

ДОДАТОК А

Графічний матеріал кваліфікаційної роботи

**Харківський національний університет
радіоелектроніки**

Кафедра ЕОМ

Кваліфікаційна робота

**Методи інтелектуального аналізу
великих даних за допомогою
машинного навчання**

Виконала:

ст. гр. СПЗМ-23-1

Фатій Л.С.

Перевірив:

ас. Кравченко П.О.

2

Об'єкт дослідження та мета роботи

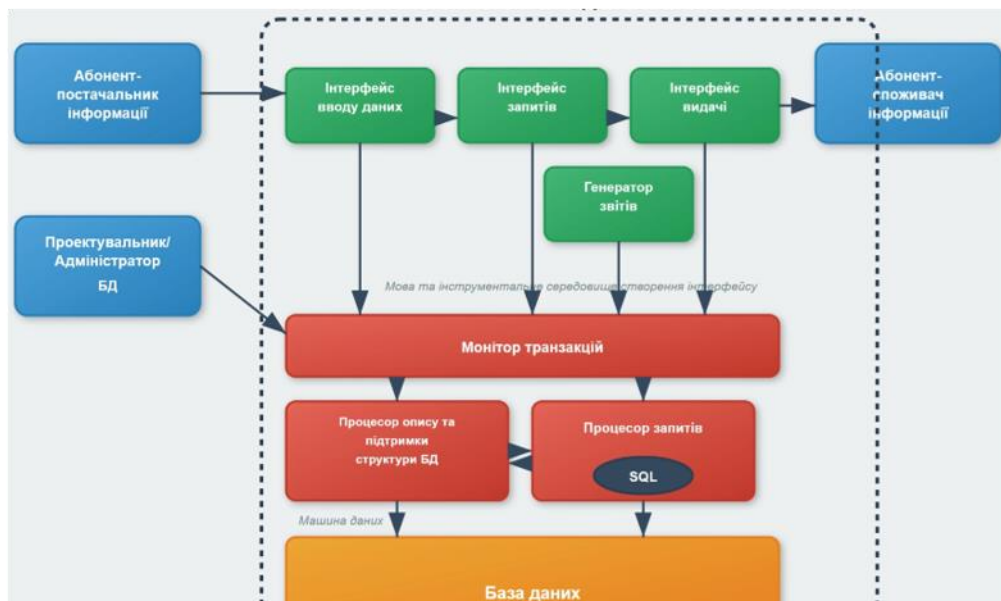
Метою кваліфікаційної роботи є аналіз методів інтелектуального аналізу великих даних в паралельних системах керування базами даних.

Об'єктом дослідження є методи ІАД в реляційних СКБД.

Завдання:

- розглянути методи та алгоритми для впровадження інтелектуального аналізу великих даних в паралельну СКБД;
- проаналізувати паралельні алгоритми рішення завдань кластеризації великих даних засобами паралельної реляційної СКБД;
- розробити архітектуру паралельної СКБД з використанням модифікованого методу.

Функціональна структура системи керування базами даних ³



Інтелектуальний аналіз великих даних ⁴

```

Algorithm: Знайти найближчу пару з обмеженням відстані
Input: S_I^m - множина елементів, n - мінімальна відстань індексів
Output: {pos_bsf, dist_bsf} - позиція та відстань найкращої пари

1: // Ініціалізація змінних
2: dist_bsf ← ∞ // найкраща відстань
3: pos_bsf ← -1 // позиція найкращої пари
4:
5: // Основний цикл по всіх елементах
6: for i = 1 to |S_I^m| do
7:   current_element ← S_I^m[i]
8:   dist_min ← ∞ // мінімальна відстань для поточного елемента
9:
10:  // Пошук найближчого елемента з обмеженням
11:  for j = 1 to |S_I^m| do
12:    if |i - j| ≥ n then // перевірка обмеження відстані індексів
13:      dist ← EditDistance(S_I^m[i], S_I^m[j])
14:
15:      // Раннє припинення, якщо знайдено дуже близьку пару
16:      if dist < dist_bsf then
17:        break // перервати внутрішній цикл
18:      end if
19:
20:      // Оновити мінімальну відстань для поточного елемента
21:      if dist < dist_min then
22:        dist_min ← dist
23:      end if
24:    end if
25:  end for
26:
27:  // Оновити глобальний мінімум, якщо знайдено кращу пару
28:  if dist_min < dist_bsf then
29:    dist_bsf ← dist_min
30:    pos_bsf ← i
31:  end if
32: end for
33:
34: return {pos_bsf, dist_bsf}

```

Приклади запитів

5

```

-----
-- БЛОК 1: Пошук асоціативних правил для пневмонії у пацієнтів старіє 60
-----
-- Знаходження асоціативних правил для пневмонії у пацієнтів старіє 60
FIND ASSOCIATION RULES AS HealthRuleSet
RELATED TO Salary, Age, isSmoker, Disease
FROM HealthDB
WHERE Disease = 'Pneumonia' AND Age > 60
WITH SUPPORT_THRESHOLD = 0.05
WITH CONFIDENCE_THRESHOLD = 0.67;

-----
-- БЛОК 2: Виводуток специфічних правил з набору даних
-----
-- Виводуток правил для аналізу зв'язку між курінням та пневмонією
MINE RULE HealthRuleSet AS
SELECT DISTINCT
  1..n Disease AS antecedent, -- передумова
  1..1 isSmoker AS consequent -- наслідок
FROM HealthDB
WHERE antecedent.Disease = 'Pneumonia'
AND antecedent.Age > 60
EXTRACTING RULES WITH
  SUPPORT: 0.1,
  CONFIDENCE: 0.3;

-----
-- БЛОК 3: Функція отримання правил з бази даних
-----

FUNCTION GetRules(HealthDB)
RETURNS HealthRuleSet AS
BEGIN
  -- Отримання правил з підтримкою та довірою
  SELECT R INTO HealthRuleSet
  FROM HealthDB AS R
  WHERE R.Antecedent IN {
    (Disease = 'Pneumonia'),
    (Age > 60),
    (Salary > 1000)
  }
  AND R.Consequent IN {(isSmoker = TRUE)};
  AND R.Support > 0.1
  AND R.Confidence > 0.7;

  RETURN HealthRuleSet;
END;

-----
-- БЛОК 4: Визначення та вивід найкращих правил
-----

FUNCTION SelectRules(HealthRuleSet)
RETURNS FilteredRuleSet AS
BEGIN
  -- Вивід правил з високим підтримкою для пневмонії
  SELECT * FROM HealthRuleSet
  WHERE Antecedent HAS ((Disease = 'Pneumonia'))
  AND Consequent HAS {
    (Salary > 0),
    (Salary < 10000)
  }
  AND Support > 0.1
  AND Confidence > 0.7
  AND Lift > 1.0; -- додатковий критерій вибору
END;

-----
-- БЛОК 5: Прогнозування асоціацій з використанням ML моделі
-----
-- Прогнозування ризику пневмонії для конкретного пацієнта
SELECT
  PredictAssociation(
    [HealthMiningModel].[AssocLines],
    INCLUDE_STATISTICS,
    3 -- top-3 асоціації
  ) AS PredictionsResults
FROM [HealthMiningModel]
NATURAL PREDICTION JOIN (
  SELECT
    68 AS [Age],
    TRUE AS [isSmoker],
    'Pneumonia' AS [Disease]
  ) AS [PatientProfile];

-----
-- БЛОК 6: Додатковий аналіз для покращення результатів
-----
-- Аналіз кореляцій всіх факторів ризику
WITH RiskFactors AS (
  SELECT
    Age,
    isSmoker,
    Disease,
    Salary,
    COUNT(*) AS frequency,
    AVG(CASE WHEN Disease = 'Pneumonia' THEN 1.0 ELSE 0.0 END) AS pneumonia_rate
  FROM HealthDB
  WHERE Age > 50
  GROUP BY Age, isSmoker, Salary
)
SELECT
  *,
  pneumonia_rate * 100 AS risk_percentage
FROM RiskFactors
WHERE pneumonia_rate > 0.1
ORDER BY pneumonia_rate DESC;

```

6

Обробка запиту великих даних

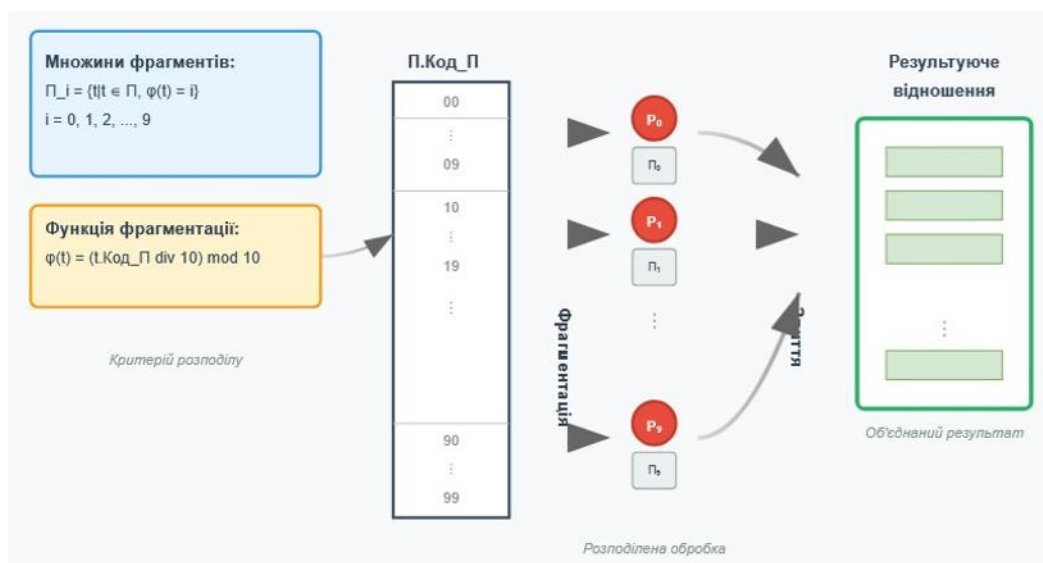
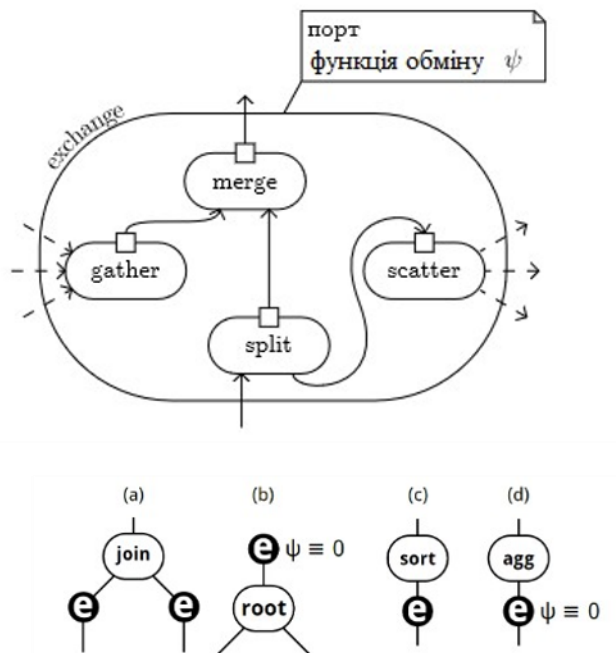


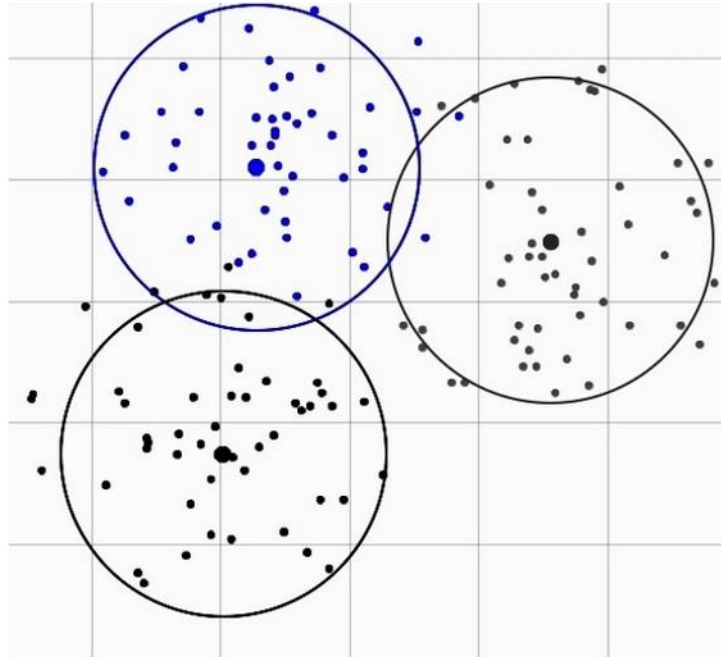
Схема обробки запиту в паралельній базі даних



Структура оператора обміну та паралельний план запиту



Кластеризація та пошук шаблонів



Метод pgFCM

```

-----
-- ОСНОВНИЙ АЛГОРИТМ FUZZY C-MEANS
-- Математичні формули та SQL реалізація
-----

/*
НАСТАНОВКИ ОСНОВНИЙ АЛГОРИТМУ:
1. Цільова функція (мінімізується):
J,FOR(X,k,n) = \sum_{i=1}^m \sum_{j=1}^k u_{ij}^m * d(x_i, c_j)^2
2. Означення центроїдів:
V_j, k, c_j = (\sum_{i=1}^m u_{ij}^m * x_i) / (\sum_{i=1}^m u_{ij}^m)
3. Означення коефіцієнтів членства:
u_{ij} = \frac{1}{\sum_{k=1}^k (d(x_i, c_k) / d(x_i, c_j))^{2/(1-m)}}
4. Критерій збіжності:
max |u_{ij}^{(n+1)} - u_{ij}^{(n)}| < \epsilon
Де:
- x = (x_1, x_2, ..., x_n) - координати об'єкта
- k - кількість кластерів
- m - параметр м'яккості (зазвичай m = 2)
- u_{ij} - ступінь належності об'єкта i до кластера j
- c_j - центроїд кластера j
- d(x_i, c_j) - відстань між об'єктом i та центроїдом j
- \epsilon - поріг збіжності
*/

-----
-- СТРУКТУРА ДАНИХ: Створення тимчасових таблиць
-----

-- Таблиця матриць членства U(i,j) = ступінь належності об'єкта i до кластера j
CREATE TEMP TABLE MEMBERSHIP_MATRIX (
  object_id INTEGER NOT NULL, -- ідентифікатор об'єкта
  cluster_id INTEGER NOT NULL, -- ідентифікатор кластера
  membership_value NUMERIC(10,4), -- коефіцієнт u_{ij}
  PRIMARY KEY (object_id, cluster_id)
);

-- Таблиця параметрів алгоритму FCM(x,n,s,delta)
CREATE TEMP TABLE ALGORITHM_PARAMETERS (
  dimensions_d INTEGER, -- розмірність простору ознак
  clusters_k INTEGER, -- кількість кластерів
  objects_n INTEGER, -- кількість об'єктів
  iteration_s INTEGER, -- номер поточної ітерації
  convergence_delta NUMERIC(10,8), -- критерій збіжності
  PRIMARY KEY (iteration_s)
);

-- Додаткова таблиця для збереження поточних значень SV(i,1)
CREATE TEMP TABLE SAVED_VALUES (
  object_id INTEGER NOT NULL, -- ідентифікатор об'єкта
  feature_id INTEGER NOT NULL, -- ідентифікатор ознаки
  feature_value NUMERIC(12,6), -- значення ознаки
  PRIMARY KEY (object_id, feature_id)

```

```

-----
-- ІНІЦІАЛІЗАЦІЯ: Початкові значення та параметри
-----

-- Збереження поточного стану матриць членства
INSERT INTO SAVED_VALUES (object_id, feature_id, feature_value)
SELECT object_id, 1, membership_value
FROM MEMBERSHIP_MATRIX;

-- Збереження додаткових параметрів (розміри ознак)
INSERT INTO SAVED_VALUES (object_id, feature_id, feature_value)
SELECT object_id, dimensions_d, extended_value
FROM EXTENDED_FEATURES;

-- Ініціалізація параметрів алгоритму
INSERT INTO ALGORITHM_PARAMETERS (dimensions_d, clusters_k, objects_n, iteration_s, convergence_delta)
VALUES (c_value, k_value, o_value, 0, 0.0);

-- Генерація випадкових початкових значень матриць членства
INSERT INTO MEMBERSHIP_MATRIX (object_id, cluster_id, membership_value)
VALUES (1, 1, RAND(0)),
       (1, 2, RAND(0)),
       ... (повторюючи для всіх комбінацій i,j)
       (c_value, k_value, RAND(0));

-----
-- ДОСКОНИЙ ЦИКЛ АЛГОРИТМУ
-----

-- КРОК 1: Означення матриць членства за формулою FCM
UPDATE MEMBERSHIP_MATRIX
SET membership_value = (
  SELECT
    new_membership.calculated_value / normalization.total_sum AS normalized_membership
  FROM (
    -- Обчислення чисельника формули означення
    SELECT
      un.object_id,
      un.cluster_id,
      SUM(un.membership_value) AS calculated_value
    FROM MEMBERSHIP_MATRIX AS un
    GROUP BY un.object_id, un.cluster_id
  ) AS new_membership,
  (
    -- Обчислення знаменника (нормалізація)
    SELECT
      object_id,
      SUM(membership_value) AS total_sum
    FROM MEMBERSHIP_MATRIX
    GROUP BY object_id
  ) AS normalization
  WHERE new_membership.object_id = normalization.object_id
  AND new_membership.object_id = MEMBERSHIP_MATRIX.object_id
  AND new_membership.cluster_id = MEMBERSHIP_MATRIX.cluster_id
);

```

Реалізація обчислень методу та оновлення таблиць

```

-- АЛГОРИТМ НЕЧІТКОЇ КЛАСТЕРИЗАЦІЇ FUZZY C-MEANS
-----
-- КРОК 1: Обчислення центроїдів кластерів
-----
-- Розрахунок координат центрів кластерів на основі зважених середніх
INSERT INTO CENTROIDS (cluster_id, feature_id, centroid_value)
SELECT
  cluster_j,
  feature_l,
  -- Формула зваженого середнього для центроїда
  SUM(membership.weight * feature_values.value) /
  SUM(membership.weight) AS centroid_value
FROM (
  -- Підзапит для отримання ваг членства та значень ознак
  SELECT
    object_i,
    cluster_j,
    POWER(membership_matrix.value, fuzziness_m) AS weight
  FROM MEMBERSHIP_MATRIX AS membership_matrix
) AS membership,
FEATURE_VALUES AS feature_values
WHERE membership.object_i = feature_values.object_i
GROUP BY cluster_j, feature_l;

-- КРОК 2: Обчислення стандартних відхилень (відстаней)
-----
-- Розрахунок евклідової відстані від об'єктів до центроїдів
INSERT INTO STANDARD_DEVIATIONS (object_i, cluster_j, distance)
SELECT
  object_i,
  cluster_j,
  -- Евклідова відстань до центроїда кластера
  SQRT(SUM(POWER(feature_values.value - centroids.centroid_value, 2))) AS distance
FROM FEATURE_VALUES AS feature_values,
CENTROIDS AS centroids
WHERE feature_values.feature_l = centroids.feature_id
AND feature_values.object_i = centroids.cluster_id
GROUP BY object_i, cluster_j;

```

```

-----
-- КРОК 3: Оновлення матриці членства (нечітких коєфіцієнтів)
-----
-- Обчислення нових значень функції членства для кожного об'єкта
INSERT INTO UPDATED_MEMBERSHIP (object_i, cluster_j, membership_value)
SELECT
  current_cluster.object_i,
  current_cluster.cluster_j,
  -- Формула оновлення членства у Fuzzy C-Means
  POWER(current_cluster.distance, POWER(2.0, (1.0 - fuzziness_parameter))) *
  denominator.density_sum AS membership_value
FROM (
  -- Підзапит для обчислення значення формули членства
  SELECT
    object_i,
    1.0 / SUM(POWER(distance, POWER(2.0, (fuzziness_parameter - 1.0)))) AS density_sum
  FROM STANDARD_DEVIATIONS
  GROUP BY object_i
) AS denominator,
STANDARD_DEVIATIONS AS current_cluster
WHERE denominator.object_i = current_cluster.object_i;

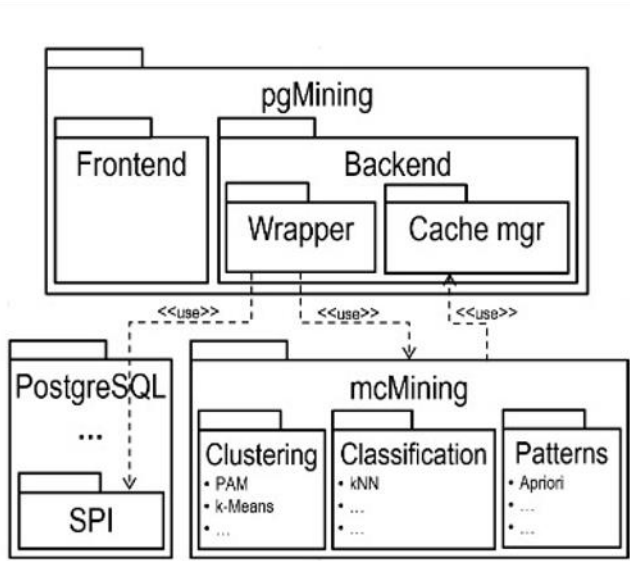
-- КРОК 4: Перевірка критерію збіжності
-----
-- Обчислення максимальної зміни у матриці членства для критерію зупинки
SELECT
  MAX(ABS(new_membership.membership_value - old_membership.membership_value))
INTO convergence_threshold
FROM OLD_MEMBERSHIP AS old_membership,
UPDATED_MEMBERSHIP AS new_membership
WHERE old_membership.object_i = new_membership.object_i
AND old_membership.cluster_j = new_membership.cluster_j;

-- КРОК 5: Збереження найкращого стану і ітерації та оновлення даних
-----
-- Логуючи параметри поточної ітерації
INSERT INTO ITERATION_LOG (
  dimensions,
  clusters_count,
  objects_count,
  iteration_step,
  convergence_value
) VALUES (
  dimension_d,
  clusters_k,
  objects_n,
  current_steps,
  convergence_threshold
);

-- Оновлення матриці членства для наступної ітерації
TRUNCATE TABLE MEMBERSHIP_MATRIX;
INSERT INTO MEMBERSHIP_MATRIX
SELECT object_i, cluster_j, membership_value
FROM UPDATED_MEMBERSHIP;

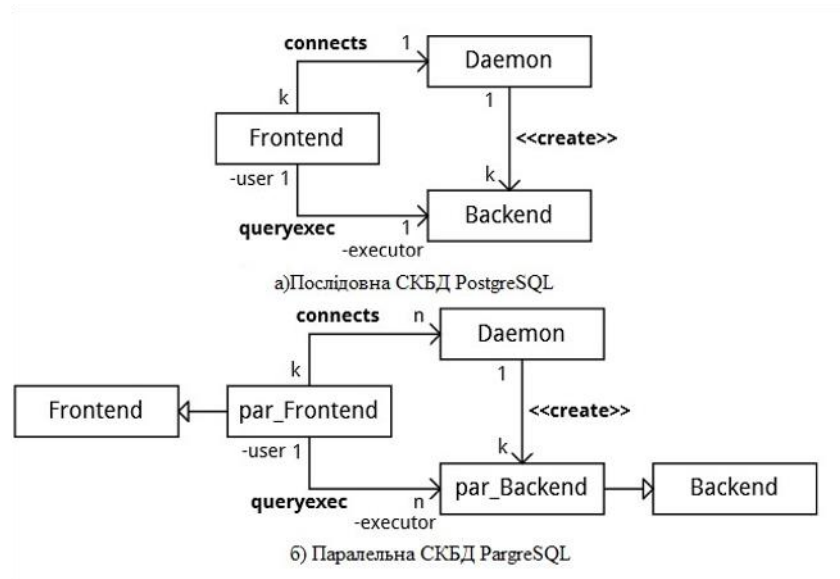
```

Системна архітектура інтеграції інтелектуального аналізу даних в СКБД PostgreSQL



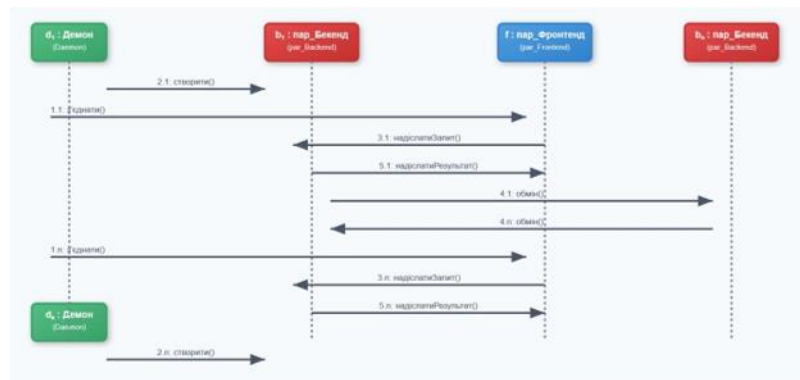
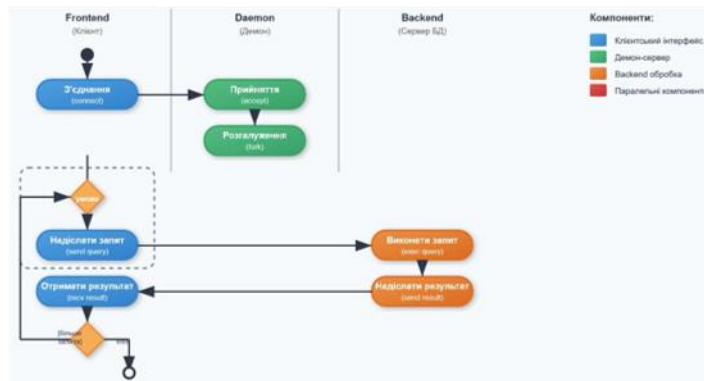
Архітектура паралельної СКБД на базі PostgreSQL. 13

Процеси СКБД

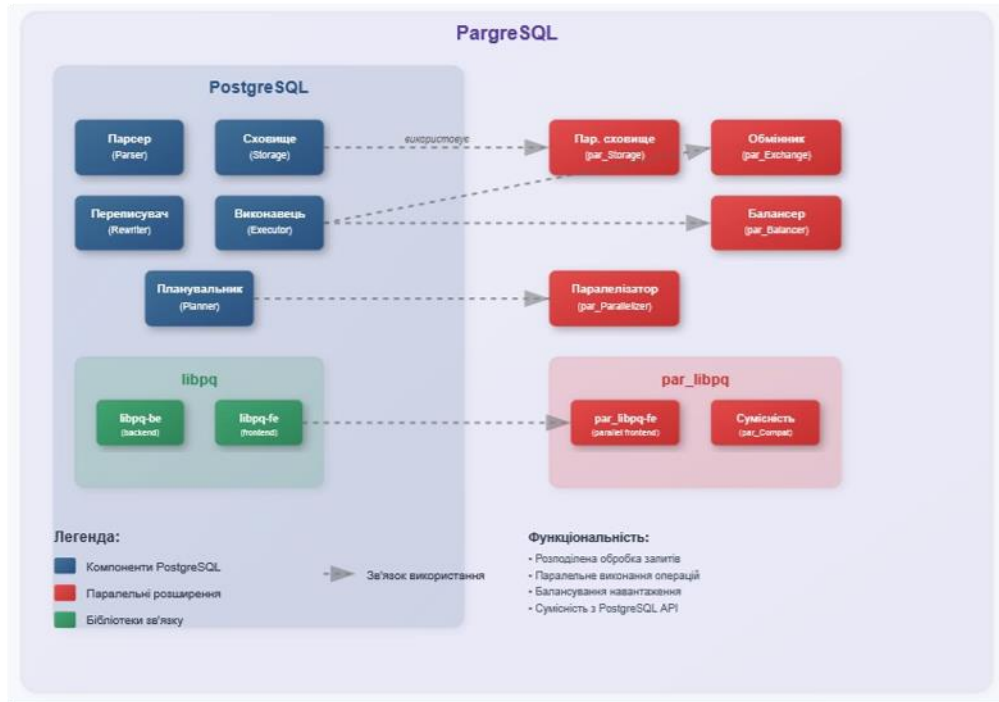


Архітектура паралельної СКБД на базі PostgreSQL. 14

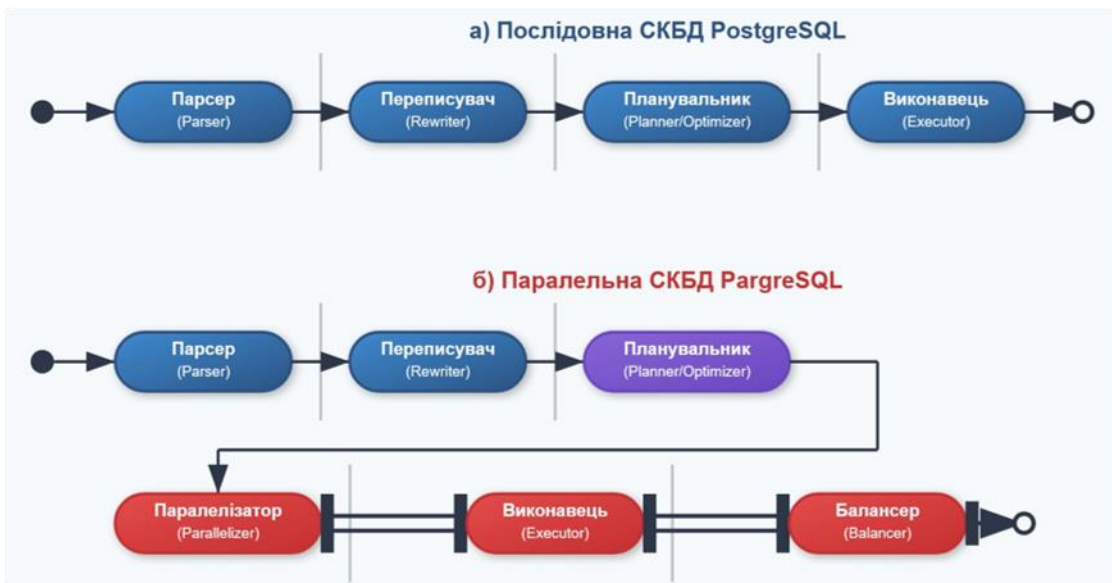
Клієнт-серверна взаємодія



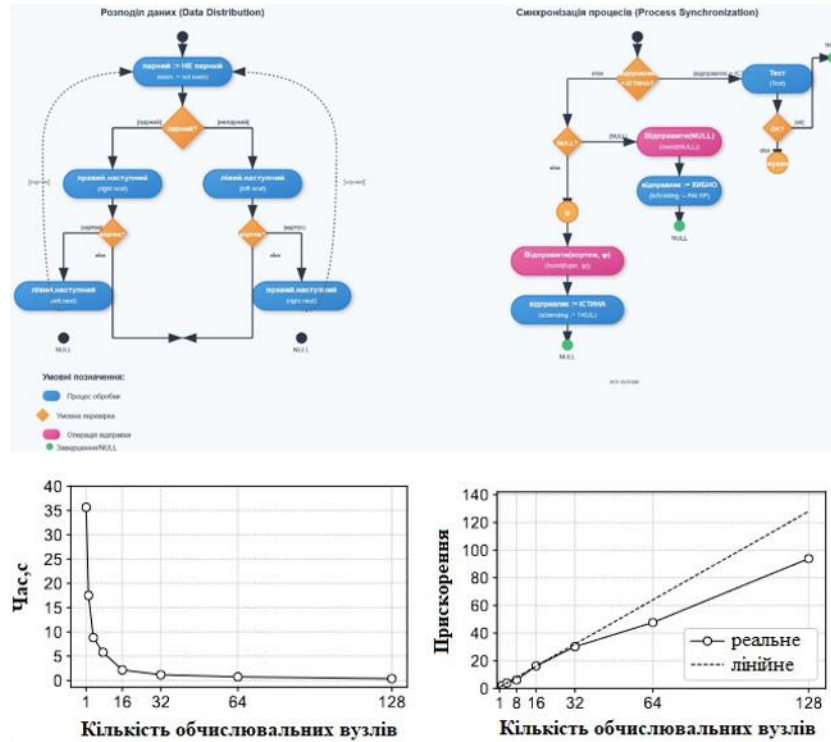
Архітектура СКБД на PostgreSQL



Архітектура паралельної СКБД на базі PostgreSQL. Обробка запиту



Реалізація деяких методів та результати



Висновки

У дослідженні здійснено глибокий аналіз інтеграції методів інтелектуального аналізу даних у середовище реляційних систем управління базами даних. Розглянуто концептуальні засади та практичні підходи до створення паралельних алгоритмів, орієнтованих на виконання у кластерних обчислювальних середовищах, що базуються на сучасних багатоядерних прискорювачах. Особливу увагу приділено архітектурним рішенням і принципам побудови інтегрованої платформи, реалізованої на основі відкритої СКБД PostgreSQL з використанням апаратного забезпечення типу Intel Many Integrated Core.