінімальне значення атрибуту на всьому співвідношенні
MaxValue	varchar(255)	Максимальне значення атрибуту на всьому співвідношенні

Схема тестової бази даних



Результати виконання Q1

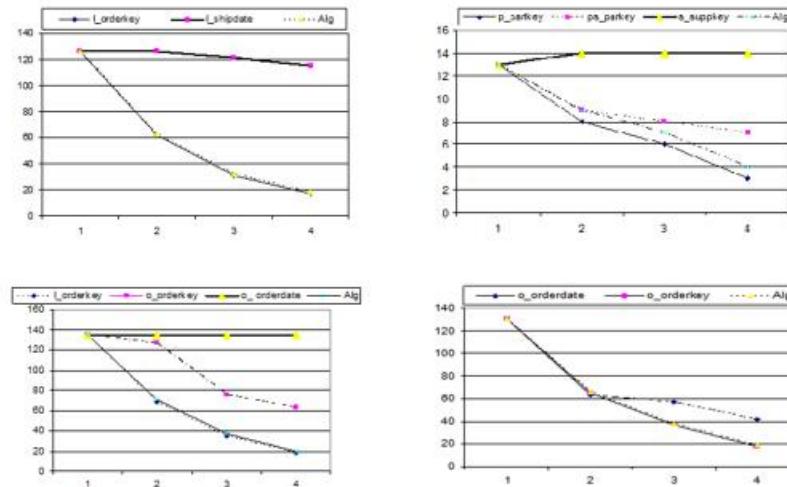
Кількість вузлів	Рядків у виборці	Проміжна вибірка	Додаткові плани запиту	Загальний час	K_y
1	4	0	0	126	1
2	4	8	1	62	1.02
4	4	16	1	31	1.02
8	4	32	1	17	0.93

Кількість вузлів	Рядків у виборці	Проміжна вибірка	Додаткові плани запиту	Загальний час	K_y
1	4	0	0	126	1
2	4	8	1	126	0.5
4	4	16	1	121	0.26
8	4	32	1	115	0.14

Кількість вузлів	Рядків у виборці	Проміжна вибірка	Додаткові плани запиту	Загальний час	K_y
1	4	0	0	126	1
2	4	8	2	63	1
3	4	16	2	32	0.98
8	4	32	2	18	0.88

Результати

11



Висновки

12

З використанням різних архітектур систем для паралельного виконання запитів з використанням надлишкової інформації було розглянуто та досліджено методи розпаралелювання SQL-запитів у СУБД на основі еквівалентного перетворення запитів до форми, що допускає паралельне виконання для застосування в гетерогенних середовищах. Внаслідок чого були отримані апріорні оцінки загального часу паралельного виконання перетворених SQL-запитів, що відрізняються урахуванням особливостей виконання запитів у гетерогенних системах з реплікацією даних, і дозволяють швидко оцінити доцільність досліджуваних методів розпаралелювання запитів на етапі конфігурування системи. Найкращі загальні показники продуктивності системи можна досягти, розробляючи систему на основі MPP архітектури, допускаючи використання SMP машин як вузлів системи, що фактично перетворює її на гібридну.

Програмний інтерфейс та результати

13

Ім'я класу	Описання
FileSet	Містить множини елементів класу File, що описує файли каталогу
File	містить ім'я, шпак і вміст файлу, а також множини елементів класу FileBlock
FileBlock	має як собою блок файлу, що включає текстовий вміст і множини відбитків
ShinglesFactory	статичний клас, що виробляє відбитки на основі заданого тексту і набору параметрів
ShinglesCollector	скрипчик відбитків і інтерфейс для зручної роботи з ними
PageProcessor	надає інтерфейс для роботи із конкретною сторінкою в частині «Виділення блоків»
Log	надає інтерфейс для роботи із конкретною сторінкою в частині «Виділення блоків»
GeneticAlg	клас для роботи з реалізацією генетичного алгоритму бібліотеки JGAP
FitnessFunc	клас, що надає функцію фітнес-функцію в відповідній зі специфікацією бібліотеки JGAP
StringExtractor	статичний клас, що надає можливість витягнення контенту з веб-сторінки на базі бібліотеки HTMLParser
VisitorRemover	клас, що надає можливість видалення вузлів і гілок DOM-дереві шляхом використання бібліотеки HTMLParser
VisitorCalculator	реалізує просторову коєфіцієнт вузлів дерева, дозволяючи працювати вагою функції алгоритму витягнення блоків, працює з використанням бібліотеки HTMLParser
BlockListTableModel, FileSetTableModel, JTableCellRenderer, MyHTMLTreeCellRenderer, NodeTableModel, ShingleListTableModel	розширення стандартних Java-класів візуальних компонентів, що мають додаткові можливості в частині візуалізації та інтерфейсу

Node	Depth	Strukt	Links	Weight	Total	W/T
DocType Tag (Def...	1	1	0	0	74	0
Tag langID = 1. Al...	1	451	14	1390	15918	0
HEAD lang...	2	6	0	39	309	0.113
TITLE Урпаво...	3	1	0	39	62	3.079
Script Node: Prog...	3	1	0	0	41	0
Tag meta lang...	3	1	0	0	85	0

Поиск	Сходимость	Блоки + Сходимость	Шпакли	Блоки + Шпакли
1.00	0.85	0.85	0.95	0.85
0.95	0.55	0.80	0.80	0.80
0.90	0.45	0.75	0.75	0.75
0.85	0.40	0.65	0.65	0.65
0.80	0.35	0.55	0.55	0.55

Висновки

14

Запропоновано модель, що дозволяє оцінювати схожість документів, використовуючи інформацію про схожість їх блоків. Відповідно до моделі розроблено метод оцінки схожості документів на основі складових їх структурно-семантичних блоків, що дозволяє поліпшити якість розпізнавання дублікатів за рахунок збільшення середнього значення показника повноти. Це дозволяє отримувати більш повну вибірку схожих документів на одних і тих же наборах даних. Проведена експериментальна перевірка запропонованих алгоритмів і методів, у розробленому програмному забезпеченні.