

Хмельницький національний університет
Факультет інформаційних технологій
Кафедра комп'ютерної інженерії та інформаційних систем

КВАЛІФІКАЦІЙНА РОБОТА МАГІСТРА

Галузь знань 12 – Інформаційні технології

Спеціальність 123 – Комп'ютерна інженерія

на тему «Метод розподілу навантаження в мультикомп'ютерних системах загального призначення згідно кластеризації вузлів за рівнем завантаженості ресурсів»

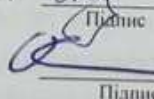
КвРКІП. 2301149.23.01.11 ПЗ

Виконала: студентка 2 курсу, група КІ2м-23-1


Підпис

Каріна ГОБА
Ім'я, прізвище

Керівник д-р. техн. наук, професор
Науковий ступінь, вчене звання


Підпис

Олег САВЕНКО
Ім'я, прізвище

До захисту допускаю:

Зав. кафедри КІС, доктор філософії, доцент

Ольгам ПАВЛОВА

12 05 2025 р.

Хмельницький, 2025

ХМЕЛЬНИЦЬКИЙ НАЦІОНАЛЬНИЙ УНІВЕРСИТЕТ

Факультет ІНФОРМАЦІЙНИХ ТЕХНОЛОГІЙ

Кафедра КОМП'ЮТЕРНОЇ ІНЖЕНЕРІЇ ТА ІНФОРМАЦІЙНИХ СИСТЕМ

Освітній рівень МАГІСТР

Галузь знань 12 ІНФОРМАЦІЙНІ ТЕХНОЛОГІЇ

Спеціальність 123 КОМП'ЮТЕРНА ІНЖЕНЕРІЯ

Освітня програма ОСВІТНЬО-НАУКОВА ПРОГРАМА «КОМП'ЮТЕРНА ІНЖЕНЕРІЯ ТА ПРОГРАМУВАННЯ»

ЗАТВЕРДЖУЮ

Зав. кафедри Ольга ПАВЛОВА

“ 01 ” 09 2024 р.

**ЗАВДАННЯ
НА КВАЛІФІКАЦІЙНУ РОБОТУ МАГІСТРА**

Каріні ГОБІ

Прізвище, ім'я, по батькові студента

1. Тема проекту (роботи) Метод розподілу навантаження в мультикомп'ютерних системах загального призначення згідно кластеризації вузлів за рівнем завантаженості ресурсів

Керівник проекту (роботи) Олег САВЕНКО д.т.н., професор

Прізвище, ім'я, по батькові, науковий ступінь, вчене звання

Затверджена наказом ректора університету від 08.01.2025 №8

2. Строк подання студентом проекту (роботи) на кафедру 01.05.2025 р.

3. Вихідні дані до проекту (роботи) Завдання на дипломне проектування

4. Зміст пояснювальної записки (перелік питань, які потрібно розробити)

Теоретичні основи розподілу навантаження в мультикомп'ютерних системах

Аналіз існуючих методів розподілу навантаження в мультикомп'ютерних системах

Система забезпечення розподілу навантаження згідно кластеризації вузлів

Оцінка ефективності методу розподілу навантаження в мультикомп'ютерних системах

5. Перелік графічного матеріалу (із зазначенням обов'язкових креслень)

6. Консультанти розділів кваліфікаційної роботи магістра

Розділ	Прізвище, ініціали та посада консультанта	Підпис, дата	
		завдання видав	завдання прийняв
Нормоконтроль	Сергій ЛИСЕНКО, професор кафедри КПС		
Антиплагіат	Андрій НіЧЕПОРУК, доцент кафедри КПС		

7. Дата видачі завдання « 01 » 09 2024р.

КАЛЕНДАРНИЙ ПЛАН

№з/п	Назва етапів (розділів) кваліфікаційної роботи магістра	Термін виконання етапів проекту (роботи)	Примітка
1	Вибір напрямку дослідження та узгодження тематики КвРМ з керівником	01.09.2024	виконано
2	Ознайомлення з предметною областю; формулювання мети та задач дослідження; визначення об'єкта та предмета дослідження	01.10.2024	виконано
3	Робота над розділом 1 – аналіз відомих моделей, методів за темою; постановка задачі	01.11.2024	виконано
4	Робота над розділом 2 – розробка моделей для вирішення поставленої задачі	01.12.2024	виконано
5	Робота над науковою статтею	01.02.2025	виконано
6	Робота над розділом 3 – розробка методів для вирішення поставленої задачі	15.02.2025	виконано
7	Робота над розділом 4 – проектування та розробка ПЗ для вирішення поставленої задачі, експериментальна частина	01.04.2025	виконано
8	Оформлення пояснювальної записки згідно вимог	18.04.2025	виконано
9	Попередній захист ДРМ	29.04.2025	виконано
10	Захист ДРМ на засіданні ЕК	До 31.05.2025	

Студент

Керівник роботи


Підпис

Підпис

Каріна ГОБА
Ім'я, прізвище
Олег САВЕНКО
Ім'я, прізвище

РЕФЕРАТ

Тема кваліфікаційної роботи магістра: Метод розподілу навантаження в мультимедійних системах загального призначення згідно кластеризації вузлів за рівнем завантаженості ресурсів

Автор роботи: Каріна ГОБА

Керівник роботи: Олег САВЕНКО

Пояснювальна записка: 72 с., 16 рис., 6 табл., 3 дод., 80 джерел.

РОЗПОДІЛ НАВАНТАЖЕННЯ, МУЛЬТИКОМП'ЮТЕРНІ СИСТЕМИ, КЛАСТЕРИЗАЦІЯ ВУЗЛІВ, ЗАВАНТАЖЕНІСТЬ РЕСУРСІВ, АЛГОРИТМИ РОЗПОДІЛУ НАВАНТАЖЕННЯ, ВУЗЛИ СИСТЕМИ, РЕСУРСНИЙ МОНІТОРИНГ, ОПТИМІЗАЦІЯ РОЗПОДІЛУ

Об'єктом дослідження є мультимедійні системи загального призначення, зокрема процеси розподілу навантаження між вузлами в таких системах.

Предметом дослідження є методи та алгоритми розподілу навантаження в мультимедійних системах, засновані на кластеризації вузлів за рівнем завантаженості ресурсів.

Метою кваліфікаційної роботи є забезпечення ефективності розподілу навантаження в мультимедійних системах загального призначення з використанням кластеризації вузлів за рівнем їхнього завантаження.

Для розв'язання поставлених задач використовувалися: аналіз існуючих методів балансування навантаження в мультимедійних системах; розробка алгоритму кластеризації вузлів на основі їхнього поточного завантаження; оцінка ефективності розподілу ресурсів при використанні запропонованого методу; моделювання та експериментальне тестування алгоритму на реальних або емуляційних даних.

Наукова новизна дослідження полягає у розробці нового методу балансування навантаження, який на відміну від відомих використовує кластеризацію вузлів за рівнем їхнього завантаження, що дозволяє знизити

енергоспоживання, підвищити рівномірність розподілу задач та зменшити затримки в мультикомп'ютерних системах.

На основі проведених досліджень розроблена архітектура і компоненти програмного забезпечення, що включає модулі для кластеризації вузлів, моніторингу завантаженості ресурсів, а також алгоритм розподілу навантаження з урахуванням результатів кластеризації.

Практична значущість - запропонований метод може бути застосований для оптимізації обчислювальних процесів у хмарних сервісах, дата-центрах, суперкомп'ютерах та корпоративних мережах, що дозволить значно покращити ефективність використання апаратних ресурсів.

У вступі подано об'єкт та предмет дослідження, мету, наукову новизну та практичну цінність роботи, а також характеристику структури роботи.

У першому розділі проведено аналіз існуючих методів балансування навантаження в мультикомп'ютерних системах.

У другому розділі здійснено розробку методу кластеризації вузлів на основі їхнього поточного завантаження.

У третьому розділі провели оцінку ефективності розподілу ресурсів при використанні запропонованого методу.

У четвертому розділі змодельювали та провели експериментальне тестування алгоритму на реальних або емуляційних даних.

У висновку підведено підсумки досягнення та ефективності методу та розв'язання задач дослідження

ЗМІСТ

СКОРОЧЕННЯ ТА УМОВНІ ПОЗНАКИ	5
ВСТУП	6
1 ТЕОРЕТИЧНІ ОСНОВИ РОЗПОДІЛУ НАВАНТАЖЕННЯ В МУЛЬТИКОМП'ЮТЕРНИХ СИСТЕМАХ	9
1.1 Поняття мультикомп'ютерних систем загального призначення	Error! Bookmark not defined.1.3
1.2	Error! Bookmark not defined.1.4
1.3	Error! Bookmark not defined.1.4
1.4	Error! Bookmark not defined.1.4
2 МОДЕЛЬ РОЗПОДІЛУ НАВАНТАЖЕННЯ В МУЛЬТИКОМП'ЮТЕРНИХ СИСТЕМАХ	24
2.1	Error! Bookmark not defined.2.2
2.2	Error! Bookmark not defined.2.3
2.3	Error! Bookmark not defined.2.4
2.4	Error! Bookmark not defined.2.5
2.5	Error! Bookmark not defined.2.5
3 СИСТЕМА ЗАБЕЗПЕЧЕННЯ РОЗПОДІЛУ НАВАНТАЖЕННЯ ЗГІДНО КЛАСТЕРИЗАЦІЇ ВУЗЛІВ	40
3.1	Error! Bookmark not defined.3.1.1
3.1.1 Обчислення навантаження за допомогою CNN та RNN	40
3.2	Error! Bookmark not defined.3.2.1
3.2.1 Оптимізована кластеризація вузлів	45
3.2.2 Кластеризація на основі алгоритму оптимізації Hybrid Lyrebird Falcon	48
3.3	Error! Bookmark not defined.3.4
3.4	Висновок до третього розділу
	54
4 ОЦІНКА ЕФЕКТИВНОСТІ МЕТОДУ РОЗПОДІЛУ НАВАНТАЖЕННЯ В МУЛЬТИКОМП'ЮТЕРНИХ СИСТЕМАХ	56
4.1 Визначення основних характеристик та експериментальні дослідження	56

	4
4.2 Архітектура програмної реалізації інформаційної технології	58
4.3 Експериментальне дослідження ефективності методу розподілу навантаження в умовах інтенсивної роботи розумного будинку	Error! Bookmark not defined.
4.5 Висновок до четвертого розділа	77
ВИСНОВКИ	78
ПЕРЕЛІК ДЖЕРЕЛ ПОСИЛАНЬ	79
ДОДАТОК А Лістинг програмного забезпечення	86
ДОДАТОК Б Наукова праця здобувача	88
ДОДАТОК В Презентація	96

СКОРОЧЕННЯ ТА УМОВНІ ПОЗНАКИ

LAN - Local Area Network (локальна мережа)

IaaS - Infrastructure as a Service (інфраструктура як послуга)

PaaS - Platform as a Service (платформа як послуга)

SaaS - Software as a Service (програмне забезпечення як послуга)

AWS - Amazon Web Services

IoT - Internet of Things (Інтернет речей)

NB-IoT - Narrowband Internet of Things (вузькосмуговий інтернет речей)

LTE-M - Long Term Evolution for Machines (мережа мобільного зв'язку для машин)

RED - Random Early Detection (випадкове раннє виявлення перевантаження)

DBSCAN - Density-Based Spatial Clustering of Applications with Noise
(кластеризація на основі щільності з урахуванням шумів)

CPU - Central Processing Unit (центральний процесор)

RAM - Random Access Memory (оперативна пам'ять)

ВСТУП

Сучасні мультимедійні системи загального призначення активно використовуються в багатьох сферах, включаючи високопродуктивні обчислення, обробку великих обсягів даних, хмарні обчислення та розподілені інформаційні системи. Одним із ключових викликів таких систем є ефективний розподіл навантаження, що безпосередньо впливає на продуктивність, швидкодію та енергоспоживання комп'ютерних кластерів.

Одним із перспективних підходів до оптимізації роботи мультимедійних систем є кластеризація вузлів за рівнем їхнього завантаження. Такий підхід дозволяє групувати обчислювальні ресурси відповідно до їх поточного стану, що дає змогу підвищити ефективність розподілу завдань, мінімізувати затримки та зменшити перевантаження окремих компонентів системи.

Актуальність дослідження є розподіл навантаження є критичною задачею для розподілених систем, оскільки неправильне управління ресурсами може призвести до неоптимального використання процесорних ядер, збільшення часу виконання задач та підвищеного енергоспоживання. Більшість сучасних методів балансування навантаження використовують статичні або динамічні алгоритми, однак вони не завжди враховують реальний стан окремих вузлів. Тому застосування кластеризації вузлів за рівнем завантаженості дозволяє адаптивно керувати ресурсами та підвищувати ефективність мультимедійних систем.

Метою даної роботи є аналіз та розробка методу розподілу навантаження в мультимедійних системах загального призначення з використанням кластеризації вузлів за рівнем їхнього завантаження.

Для досягнення поставленої мети необхідно вирішити такі завдання:

1. Аналіз існуючих методів балансування навантаження в мультимедійних системах.
2. Розробка методу кластеризації вузлів на основі їхнього поточного завантаження.

3. Оцінка ефективності розподілу ресурсів при використанні запропонованого методу.

4. Моделювання та експериментальне тестування алгоритму на реальних або емуляційних даних.

Об'єктом дослідження є процеси розподілу обчислювальних задач у мультикомп'ютерних системах загального призначення.

Предметом дослідження є методи балансування навантаження, які базуються на кластеризації вузлів за рівнем їхнього завантаження.

Наукова новизна дослідження полягає у розробці нового методу балансування навантаження, який на відміну від відомих використовує кластеризацію вузлів за рівнем їхнього завантаження, що дозволяє знизити енергоспоживання, підвищити рівномірність розподілу задач та зменшити затримки в мультикомп'ютерних системах.

На основі проведених досліджень розроблена архітектура і компоненти програмного забезпечення, що включає модулі для кластеризації вузлів, моніторингу завантаженості ресурсів, а також алгоритм розподілу навантаження з урахуванням результатів кластеризації.

Практична значущість - запропонований метод може бути застосований для оптимізації обчислювальних процесів у хмарних сервісах, дата-центрах, суперкомп'ютерах та корпоративних мережах, що дозволить значно покращити ефективність використання апаратних ресурсів.

Для розв'язання поставлених задач використовувалися методи кластеризації вузлів за рівнем завантаженості ресурсів, що дозволяє оптимально розподіляти навантаження між різними частинами системи на основі їх поточного стану, алгоритмів розподілу навантаження, спрямованих на мінімізацію затримок і максимізацію продуктивності системи шляхом ефективного використання наявних ресурсів, математичного моделювання, для оцінки ефективності запропонованих методів на основі симуляцій різних сценаріїв роботи мультикомп'ютерних систем, аналізу та оптимізації ресурсів, що включає моніторинг завантаженості вузлів і динамічне перенаправлення навантаження для підтримки збалансованості системи,

статистичних методів, для оцінки продуктивності та ефективності роботи системи в умовах різних рівнів навантаження.

За темою кваліфікаційної роботи опубліковано одну статтю в науковому журналі Гоба К., Жуковський П., Ковальчук В., Клейн О. (2025). МЕТОД РОЗПОДІЛУ НАВАНТАЖЕННЯ В МУЛЬТИКОМП'ЮТЕРНИХ СИСТЕМАХ ЗАГАЛЬНОГО ПРИЗНАЧЕННЯ ЗГІДНО КЛАСТЕРИЗАЦІЇ ВУЗЛІВ ЗА РІВНЕМ ЗАВАНТАЖЕНОСТІ РЕСУРСІВ. MEASURING AND COMPUTING DEVICES IN TECHNOLOGICAL PROCESSES, (2).

1 ТЕОРЕТИЧНІ ОСНОВИ РОЗПОДІЛУ НАВАНТАЖЕННЯ В МУЛЬТИКОМП'ЮТЕРНИХ СИСТЕМАХ

1.1 Поняття мультимікомп'ютерних систем загального призначення

Мультимікомп'ютерні системи – це обчислювальні системи, які складаються з кількох автономних вузлів (комп'ютерів), об'єднаних у єдину інфраструктуру для виконання спільних завдань. Їх головною особливістю є розподілений характер обчислень, що дозволяє виконувати задачі паралельно на різних вузлах, оптимізуючи час і використання ресурсів [1].

Основні характеристики мультимікомп'ютерних систем:

Архітектура системи - мультимікомп'ютерні системи можуть бути побудовані на основі однорідних (з однаковими ресурсами) або гетерогенних (різних за потужністю) вузлів. Зазвичай такі системи використовують мережу передачі даних для взаємодії між вузлами [2].

Масштабованість - важливою перевагою таких систем є можливість збільшення кількості вузлів для підвищення продуктивності без істотних змін в архітектурі [3].

Призначення - мультимікомп'ютерні системи загального призначення створюються для вирішення широкого спектра завдань, таких як наукові обчислення, обробка великих обсягів даних, виконання хмарних сервісів тощо [4].

Кластерні системи — це обчислювальні комплекси, вузли яких з'єднані через локальну мережу (LAN). Вони здебільшого використовуються для вирішення завдань, що вимагають значних ресурсів і високої обчислювальної потужності, таких як високопродуктивні обчислення та складні наукові розрахунки [5]. Хмарні обчислення, у свою чергу, забезпечують взаємодію вузлів через Інтернет, надаючи доступ до ресурсів за моделями IaaS (інфраструктура як послуга), PaaS (платформа як послуга) та SaaS (програмне забезпечення як послуга). Прикладом таких систем є платформи Amazon Web Services (AWS) та Google Cloud Platform [6]. Грід-системи об'єднують вузли, розташовані в різних географічних зонах, для виконання великих обчислювальних завдань, але характеризуються низькою інтеграцією між

вузлами [7]. Суперкомп'ютери ж складаються з великої кількості вузлів, оптимізованих для виконання високопродуктивних обчислень, таких як моделювання кліматичних процесів або обробка космічних даних [8].

Серед прикладів мультикомп'ютерних систем можна виділити кластери для обробки великих даних, такі як Hadoop або Spark, хмарні платформи, зокрема Amazon Web Services і Google Cloud Platform, а також суперкомп'ютери, наприклад, Titan чи Fugaku, які використовуються для наукових досліджень [9,10].

Hadoop є платформою для розподіленої обробки великих обсягів даних. Вона базується на моделі MapReduce і дозволяє розподіляти завдання між численними вузлами системи, забезпечуючи паралельну обробку даних. Основні сфери використання включають аналіз даних, обробку журналів подій, побудову рекомендаційних систем та аналіз соціальних мереж [11].

Apache Spark - це розподілена обчислювальна платформа для швидкої обробки даних. Вона підтримує обробку великих обсягів даних у пам'яті, що значно прискорює виконання завдань у порівнянні з Hadoop. Spark використовується для машинного навчання, обробки потокових даних, аналізу великих масивів даних та побудови моделей прогнозування [12].

AWS - це хмарна платформа, яка надає інфраструктуру, платформи та програмне забезпечення як послугу (IaaS, PaaS, SaaS). Вона забезпечує ресурси для обчислень, зберігання даних, машинного навчання, аналітики, а також створення масштабованих і високодоступних систем. AWS широко використовується для розробки веб-додатків, хостингу баз даних та побудови хмарних сервісів [13].

Titan - суперкомп'ютер, який використовується для наукових досліджень, таких як моделювання кліматичних змін, обробка великих обсягів даних у геноміці, моделювання поведінки матеріалів і аналіз фізичних явищ. Завдяки високій продуктивності він дозволяє виконувати складні обчислення з використанням паралельних алгоритмів [14, 15].

Fugaku - один із найпотужніших суперкомп'ютерів у світі, який застосовується для розробки ліків, симуляцій кліматичних процесів, моделювання землетрусів, досліджень у сфері фізики і машинного навчання. Його можливості

дозволяють вирішувати завдання, які потребують обробки екстремально великих обсягів даних [16].

Розподіл навантаження є критично важливим аспектом функціонування мультикомп'ютерних систем. Від ефективності використання обчислювальних ресурсів залежить продуктивність системи, її масштабованість та стабільність. Залежно від принципів організації процесу, методи розподілу навантаження поділяються на статичні та динамічні [17].

Статичні методи передбачають попередньо визначений розподіл завдань між вузлами системи, що не змінюється під час виконання задач. Розподіл базується на фіксованих правилах, таких як черговість чи випадковість. Ці методи не враховують реального стану системи, але забезпечують простоту реалізації [18].

Round-Robin (циклічний розподіл): завдання послідовно розподіляються між вузлами за круговою схемою, кожен вузол отримує чергове завдання незалежно від його поточного стану [19, 20]. Наприклад, у системі з трьома вузлами перше завдання надходить на вузол 1, друге — на вузол 2, третє — на вузол 3, а четверте повертається на вузол 1 [21]. Переваги: рівномірний розподіл завдань між вузлами. Недоліки: не враховує поточного стану вузлів, що може призводити до перевантаження.

Random Assignment (випадковий розподіл): завдання призначаються вузлам випадковим чином. Цей метод є простим у реалізації та забезпечує рівномірність у великій кількості завдань [22]. Переваги: простота реалізації. Недоліки: можливість нерівномірного розподілу в малих вибірках.

Переваги статичних методів - простота реалізації та мінімальні витрати на обчислення, низькі вимоги до інфраструктури, відсутність необхідності моніторингу системи в реальному часі [23].

Недоліки статичних методів - неадаптивність до змін у системі, таких як збільшення навантаження або вихід вузлів з ладу, імовірність нерівномірного розподілу навантаження, що може призвести до перевантаження окремих вузлів [24].

Динамічні методи розподілу навантаження враховують поточний стан системи. Завдання розподіляються на основі моніторингу завантаженості вузлів у реальному часі. Це дозволяє адаптувати розподіл до змін, забезпечуючи більш ефективне використання ресурсів [25].

Роль моніторингу - моніторинг завантаженості вузлів у реальному часі є ключовим елементом динамічних методів. Він забезпечує збір інформації про поточний стан вузлів, включаючи використання процесора, пам'яті, мережевих ресурсів тощо [26].

Least Connection (найменша кількість з'єднань): завдання направляється на вузол із найменшою кількістю активних з'єднань. Переваги: балансування навантаження за кількістю одночасних задач. Недоліки: підвищені вимоги до моніторингу [27].

Load-Aware (завантаження-орієнтоване): завдання направляються на вузол із найменшим рівнем завантаженості. Переваги: ефективне використання ресурсів. Недоліки: потребує високоточного моніторингу системи [28].

Adaptive Balancing (адаптивне балансування): метод динамічно коригує стратегію розподілу на основі змін у системі. Переваги: максимальна гнучкість і адаптивність. Недоліки: високі витрати на обчислення [29].

Переваги динамічних методів - адаптивність до змін у системі, ефективне використання ресурсів і мінімізація ризиків перевантаження.

Недоліки динамічних методів - додаткові витрати на обчислення та передачу даних для моніторингу і управління, підвищені вимоги до мережевої інфраструктури [30, 31].

Для оцінки ефективності методів розподілу навантаження проведемо їх порівняння за основними критеріями (таблиця 1.1) [32].

Комбіновані методи - поєднують елементи статичних і динамічних підходів. Наприклад, на початковому етапі завдання можуть бути розподілені за статичною схемою, а в процесі виконання використовуються динамічні механізми для оптимізації. Це дозволяє досягти компромісу між швидкістю, гнучкістю і ресурсозатратністю [33].

Таблиця 1.1 – Порівняльний аналіз статичних і динамічних методів

Критерії	Статичні методи	Динамічні методи
Швидкість	Висока (немає додаткових обчислень)	Середня (потребує моніторингу)
Гнучкість	Низька	Висока
Ресурсозатратність	Низька	Висока
Стійкість до змін	Низька	Висока

Перевантаження мережі може призвести до збоїв у роботі, втрати даних або навіть до кібератак. Для забезпечення стабільної роботи IoT-інфраструктури використовуються такі ключові підходи [34]:

Горизонтальне масштабування баз даних (рисунок 1.1). Оскільки IoT-пристрої генерують великі обсяги даних, традиційні централізовані бази даних можуть не справлятися із навантаженням. Горизонтальне масштабування передбачає додавання нових серверів, що розподіляють зберігання та обробку даних між собою [35].

Основні переваги підходу: балансування навантаження – навантаження рівномірно розподіляється між кількома серверами, висока доступність – навіть якщо один сервер виходить з ладу, інші продовжують працювати. Швидкість обробки даних – розподілені бази даних дозволяють швидше виконувати запити [36].

Оскільки IoT-пристрої генерують великі обсяги даних, традиційні централізовані бази даних можуть не справлятися із навантаженням. Горизонтальне масштабування передбачає додавання нових серверів, що розподіляють зберігання та обробку даних між собою [37].

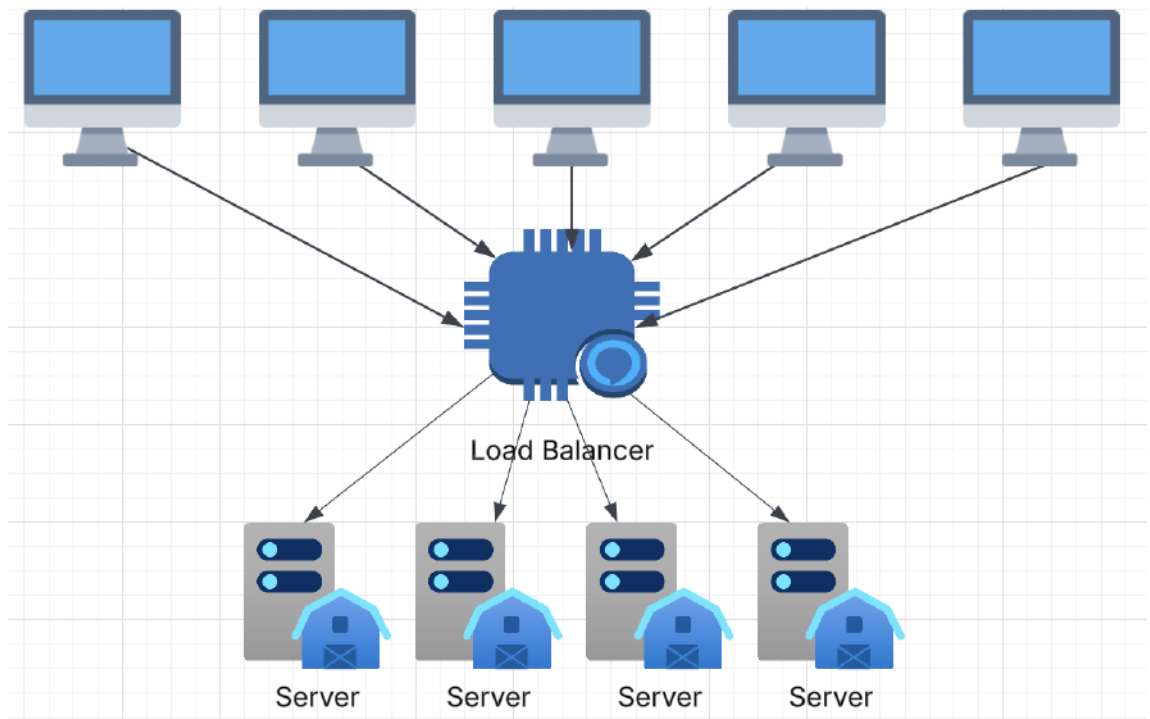


Рисунок 1.1 - Горизонтальне масштабування

Основні переваги підходу: балансування навантаження – навантаження рівномірно розподіляється між кількома серверами, висока доступність – навіть якщо один сервер виходить з ладу, інші продовжують працювати. Швидкість обробки даних – розподілені бази даних дозволяють швидше виконувати запити [38].

Блокчейн-технологія забезпечує децентралізацію управління мережею IoT, що значно підвищує безпеку та відмовостійкість системи. У традиційних IoT-мережах центральний сервер є критичною точкою, яка може зазнати кібератаки або відмови [39].

Захист від атак – кожен вузол зберігає копію транзакцій, що унеможливорює підробку даних. Децентралізація – немає єдиної точки відмови, система працює навіть у разі виходу з ладу деяких пристроїв. Прозорість – усі зміни в даних зберігаються у блокчейні, що унеможливорює несанкціоновані зміни [40].

ІОТА Tangle – розподілена реєстрова технологія, що використовується для IoT-пристроїв, яка усуває проблему високих енергетичних витрат та затримок, що притаманні традиційним блокчейн-системам.

Використання стільникового IoT для безперебійного підключення (рисунок 1.2)

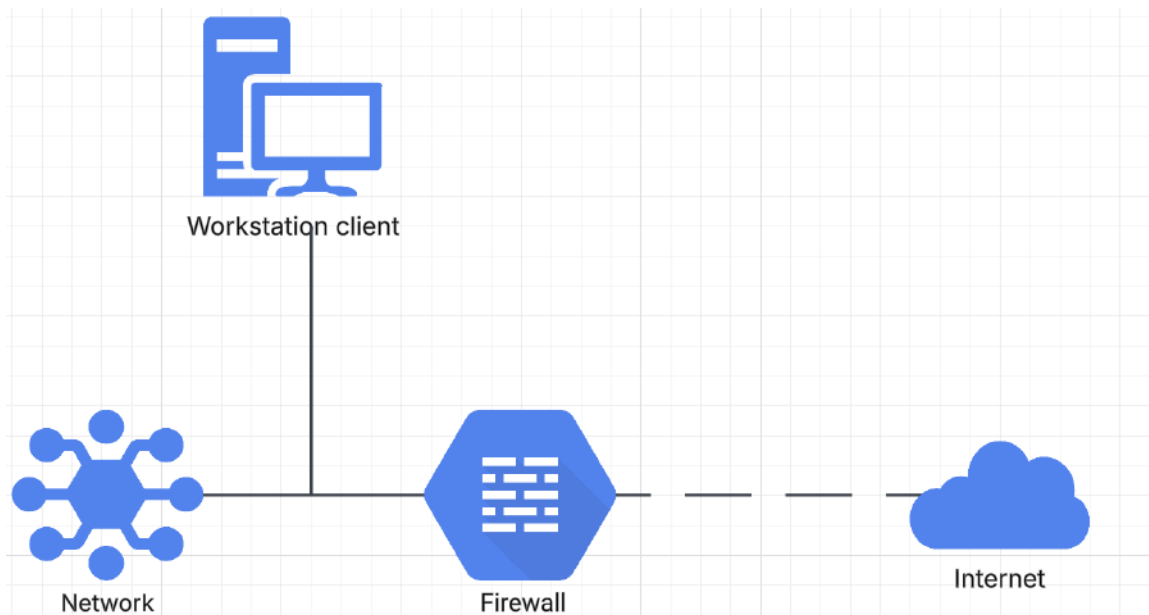


Рисунок 1.2 - Стільниковий IoT для безперебійного підключення

Стільниковий IoT (LTE-M, NB-IoT) забезпечує стабільне підключення пристроїв IoT навіть у важкодоступних зонах. Це особливо важливо для розподілу навантаження в мобільних мережах та динамічних системах [41].

Технологія NB-IoT (Narrowband IoT) використовується для підключення "розумних" лічильників, датчиків у сільському господарстві та інших пристроїв, які передають невеликі обсяги даних на великі відстані [42].

Однією з основних загроз у мережах IoT є аномальна поведінка трафіку, спричинена як збоєм пристрою, так і кібератаками. Щоб уникнути перевантаження мережі, застосовуються методи аналізу аномалій.

Машинне навчання – нейромережі аналізують поведінку пристроїв та виявляють аномалії [43].

Метод порогових значень – система визначає критичні рівні навантаження та реагує на їх перевищення. Використання евристичних методів – аналіз шаблонів поведінки трафіку для виявлення підозрілих дій. Random Early Detection (RED) –

алгоритм виявлення перевантажень, що відстежує напливи трафіку та запобігає їхньому впливу на всю систему [44].

Загальний підхід до забезпечення розподілу навантаження у мережах Інтернету речей (IoT) базується на комбінації різних методів, які забезпечують ефективне використання ресурсів та безперебійну роботу пристроїв. Одним із ключових рішень є застосування децентралізованих алгоритмів балансування, коли кожен вузол мережі самостійно оцінює рівень свого завантаження та взаємодіє з іншими пристроями для рівномірного розподілу обчислювальних завдань [45].

1.2 Кластеризація вузлів у мультикомп'ютерних системах

Кластеризація вузлів є важливим підходом у мультикомп'ютерних системах, який забезпечує групування вузлів за певними критеріями. Це дозволяє підвищити ефективність управління ресурсами та оптимізувати процеси розподілу навантаження [46].

Кластеризація вузлів у мультикомп'ютерних системах означає групування вузлів у межах системи відповідно до певних характеристик, таких як завантаженість ресурсів, географічна близькість чи типи виконуваних завдань. Основна мета кластеризації полягає у створенні логічно пов'язаних груп вузлів для більш ефективного використання ресурсів та підвищення продуктивності системи [47].

Для реалізації кластеризації вузлів використовуються різні критерії, що дозволяють забезпечити релевантність групування залежно від потреб системи.

Завантаженість ресурсів : вузли групуються за рівнем завантаженості обчислювальних ресурсів, таких як процесор, оперативна пам'ять або пропускну здатність мережі. Наприклад, вузли з високим рівнем завантаженості можуть бути виділені в окрему групу для перенаправлення частини задач до менш завантажених вузлів [48].

Географічна близькість вузлів: вузли, які розташовані в межах однієї географічної зони або мають низькі затримки у зв'язку між собою, можуть бути об'єднані в окремий кластер для зменшення часу передачі даних [49].

Пріоритетність задач та типи обчислень: вузли можуть групуватися залежно від типу задач, які вони виконують, наприклад, вузли, що спеціалізуються на обробці графічних даних, можуть бути виділені в окремий кластер [50].

Ієрархічні методи (рисунок 1.3) - ці методи передбачають поступове об'єднання вузлів у групи на основі схожості їхніх характеристик. Перевага цього підходу полягає у створенні багаторівневої структури кластерів, що може бути використана для детального аналізу [51].

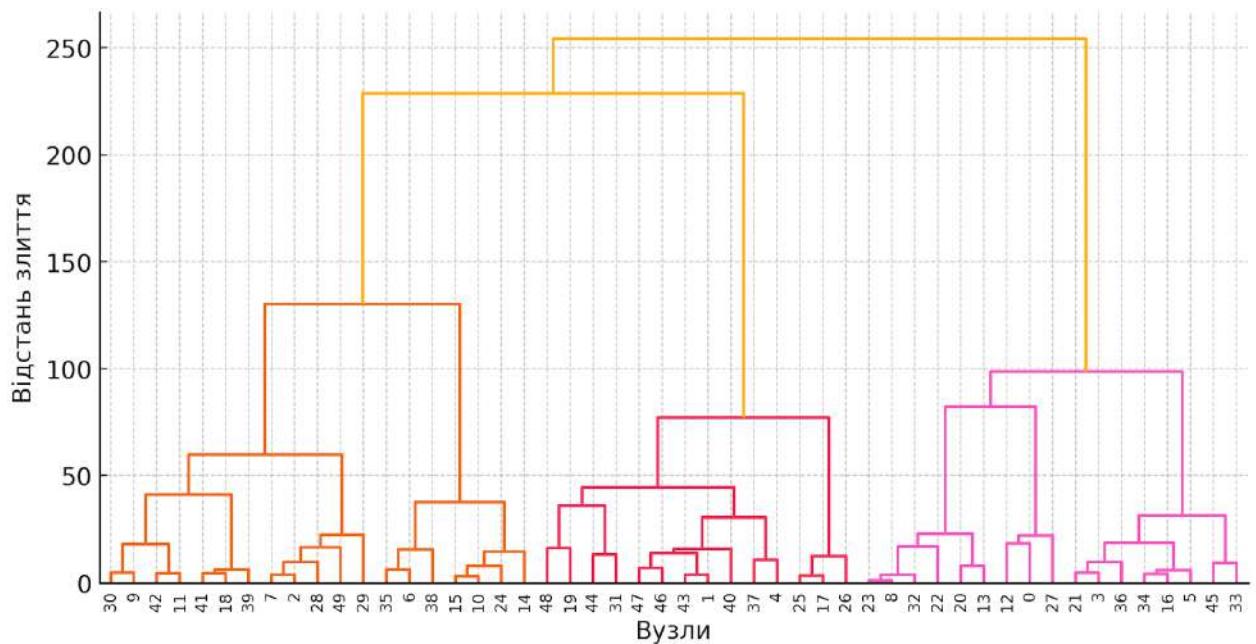


Рисунок 1.3 - Ієрархічна кластеризація вузлів

Алгоритм методу можна описати наступним чином [52, 53]:

1. Побудувати матрицю схожості, яка відображає відстань між кожною парою об'єктів у вибірці. На початковому етапі кожен об'єкт розглядається як окремий кластер.

2. Визначити пару кластерів, що є найбільш подібними відповідно до значень у матриці схожості, та об'єднати їх у новий, більший кластер. Після цього оновити матрицю схожості, враховуючи об'єднання двох кластерів.

3. Якщо всі об'єкти об'єднані в єдиний кластер, завершити роботу алгоритму. В іншому випадку, повторити процес із другого кроку.

Другий графік демонструє ієрархічну кластеризацію вузлів, яка представлена у вигляді дендрограми (дерева кластерів) [54]. Горизонтальна вісь – індекси вузлів у системі. Вертикальна вісь – відстань між вузлами, яка використовується для об'єднання в групи. Гілки дерева демонструють, як вузли об'єднуються у все більші групи, доки не сформується остаточні кластери.

Алгоритми k-means (рисунок 1.4) - це один із найпопулярніших алгоритмів кластеризації, який передбачає розподіл вузлів на задану кількість кластерів на основі мінімізації відстані між вузлами та центрами кластерів.

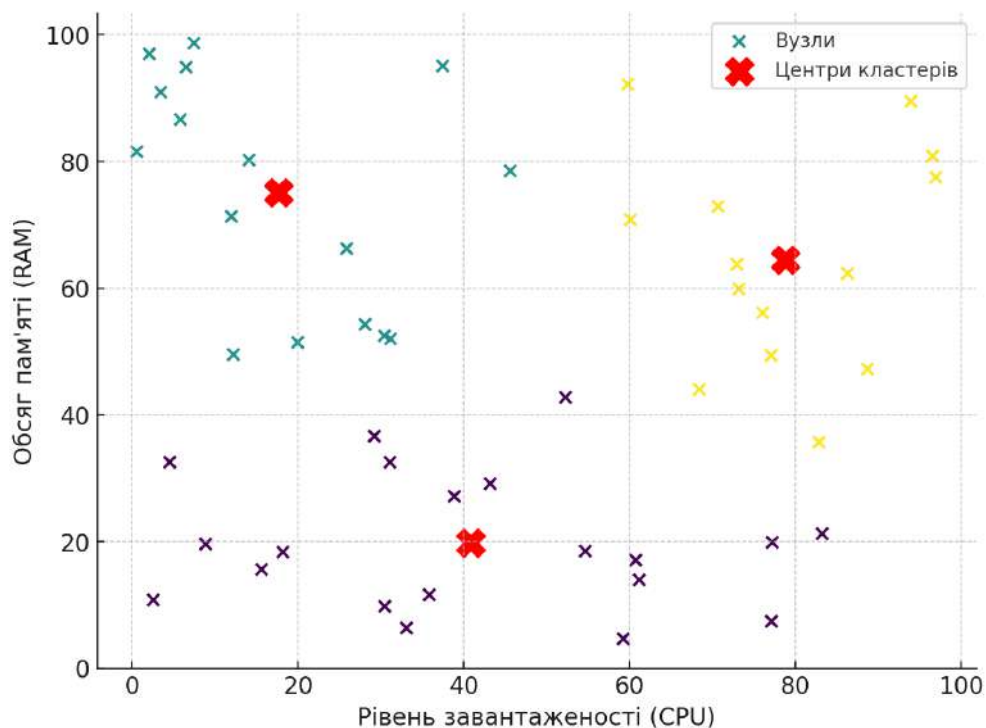


Рисунок 1.4 - Кластеризація вузлів методом K-Means

Кластеризація часто застосовується, коли потрібно більше дізнатися про структуру набору даних. [55]. У таких випадках метод k-середніх допомагає

виявити цікаві кластери, які можна дослідити окремо. Однак кластеризація не завжди є надійним автоматичним методом класифікації, оскільки її результати залежать від початкових умов і можуть бути нестабільними. Тому рекомендовано повторювати процес кілька разів для досягнення стабільного результату [56].

На графіку зображені вузли мультикомп'ютерної системи (кожна точка — окремий вузол), які згруповані у три кластери за допомогою алгоритму K-Means. Кольори точок представляють різні кластери, до яких належать вузли залежно від їхнього рівня завантаження та доступної пам'яті. Червоні "X" позначають центри кластерів, які є усередненими значеннями характеристик вузлів у кожній групі [57]. Вісі графіка: X-вісь – рівень завантаженості процесора (CPU), Y-вісь – обсяг доступної оперативної пам'яті (RAM). DBSCAN (Density-Based Spatial Clustering of Applications with Noise) - метод, який групує вузли на основі щільності їх розташування, що дозволяє ефективно обробляти дані зі шумом [58].

Fuzzy Clustering [59] - алгоритм, який дозволяє вузлам належати до кількох кластерів із різними ступенями приналежності, що є корисним для систем із динамічно змінними характеристиками.

На першому кроці для кожної точки обчислюється відстань до центрів кластерів, зазвичай використовуючи евклідову метрику. Виходячи з цих відстаней, визначаються значення матриці U , яка показує, наскільки сильно кожна точка належить до певного кластера. Далі центри кластерів оновлюються, використовуючи середньозважене значення координат точок, враховуючи їхні ступені приналежності до кластерів [60].

Цей процес повторюється ітеративно: обчислюються відстані, оновлюється матриця U , коригуються центри кластерів. Алгоритм завершується тоді, коли зміни в значеннях матриці U або координатах центрів кластерів стають меншими за заданий поріг або досягається максимальна кількість ітерацій [61].

У результаті кожна точка має розподілений ступінь приналежності до кластерів, що дозволяє враховувати розмитість меж між групами даних.

Роль кластеризації у розподілі навантаження - кластеризація вузлів відіграє ключову роль у підвищенні ефективності розподілу навантаження у мультикомп'ютерних системах [62].

Кластеризація дозволяє виявити вузли з високим і низьким рівнем завантаженості (так звані "гарячі" і "холодні" зони), що допомагає оптимізувати розподіл завдань.

Завдяки групуванню вузлів за рівнем завантаженості або іншими характеристиками стає можливим більш ефективно перенаправлення задач між вузлами, що мінімізує ризики перевантаження окремих вузлів.

Кластеризація вузлів знаходить застосування у різних типах систем [63,64].

У хмарних обчисленнях для групування вузлів за географічною близькістю з метою мінімізації затримок.

У кластерних системах для рівномірного розподілу завдань між вузлами з урахуванням їх поточного стану.

У ґрид-системах для об'єднання вузлів, що спеціалізуються на певних типах задач, наприклад, обробка великих даних.

Балансування навантаження є однією з основних задач у мультикомп'ютерних системах, яка безпосередньо впливає на продуктивність і надійність їхньої роботи [65]. Суть проблеми полягає в забезпеченні рівномірного розподілу обчислювальних ресурсів між вузлами, з урахуванням їхніх характеристик і поточного стану. Нерівномірність розподілу може призводити до втрати ефективності системи, що вимагає розробки спеціальних алгоритмів і методів для її вирішення [66].

Основною проблемою є нерівномірний розподіл обчислювальних ресурсів у системі [67]. У результаті виникає ситуація, коли частина вузлів перевантажена, що збільшує час виконання задач, тоді як інші вузли залишаються недовантаженими, що призводить до неефективного використання ресурсів. Така ситуація може знижувати загальну продуктивність і стійкість системи [68].

Гетерогенність ресурсів вузлів: у мультикомп'ютерних системах вузли можуть мати різні апаратні характеристики, такі як кількість ядер процесора, обсяг

оперативної пам'яті чи пропускну здатність мережі. Гетерогенність ускладнює рівномірний розподіл завдань, оскільки однакове навантаження може по-різному впливати на вузли з різними ресурсами [69].

Динамічність завдань і зміна обчислювальних вимог: завдання у системі можуть змінювати свої вимоги до ресурсів у процесі виконання. Наприклад, спочатку задача може потребувати великих обсягів оперативної пам'яті, а пізніше інтенсивно завантажувати процесор. Ця динамічність створює труднощі для алгоритмів статичного розподілу навантаження [70].

Проблеми масштабування системи: зі збільшенням кількості вузлів у системі зростає складність управління ресурсами. Більші системи вимагають більш складних алгоритмів розподілу, здатних ефективно працювати у масштабованому середовищі [71].

Проблеми, пов'язані з балансуванням навантаження, є комплексними та впливають на різні аспекти роботи мультикомп'ютерних систем, включаючи їхню продуктивність, стабільність, надійність та економічність [72]. Якщо алгоритми розподілу завдань недостатньо гнучкі або не враховують особливостей окремих вузлів (наприклад, їхню гетерогенність за характеристиками обладнання чи поточним станом завантаженості), це призводить до того, що певні вузли отримують надмірну кількість завдань [73]. У результаті вони працюють довше, знижуючи загальну швидкість реагування системи, особливо для завдань, що потребують швидкої обробки або функціонування в режимі реального часу.

Інша частина вузлів при цьому може залишатися недовантаженою, що також є серйозною проблемою, оскільки ресурси таких вузлів фактично простоюють, хоча вони все одно продовжують споживати енергію [74]. Це не лише призводить до зростання витрат на електроенергію, але й збільшує витрати на обслуговування та експлуатацію обладнання, що не приносить достатньої користі [75].

Додатковою проблемою є витрати на передачу даних між вузлами системи. Часте або неправильне перенаправлення завдань через недостатньо оптимізовані алгоритми, які не враховують географічне розташування, мережеву близькість вузлів чи пропускну здатність каналів зв'язку, створює зайве навантаження на

мережеву інфраструктуру [76]. Це може спричинити підвищення затримок, уповільнення передачі інформації та навіть частковий або повний вихід з ладу мережі у пікові періоди завантаження [77].

Із зростанням кількості вузлів мультикомп'ютерної системи суттєво загострюються проблеми масштабування. Алгоритми розподілу завдань, які були ефективними для невеликих чи середніх систем, можуть не справлятися зі збільшенням кількості вузлів [78]. Це ускладнює централізоване керування, синхронізацію та моніторинг роботи системи в цілому. Через недостатню адаптивність алгоритмів та обмеженість можливостей централізованого управління зростає час, необхідний для прийняття рішень про розподіл завдань, що знову-таки знижує ефективність роботи великомасштабних систем [79].

У критичних ситуаціях перевантаження окремих вузлів може спричинити не лише уповільнення роботи, але й втрату важливих даних. Особливо це актуально для систем, що працюють з великими обсягами даних, де механізми збереження, резервування та відновлення інформації не є достатньо надійними або не передбачені належним чином. Втрата інформації під час інтенсивних операцій або аварійних ситуацій знижує надійність системи, погіршує довіру користувачів та може спричинити суттєві фінансові й репутаційні втрати [80].

1.3 Постановка задачі

Поставлена мета досягається виконанням таких основних завдань:

- проаналізувати існуючі методи розподілу навантаження в мультикомп'ютерних системах загального призначення;
- розробити метод розподілу навантаження з використанням кластеризації вузлів за рівнем завантаженості ресурсів;
- визначити критерії оцінювання ефективності запропонованого методу розподілу навантаження;
- реалізувати програмний прототип запропонованого методу та провести експериментальні дослідження його ефективності.

1.4 Висновки до першого розділу

У розділі проаналізовано типи мультикомп'ютерних систем (кластери, хмарні платформи, ґрід-системи, суперкомп'ютери) та приклади їх практичного застосування (Hadoop, Spark, AWS, Titan, Fugaku). Проведено порівняльний аналіз методів розподілу навантаження (статичних, динамічних, комбінованих), виділено перевагу динамічних через їхню адаптивність. Описано алгоритми кластеризації вузлів (ієрархічні методи, K-Means, DBSCAN, Fuzzy Clustering) для покращення балансування навантаження. Ідентифіковано основні проблеми: перевантаження вузлів, неефективність ресурсів, комунікаційні витрати, проблеми масштабування та енергоспоживання. Отримані результати є основою для розробки ефективного методу керування потужністю, представленого в наступних розділах.

2 МОДЕЛЬ РОЗПОДІЛУ НАВАНТАЖЕННЯ В МУЛЬТИ - КОМП'ЮТЕРНИХ СИСТЕМАХ

2.1 Забезпечення розподілу навантаження в мультикомп'ютерних системах загального призначення. Показники розподілу навантаження

Забезпечення рівномірного розподілу навантаження в мультикомп'ютерних системах є ключовим завданням для підвищення ефективності їх роботи. Оптимальне балансування дозволяє мінімізувати простой процесорів, знизити енергоспоживання, уникнути перевантаження окремих вузлів та покращити якість обслуговування користувачів. Для оцінки ефективності балансування навантаження використовуються різні показники, що відображають рівень використання ресурсів, швидкість обробки запитів та загальну продуктивність системи.

Одним із найважливіших показників ефективності балансування навантаження є рівень використання процесорних ресурсів, що визначається коефіцієнтом завантаженості процесора:

$$U_{CPU} = \frac{T_{busy}}{T_{total}} \times 100\%, \quad (2.1)$$

де T_{busy} – час, протягом якого процесор був зайнятий обчисленнями, а T_{total} – загальний час роботи системи.

Аналогічно, рівень використання пам'яті (U_{RAM}) визначається як:

$$U_{RAM} = \frac{M_{used}}{M_{total}} \times 100\%, \quad (2.2)$$

де M_{used} – час, протягом якого процесор був зайнятий обчисленнями, M_{total} – загальний час роботи системи.

Якщо ці показники наближаються до 100%, це означає, що вузол перевантажений і потребує балансування навантаження.

Затримка передачі даних (Δt) - є важливим показником для IoT-мереж, особливо в системах реального часу. Вона включає кілька складових:

$$\Delta t = t_{proc} + t_{queue} + t_{trans} + t_{prop}, \quad (2.3)$$

де t_{proc} - час обробки даних;

t_{queue} - час очікування в черзі на передачу;

t_{trans} - час передавання пакета через канал зв'язку;

t_{prop} - час поширення сигналу в мережі.

Мінімізація затримки є критичною для IoT-додатків, що працюють у режимі реального часу, таких як системи моніторингу здоров'я чи автономні транспортні засоби.

Пропускна здатність мережі (B) визначає максимальну кількість даних, що може передаватися через мережу за одиницю часу:

$$B = \frac{D}{T}, \quad (2.4)$$

де D – обсяг переданих даних (у бітах або байтах);

T – час передачі.

Якщо фактичний рівень використання пропускнуої здатності наближається до її межі, необхідно розподіляти навантаження між кількома каналами або застосовувати методи оптимізації трафіку.

Енергоспоживання пристроїв - у бездротових IoT-мережах оптимізація енергоспоживання є критично важливою. Споживана потужність (P) може бути розрахована як:

$$P = V \times I, \quad (2.5)$$

де V – напруга живлення;

I – споживаний струм.

Загальний час обробки запитів у системі включає:

$$T_{response} = T_{compute} + T_{comm} + T_{queue}, \quad (2.6)$$

де $T_{response}$ – загальний час відповіді системи (response time), тобто час від надходження запиту до отримання результату;

$T_{compute}$ – час, витрачений на обчислення, безпосереднє виконання задачі обчислювальним вузлом;

T_{comm} – час, витрачений на передавання даних між вузлами або між клієнтом і системою (communication time);

T_{queue} – час очікування в черзі до обробки (queueing time), тобто затримка через перевантаження або черговість обробки запитів.

Ефективне балансування навантаження допомагає зменшити енергоспоживання, наприклад, шляхом чергування активних вузлів або адаптації швидкості передавання даних.

Коефіцієнт завантаженості вузлів - цей показник визначає рівномірність розподілу навантаження між вузлами мережі. Його можна розрахувати за допомогою коефіцієнта варіації:

$$CV = \frac{\sigma}{\mu}, \quad (2.7)$$

де σ – стандартне відхилення рівнів завантаженості вузлів, μ – середнє значення завантаженості всіх вузлів.

Якщо CV близьке до нуля, то навантаження розподілене рівномірно, а якщо велике – система потребує балансування. Якість обслуговування (QoS) - якість обслуговування залежить від кількох параметрів, зокрема рівня втрати пакетів (PLR):

$$PLR = \frac{N_{lost}}{N_{sent}} \times 100\%, \quad (2.8)$$

де N_{lost} – кількість втрачених пакетів, N_{sent} – загальна кількість відправлених пакетів.

Високий рівень втрати пакетів може свідчити про перевантаження мережі та необхідність її оптимізації.

Частота зміни навантаження визначає, наскільки динамічно змінюється трафік у мережі. Її можна оцінити як похідну від функції завантаженості:

$$R = \frac{dL}{dt}, \quad (2.9)$$

де L – рівень навантаження мережі, t – час.

Якщо R високе, необхідно застосовувати адаптивні методи балансування, які швидко реагують на зміну умов.

Процес кластеризації вузлів у мультикомп'ютерних системах базується на аналізі поточного завантаження кожного вузла та групуванні їх у кластери за рівнем використання ресурсів.

Після розрахунку показників вище вузли мультикомп'ютерної системи розподіляються за рівнем завантаженості у три основні групи (кластери):

- перевантажені вузли ($U > U_{MAX}$) – вузли, що працюють на межі продуктивності та потребують розвантаження;
- збалансовані вузли ($U_{MIN} \leq U \leq U_{MAX}$) – вузли з оптимальним рівнем завантаження;
- недовантажені вузли ($U < U_{MIN}$) – вузли, які можуть взяти на себе додаткові завдання.

2.2 Модель розподілу навантаження в мультикомп'ютерних системах загального призначення

Забезпечення рівномірного розподілу навантаження в мультикомп'ютерних системах є критично важливим завданням для підвищення ефективності обчислень та зменшення простоїв процесорних ресурсів. У разі нераціонального використання обчислювальних потужностей окремі вузли можуть бути перевантаженими, у той час як інші простоюють. Це призводить до неефективного використання ресурсів та збільшення загального часу виконання процесів.

Запропонований метод балансування навантаження базується на кластеризації вузлів за рівнем їхньої завантаженості та динамічному управлінні міграцією процесів. На відміну від традиційних статичних методів, що визначають розподіл завдань лише під час їх ініціалізації, наш підхід дозволяє адаптивно перерозподіляти обчислювальні ресурси залежно від поточного стану системи.

Метод передбачає постійний моніторинг завантаженості вузлів і визначення середнього рівня використання ресурсів у системі. Для цього на кожному вузлі розраховується поточне навантаження, а глобальне середнє обчислюється за формулою:

$$L_{avg} = \frac{1}{N} \sum_{i=1}^N L_i, \quad (2.10)$$

де L_{avg} – середнє навантаження системи;

N – загальна кількість вузлів;

L_i – завантаженість i -го вузла.

Вузол вважається перевантаженим, якщо його рівень завантаження перевищує середнє значення на певний поріг:

$$L_i > L_{avg} + \delta. \quad (2.11)$$

У такому випадку вузол ініціює процес міграції частини своїх обчислень на інші вузли системи. Вибір вузла-приймача здійснюється на основі аналізу поточного стану інших вузлів. Оптимальним вважається вузол, що має найнижчий рівень завантаженості:

$$j = \arg \min (L_i), \quad L_i < L_{avg}. \quad (2.12)$$

Для забезпечення стабільності роботи системи та уникнення зайвих витрат на передачу даних міграція здійснюється лише у випадку, коли вона дозволяє суттєво зменшити перевантаження.

При визначенні процесу, що підлягає переміщенню, враховується кілька факторів. Зокрема, аналізується поточний розмір процесу, його використання процесорних ресурсів та комунікаційна активність. Найбільш доцільно переносити ті процеси, що створюють мінімальне навантаження на канал зв'язку і водночас значно розвантажують вузол, з якого вони мігрують.

Критерій вибору процесу можна представити таким чином:

$$P_{opt} = \arg \max (L_i - L_j - C(P_k)), \quad (2.13)$$

де $C(P_k)$ - оцінка складності міграції процесу P_k . Процеси, що мають високий рівень комунікаційної активності, мігрують лише в тому випадку, якщо це значно покращить загальне балансування системи.

Для узгодження міграції між вузлами використовується механізм обміну повідомленнями. Якщо вузол виявляє, що його навантаження перевищує допустиме значення, він надсилає широкомовне повідомлення Too High, сповіщаючи інші вузли про свою перевантаженість. У відповідь вузли з низьким рівнем завантаженості можуть надіслати сигнал Assent, що означає їхню готовність прийняти процес для виконання.

Якщо протягом певного часу вузол-приймач не знайдено, система оновлює середнє значення навантаження через повідомлення Change Average. Це дозволяє

коригувати параметри розподілу і уникати ситуацій, коли балансування не приносить бажаного ефекту.

Вузол, що отримав повідомлення Too High, аналізує власне завантаження і приймає рішення про можливість міграції процесу. Якщо він готовий прийняти додаткове навантаження, він надсилає відповідь і чекає на передачу процесу.

Якщо протягом визначеного інтервалу часу вузол-приймач не відповідає, перевантажений вузол може збільшити глобальне середнє значення навантаження, що дозволяє краще адаптуватися до нерівномірного розподілу ресурсів.

Приблизна модель методу балансування навантаження в мультикомп'ютерній системі наведена на рисунку 2.1.

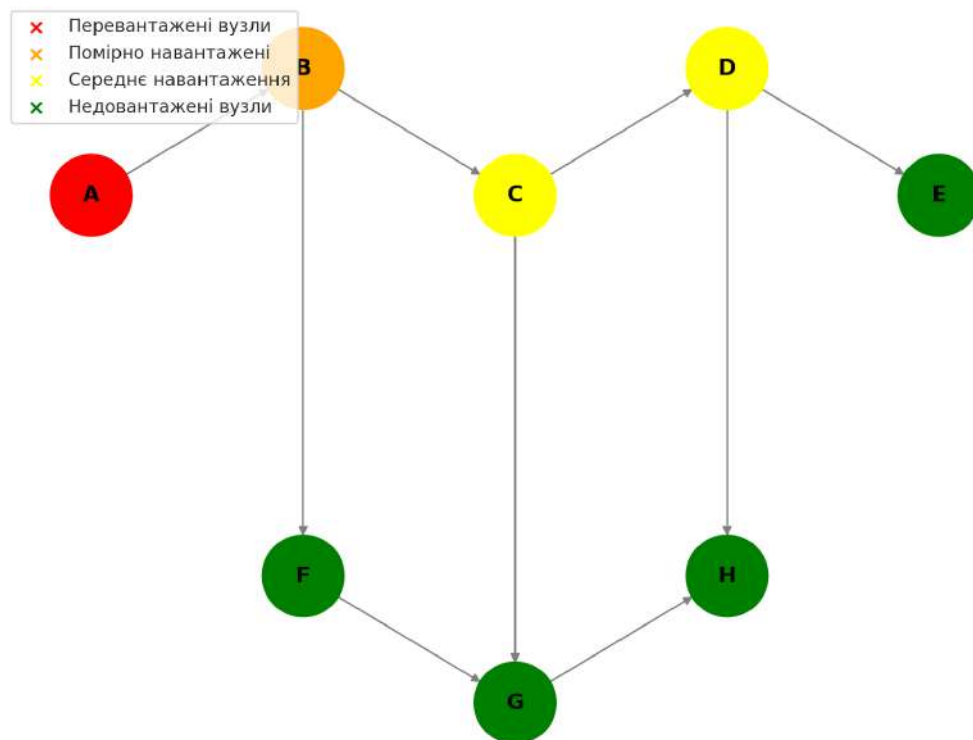


Рисунок 2.1 – Модель методу балансування навантаження

Червоний вузол (A) – перевантажений сервер, який ініціює процес міграції. Помаранчеві та жовті вузли (B, C, D) – вузли з різним рівнем навантаженості. Зелені вузли (E, F, G, H) – недовантажені вузли, які можуть прийняти процеси.

Стрілки показують комунікаційні зв'язки між вузлами, що дозволяють визначити оптимальний напрямок міграції процесів для вирівнювання

навантаження в системі.

Перевантажений вузол *A* надсилає сигнал про високе навантаження (Too High) до сусідніх вузлів.

Вузли з меншим завантаженням (*E, F, G, H*) можуть надіслати відповідь (Ассерт) і прийняти процеси для виконання.

Якщо міграція процесів відбувається ефективно, загальна завантаженість системи вирівнюється, і червоний вузол перестає бути критично перевантаженим.

У разі відсутності недовантажених вузлів середнє навантаження переглядається (Change Average) і система адаптується до нових умов.

2.3 Алгоритм забезпечення методу розподілу навантаження

Оскільки метою алгоритму вирівнювання навантаження є зменшення різниці у навантаженнях між вузлами мультикомп'ютера, то відповідним моментом для міграції процесу на інший вузол є момент, коли навантаження на його вузлі перевищує середнє значення для всіх вузлів мультикомп'ютера, що робить цей вузол перевантаженим. На цьому етапі здійснюється аналіз поточного стану завантаженості системи для визначення, чи потрібно перерозподілити навантаження між вузлами.

Алгоритм працює таким чином, що коли вузол перевантажений, він запускає процес міграції на інший вузол, що має менше завантаження. Аналогічно, відповідним вузлом для міграції процесу є вузол, навантаження якого є нижчим за загальне середнє значення. Це дозволяє вибрати вузол, який не тільки має менше навантаження, але й може ефективно впоратися з додатковими завданнями, не створюючи ризику для стабільності всієї системи. Тому вибір вузла для міграції заснований на аналізі поточного навантаження всіх вузлів мультикомп'ютера та пошуку найбільш оптимального для прийому нового процесу.

Для коректного виконання міграції процесу важливо мати точну інформацію про середнє навантаження на всі вузли мультикомп'ютера. Оцінка середнього навантаження дозволяє своєчасно визначити моменти, коли процеси необхідно

перемістити, аби уникнути перевантаження окремих вузлів та підтримати рівномірний розподіл навантаження серед усіх вузлів. Для цього кожен вузол постійно оновлює свої показники навантаження, які передаються в центральний блок моніторингу системи. Після аналізу цих даних визначається, чи потрібно ініціювати міграцію.

Окрім оцінки поточного навантаження на кожному з вузлів мультикомп'ютерної системи, при прийнятті рішень щодо міграції задач додатково враховуються такі критично важливі фактори, як пропускна здатність мережі між вузлами, доступні обчислювальні потужності та поточні системні ресурси, зокрема обсяг оперативної пам'яті, завантаженість процесора, дисковий простір і рівень використання пристроїв введення-виведення.

Такий комплексний підхід дозволяє точно моделювати ефективність потенційної міграції, знижуючи ризик перевантаження або втрати продуктивності. Завдяки цьому алгоритм прийняття рішень набуває адаптивності, забезпечуючи динамічне переналаштування конфігурації системи в реальному часі, що дозволяє кожному вузлу працювати з максимальною ефективністю відповідно до своїх поточних можливостей.

Це дає змогу оптимізувати розподіл обчислювального навантаження, мінімізувати час відгуку на запити користувача, зменшити затримки в обробці даних та підвищити загальну продуктивність і надійність усієї мультикомп'ютерної архітектури. На основі таких принципів реалізується механізм адаптивного перерозподілу навантаження шляхом кластеризації вузлів, який дозволяє формувати динамічні групи ресурсів залежно від їхнього поточного стану та взаємозв'язків. Це дає змогу ефективно керувати обчислювальними ресурсами навіть за умов швидкої зміни навантаження або часткових відмов, зберігаючи стабільність функціонування всієї системи.

Таким чином, запропонований підхід забезпечує високий рівень масштабованості, адаптивності та відмовостійкості, що є ключовими характеристиками сучасних розподілених обчислювальних середовищ, зокрема хмарних платформ, високопродуктивних обчислювальних кластерів та

інфраструктури обробки великих даних. Вигляд такого алгоритму можна побачити на рисунку 2.2.

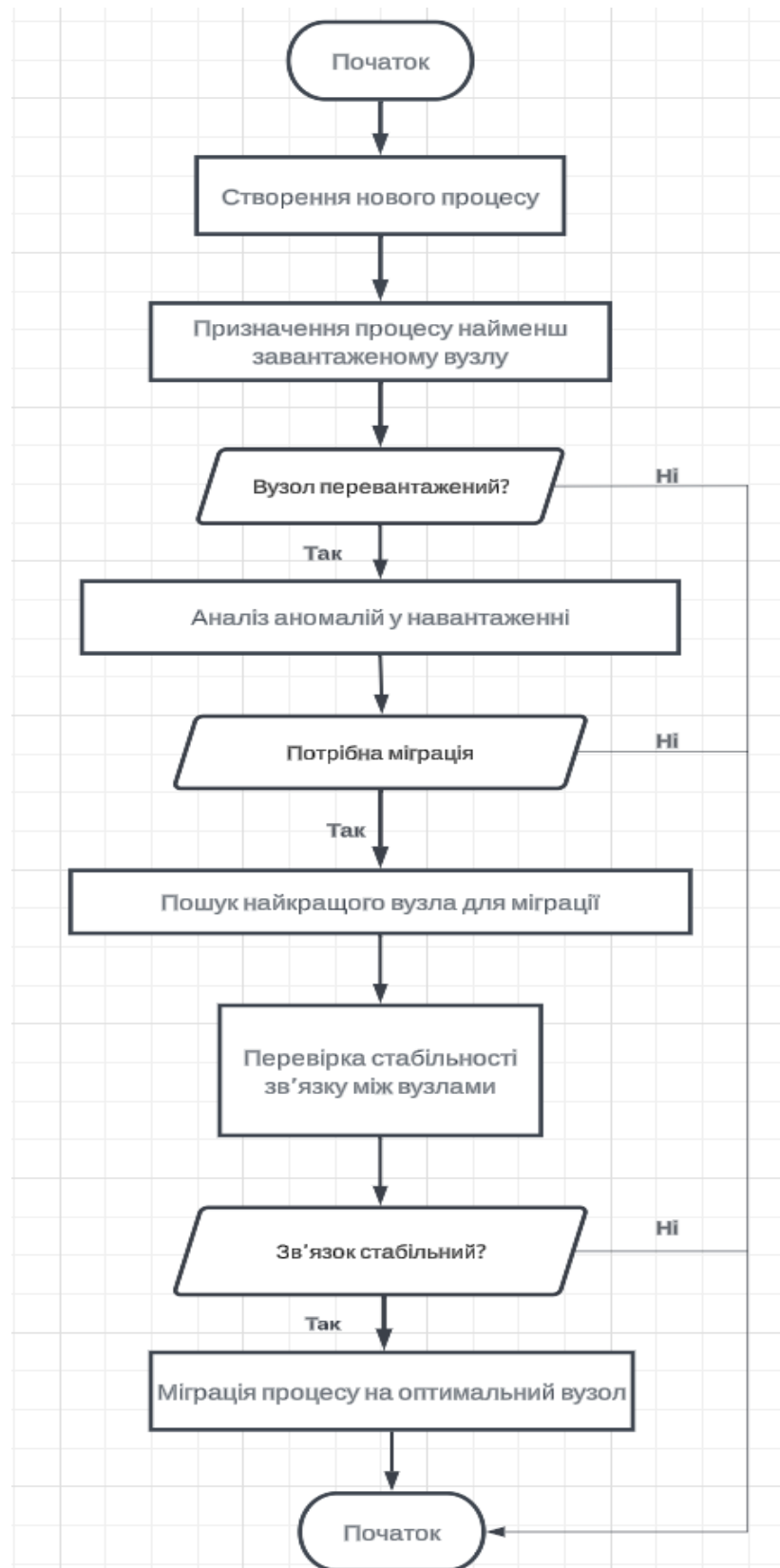


Рисунок 2.2 - Алгоритм оцінка середнього навантаження в вузлах

Розглянемо опис вище зазначеної схеми:

1. Початок – ініціалізація алгоритму.
2. Створення нового процесу – запуск нового завдання у системі.
3. Призначення процесу найменш завантаженому вузлу – вибір вузла з мінімальним поточним навантаженням.
4. Вузол перевантажений? Так - аналіз аномалій у навантаженні. Ні - повернення до стабільного стану.
5. Аналіз аномалій у навантаженні – перевірка рівномірності завантаженості вузлів.
6. Потрібна міграція? Так - пошук найкращого вузла для міграції. Ні - повернення до стабільного стану.
7. Пошук найкращого вузла для міграції – вибір найменш завантаженого вузла.
8. Перевірка стабільності зв'язку між вузлами – перевірка доступності мережевого каналу.
9. Зв'язок стабільний? Так - Міграція процесу. Ні - повернення до початкового стану.
10. Міграція процесу на оптимальний вузол – балансування навантаження шляхом перенесення процесу.
11. Фініш – завершення алгоритму.

Алгоритм ефективно реагує на зміни у завантаженості вузлів, дозволяючи динамічно перерозподіляти ресурси і підтримувати рівномірний розподіл навантаження в межах системи. Це дозволяє зменшити ймовірність перевантаження окремих вузлів і підтримувати високу продуктивність системи.

Завдяки швидкому реагуванню на зміну навантаження, процеси не простоюють у чергах, що дозволяє підвищити загальну ефективність системи. Алгоритм активно перерозподіляє ресурси, щоб підтримувати стабільну роботу навіть у випадку значних змін у завантаженості.

Алгоритм використовує децентралізовану архітектуру, що дозволяє кожному вузлу самостійно приймати рішення про розподіл навантаження. Це не тільки

зменшує навантаження на центральний контролер, але й дозволяє забезпечити швидке реагування на зміни в системі.

Система здатна адаптуватися до зміни кількості вузлів у мережі, що дозволяє використовувати алгоритм у різних типах обчислювальних середовищ, включаючи хмарні платформи та високопродуктивні комп'ютерні системи. Це робить алгоритм універсальним і гнучким для широкого спектру застосувань.

Алгоритм дозволяє зменшити загальне енергоспоживання системи завдяки оптимізації використання ресурсів і зменшенню кількості зайвих міграцій процесів. Це є важливим аспектом для сучасних дата-центрів та обчислювальних кластерів, де енергоефективність є критичним фактором.

Алгоритм включає механізми перевірки стабільності зв'язку між вузлами, що дозволяє запобігти помилкам, які можуть виникнути під час міграції процесів у разі нестабільного з'єднання. Це підвищує надійність і стійкість системи до зовнішніх збоїв.

Однією з ключових проблем є затримка, пов'язана з переміщенням процесів між вузлами. Якщо процес є об'ємним або активно використовує ресурси вводу-виводу, його перенесення може спричинити додаткові накладні витрати. Для вирішення цієї проблеми можна застосувати механізми попереднього копіювання даних або буферизацію, що мінімізують затримку при переміщенні великих обсягів інформації.

Процес передачі даних між вузлами мультикомп'ютера може призвести до перевантаження мережі, особливо якщо велика кількість процесів одночасно переміщується. Це може суттєво вплинути на загальну продуктивність системи. Одним з можливих рішень є оптимізація механізму вибору вузла-приймача, що дозволяє зменшити кількість зайвих переміщень і знизити навантаження на мережу.

Система може зіткнутися з проблемою неточності оцінки завантаженості вузлів через коливання навантаження в реальному часі. Це може призвести до помилок у прийнятті рішення про міграцію процесів. Використання прогнозування або адаптивних методів оцінки навантаження дозволить отримати більш точні

результати і зменшити ймовірність помилок у роботі алгоритму.

Якщо алгоритм занадто часто ініціює міграцію процесів, це може призвести до ситуації, коли ресурси витрачаються не на виконання обчислень, а на постійне перенесення процесів. Введення граничних значень для ініціації міграції дозволить зменшити ймовірність виникнення такого стану, зберігаючи ефективність розподілу навантаження.

Алгоритм є особливо ефективним у випадку однорідної архітектури вузлів, але при роботі з гетерогенними системами можуть виникати труднощі з уніфікацією критеріїв оцінки продуктивності вузлів. Для вирішення цієї проблеми можна використовувати додаткові вагові коефіцієнти, які враховуватимуть різницю в продуктивності вузлів і дозволять забезпечити ефективний розподіл навантаження.

Якщо зв'язок між вузлами є нестабільним, процес міграції може не завершитися успішно або спричинити втрату даних. Використання механізмів резервного копіювання даних, а також повторного відправлення запитів у разі помилок, дозволить значно підвищити стійкість алгоритму до проблем із мережею.

2.4 Політики або стратегії в динамічному балансуванні навантаження

В алгоритмах динамічного балансування навантаження існує чотири основні політики (стратегії), які визначають, як і в який спосіб здійснюється балансування навантаження між різними вузлами системи. Кожна з цих політик виконує окремі функції, що забезпечують ефективність і швидкість процесу балансування. Ось основні з них:

Політика переносу (Transfer Policy): ця частина алгоритму відповідає за вибір завдання для передачі з локального вузла на віддалений. Вона регулює, яке саме завдання має бути передано між вузлами, щоб оптимізувати використання ресурсів в системі. Завдання можуть бути передані в залежності від різних факторів, таких як навантаження на поточному вузлі, час обробки або пріоритет завдання.

Політика вибору (Selection Policy): політика вибору визначає, які процесори

або вузли будуть залучені до обміну навантаженням. Це означає, що система повинна мати механізм для відповідного співвіднесення задач з певними процесорами, що мають необхідні ресурси для виконання завдання. Вибір процесора здійснюється на основі доступних характеристик процесорів та поточного стану навантаження.

Політика місця розташування (Location Policy): ця частина алгоритму вибирає цільовий вузол для переданого завдання. Політика місця розташування визначає, куди саме слід передати завдання після того, як воно буде вибрано для переносу. Це залежить від доступності віддалених вузлів, а також від поточного стану мережі та навантаження на вузлах.

Інформаційна політика (Information Policy): ця частина алгоритму відповідає за збір і обробку інформації про стан вузлів у системі. Інформаційна політика дозволяє системі отримувати актуальні дані про навантаження на кожному з вузлів, що дає змогу коректно планувати подальші стратегії переносу завдань і забезпечити оптимальне використання ресурсів.

У процесі динамічного балансування ці стратегії взаємодіють між собою, що допомагає визначити, чи потрібно передавати завдання на інший вузол, чи виконувати його локально, або ж вибирати найбільш підходящий вузол для перенесення завдання з урахуванням наявних ресурсів.

Процеси, що відбуваються в системі динамічного балансування навантаження, зображені на рисунку 2.3. зокрема стратегії переносу, розташування, інформаційну стратегію та взаємодію між ними. Діаграма ілюструє послідовність ухвалення рішень, що допомагають ефективно розподіляти навантаження між вузлами в мережі.

Вхідна задача — це початковий етап процесу, на якому система отримує нову задачу для обробки. Задача надходить на один з вузлів в системі для подальшого вирішення.

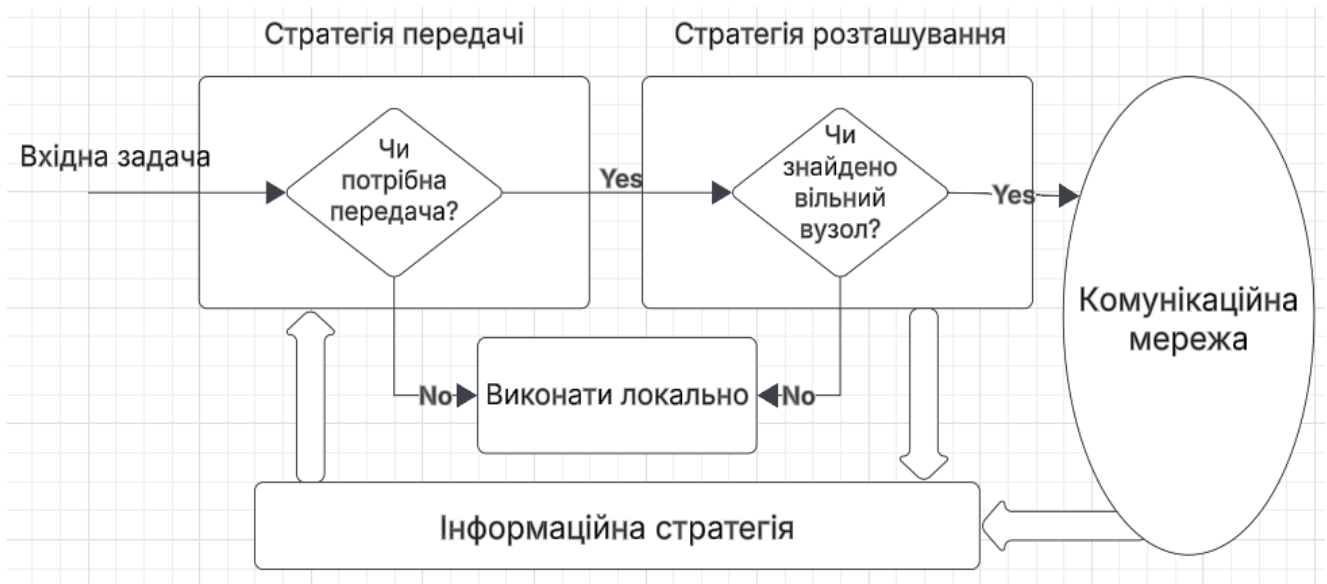


Рисунок 2.3 - Діаграма взаємодії основних компонентів алгоритму динамічного балансування навантаження

Стратегія переносу - на цьому етапі перевіряється, чи потребує задача переносу з поточного вузла на віддалений. Якщо задача підлягає перенесенню (відповідь "Так"), вона передається до наступного етапу - стратегії розташування.

Стратегія розташування - на цьому етапі система перевіряє, чи є вільні віддалені вузли, на які можна передати задачу. Якщо вільний вузол знайдений (відповідь "Так"), задача передається через Комунікаційну мережу до віддаленого вузла для виконання. Якщо ж вільний вузол не знайдений (відповідь "Ні"), система приймає рішення про виконання задачі локально, тобто на поточному вузлі.

Інформаційна стратегія - паралельно з іншими етапами виконується збір та обробка інформації про стан вузлів у системі. Це включає дані про поточне навантаження на вузли, що дозволяє ефективно управляти ресурсами та визначати, де краще виконати задачу - локально або на віддаленому вузлі.

2.5 Висновки до другого розділу

У цьому розділі досліджено важливість забезпечення рівномірного розподілу навантаження в мультикомп'ютерних системах для підвищення їх ефективності.

Оптимальне балансування навантаження є необхідним для мінімізації простоїв процесорів, зменшення енергоспоживання та уникнення перевантаження окремих вузлів. Розглянуті показники ефективності балансування, серед яких рівень використання процесорних та пам'ятових ресурсів, затримка передачі даних, пропускна здатність мережі, енергоспоживання пристроїв, а також коефіцієнт завантаженості вузлів та якість обслуговування, дозволяють здійснити детальну оцінку стану системи та вчасно вжити заходів для вирівнювання навантаження.

Процес кластеризації вузлів у мультикомп'ютерних системах, заснований на рівні завантаженості, дозволяє ефективно групувати вузли за різними критеріями та застосовувати відповідні методи балансування. Розподіл вузлів на перевантажені, збалансовані та недовантажені дає змогу реалізувати адаптивне управління ресурсами, що підвищує продуктивність та стабільність системи.

Запропонований метод забезпечення розподілу навантаження на основі динамічного моніторингу та міграції процесів дозволяє досягти рівномірного розподілу ресурсів та уникнути перевантаження окремих вузлів. Алгоритм балансування навантаження демонструє ефективність у вирішенні проблем, пов'язаних з неефективним використанням ресурсів, зокрема, через наявність перевантажених вузлів або недовантажених вузлів.

Стратегії в динамічному балансуванні навантаження, такі як політика переносу, вибору, розташування та інформаційна політика, взаємодіють для забезпечення ефективного розподілу завдань між вузлами та досягнення оптимальної продуктивності. Алгоритм, заснований на цих політиках, дозволяє забезпечити стабільну роботу системи, адаптуючи її до змінних умов та підтримуючи баланс між вузлами.

Загалом, підходи та методи, представлені в цьому розділі, мають велике значення для реалізації ефективних мультикомп'ютерних систем, здатних до адаптивного управління та оптимізації ресурсів у реальному часі, що є критично важливим для забезпечення високої продуктивності та енергоефективності сучасних обчислювальних платформ.

3 СИСТЕМА ЗАБЕЗПЕЧЕННЯ РОЗПОДІЛУ НАВАНТАЖЕННЯ ЗГІДНО КЛАСТЕРИЗАЦІЇ ВУЗЛІВ

3.1 Збір даних про завантаження віртуальної машини

Запропоноване мною рішення зосереджене на розподілі обчислювального навантаження у багатокомп'ютерних системах загального призначення шляхом кластеризації обчислювальних вузлів відповідно до рівня їхнього поточного використання ресурсів. Для реалізації цього підходу я застосовую комбіновану (гібридну) модель машинного навчання, що об'єднує як керовані методи (supervised learning) — зокрема штучні нейронні мережі (ANN), так і некеровані методи (unsupervised learning), такі як ВОК (Best of K-кластеризація), для оптимального формування груп вузлів з подібними характеристиками навантаження.

Система складається з множини вузлів обчислення $P = \{PM1, PM2 \dots PMM\}$, кожен з яких містить набір програмно або апаратно ізольованих обчислювальних одиниць (віртуалізованих або фізичних), які в рамках моделі позначено як $VM = \{VM1, VM2, \dots, VMN\}$. Система оперує множиною завдань користувача $T = \{T1, T2, \dots, TS\}$, які надходять для обробки в режимі реального часу або пакетної обробки. Для балансування навантаження між вузлами та їх обчислювальними одиницями задіяно дві ключові програмні сутності:

- диспетчер вузлів – відповідає за динамічне призначення завдань відповідним вузлам відповідно до їхнього поточного стану;
- балансувальник навантаження – координує міжкластерну рівновагу, забезпечуючи стабільність і ефективність розподілу.

3.1.1 Обчислення навантаження за допомогою CNN та RNN

Для точного визначення рівня навантаження кожного вузла використовується гібридний підхід на основі глибокого навчання із залученням двох типів нейронних мереж:

- CNN (Convolutional Neural Network) — використовується для

виділення просторових ознак з вхідних даних (наприклад, телеметрії, показників завантаження CPU/пам'яті/мережі). Цей тип мережі ефективно визначає закономірності у структурованих даних сіткового типу (наприклад, розподіл навантаження в масиві вузлів);

- RNN (Recurrent Neural Network) — забезпечує врахування часових залежностей у даних, що дозволяє враховувати динаміку змін навантаження протягом певного періоду часу. Це особливо важливо для прогнозування станів системи в умовах динамічного середовища.

Convolutional Neural Networks є одними з найефективніших методів у сфері аналізу структурованих і візуально-кодованих даних. Архітектура CNN передбачає послідовне чергування згорткових шарів, шарів підвибірки (pooling layers) та повнозв'язних шарів (fully connected layers). У контексті запропонованої системи CNN виконує роль модуля попереднього аналізу даних, що дозволяє виявити типові закономірності використання ресурсів на рівні вузлів або їхніх обчислювальних одиниць.

Для попереднього аналізу і структурного виділення просторових ознак у даних про навантаження обчислювальних вузлів використовується згорткова нейронна мережа (Convolutional Neural Network, CNN). Нижче розглянуто ключові компоненти її архітектури та їхню роль у контексті поставленого завдання.

Основа обробки в CNN становить згортковий шар, який застосовує набір фільтрів (ядер) до вхідних даних з метою виявлення локальних закономірностей.

У нашому випадку ці дані можуть відповідати часовим рядам показників ресурсного навантаження (наприклад, CPU, пам'ять, мережеве навантаження) кожного вузла. Кожне згорткове ядро аналізує невелику ділянку вхідного простору — так зване рецептивне поле, що дозволяє виявляти локальні особливості навантаження.

Математично це описується як згорткова операція між вхідним сигналом $i_c(x, y)$ та ядром $e_{kl}(u, v)$ що формує карту ознак $f_{kl}(p, q)$. Повна карта ознак записується як:

$$F_{kl} = [f_{kl}(1,1), \dots, f_{kl}(p, q), \dots, f_{kl}(P, Q)], \quad (3.1)$$

де F_{kl} – це фільтр (ядро згортки) у згортковій нейронній мережі з індексами k та l , який представлений як матриця розміром $P \times Q$;

$f_{kl}(p, q)$ – елемент фільтра, що знаходиться в p -му рядку та q -му стовпчик;

P, Q – відповідно кількість рядків і стовпчиків фільтра (розміри ядра згортки); (p, q) – конкретна позиція всередині фільтра, де $1 \leq p \leq P$ і $1 \leq q \leq Q$.

Pooling або шар об'єднання зменшує розмірність просторових ознак, агрегуючи найважливішу інформацію в межах кожного рецептивного поля. Це дозволяє мережі бути стійкою до незначних зміщень або спотворень у вхідних даних. У межах нашого завдання такий шар допомагає зменшити переобладнання (overfitting) та підвищити здатність моделі до узагальнення.

Процедура pooling описується функцією:

$$Z_{kl} = g_p(F_{kl}), \quad (3.2)$$

де g_p - pooling-операції (наприклад, max-pooling або average-pooling).

Для того, щоб модель могла навчатися складних закономірностей, застосовується функція активації, яка вводить нелінійність у модель. На практиці це дозволяє відображати залежності між різними аспектами навантаження, які не можна описати лише лінійними співвідношенням:

$$T_{kl} = g_a(F_{kl}), \quad (3.3)$$

де g_a – функція активації (наприклад, ReLU або tanh).

На завершальному етапі CNN використовує повністю зв'язаний шар, який об'єднує всі виділені раніше ознаки та формує підсумкове уявлення про стан системи. Це дозволяє моделі враховувати взаємозв'язки між різними фрагментами інформації (наприклад, між навантаженням процесора та мережевою активністю),

що важливо для точної класифікації або регресії.

Вихід цього шару, як правило, є основою для:

- класифікації вузлів за рівнем навантаження (наприклад, кластеризація);
- прогнозування навантаження в наступні моменти часу (у зв'язці з RNN);
- передавання результатів у балансувальник для прийняття рішень.

Для моделювання динаміки зміни навантаження в багатокомп'ютерній системі у реальному часі було застосовано рекурентну нейронну мережу (RNN), яка відома своєю ефективністю при роботі з часовими рядами. На відміну від класичних багат шарових перцептронів (MLP), RNN мають внутрішні зворотні зв'язки між прихованими шарами, що дозволяє мережі зберігати та передавати інформацію про попередні стани.

Завдяки цій властивості RNN здатна моделювати довготривалі часові залежності між подіями, що є критично важливим при оцінці навантаження вузлів, де історія попередніх станів може впливати на поточне навантаження.

У моїй реалізації, вхідні послідовності, що складаються з показників використання ресурсів (CPU, RAM, мережа тощо), передаються у ковзному часовому вікні через рекурентну структуру. Для обчислення вихідного стану прихованого шару на момент часу t використовується функція активації гіперболічного тангенса:

$$h_t = \tanh(W_x h_t + W_h h_{t-1} + b), \quad (3.4)$$

де h_t – поточний стан прихованого шару;

x_t – вхідний вектор на момент часу t ;

W_x – вагова матриця між входом і прихованим шаром;

W_h – вагова матриця між попереднім і поточним станом прихованого шару;

b – вектор зсуву (bias).

Ваги мережі оптимізуються шляхом зворотного поширення помилки в часі

(Backpropagation Through Time) з урахуванням похідних від функції втрат щодо кожного шару.

Особливістю моделі є поєднання CNN та RNN, де згорткова мережа (CNN) виконує попереднє виділення просторових ознак із вхідних метрик, а RNN – аналіз часових змін. Після об'єднання виходів обох мереж формується оцінка навантаження для кожного вузла.

Отримані значення використовуються для кластеризації вузлів системи за рівнем навантаження, що дозволяє ефективно розподіляти ресурси відповідно до поточної ситуації в системі.

3.2. Групування вузлів за рівнем навантаження на основі гібридного методу кластеризації

На основі попередньо розрахованих значень навантаження, отриманих за допомогою моделі гібридної нейронної мережі (CNN + RNN), виконується кластеризація обчислювальних вузлів системи на перевантажені та недовантажені. Такий підхід дозволяє приймати обґрунтовані рішення щодо перенесення задач або балансування навантаження в реальному часі.

Для класифікації вузлів застосовується стратегія порогового поділу: визначається певне контрольне значення навантаження — порогове значення. Якщо навантаження вузла перевищує цей поріг, вузол вважається перевантаженим, і включається до відповідного кластера. Якщо ж навантаження менше або дорівнює цьому значенню, вузол належить до кластера недовантажених.

Позначимо через L_i - розраховане навантаження для вузла i , а $LoadThreshold$ - задане порогове значення. Тоді математичний поділ здійснюється за такими правилами: $OverloadNodes = \{i \mid L_i > LoadThreshold\}$ $UnderloadNodes = \{i \mid L_i \leq LoadThreshold\}$.

Такий підхід дозволяє динамічно адаптувати поведінку системи, зменшуючи ризик локального перевантаження окремих вузлів при одночасному недовантаженні інших. Значення порогу може бути гнучко змінене відповідно до

конкретних вимог — наприклад, у системах з жорсткими часовими обмеженнями або у випадках обмеженого енергоспоживання.

У рамках реалізованого рішення кластеризація виконується за допомогою гібридного механізму, який поєднує механізм навчання з підкріпленням (reinforcement learning) і оптимізаційні евристики, що дозволяє ефективно контролювати формування кластерів навіть у динамічному середовищі з великою кількістю змінних параметрів.

3.2.1. Оптимізована кластеризація вузлів

Для досягнення ефективного групування обчислювальних вузлів за рівнем навантаження в даному рішенні використано модель підкріплювального навчання (Reinforcement Learning, RL). Основна ідея (на рисунку 3.1) цього підходу полягає у взаємодії між агентом (керуючим елементом, що приймає рішення) та середовищем (тобто багатовузловою обчислювальною системою, що постійно змінюється).

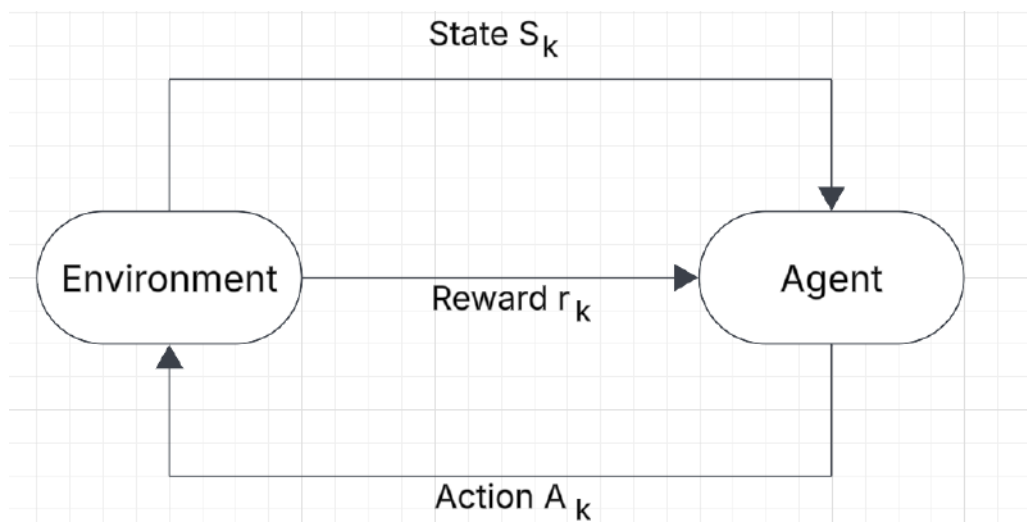


Рисунок 3.1 - Базова схема RL

Агент отримує інформацію про поточний стан системи — наприклад, рівень навантаження кожного вузла — і на основі цього обирає дію: перемістити завдання,

перегрупувати вузли або залишити все без змін.

Після виконання дії середовище реагує, оновлюючи свій стан, і видає винагороду — цифровий показник, що відображає ефективність дії (наприклад, зниження перевантаження або енергоспоживання). Процес є циклічним: кожна дія агента впливає на систему, яка у відповідь змінює свій стан, і надає зворотний зв'язок. Завдання агента – навчитися оптимальній політиці дій, яка забезпечить максимальну сумарну (накопичену) винагороду з урахуванням майбутніх наслідків.

У підкріплювальному навчанні функція $Q(s, a)$ – це числова оцінка того, наскільки вигідно виконати дію a у стані s . Іншими словами, це прогноз майбутньої вигоди з урахуванням знижки на майбутнє (дисконт). У випадку безперервного часу, функція Q визначається як математичне очікування інтегралу від потоку винагород у часі:

$$Q(s, a) = E \left[\int_{t_0}^{\infty} e^{-\beta(t-t_0)} r(t) dt \mid s_0 = s, a_0 = a \right], \quad (3.4)$$

де $r(t)$ – функція миттєвої винагороди в момент часу t ;

β – коефіцієнт дисконтування, що зменшує вагу майбутніх подій;

E – математичне очікування.

Ця формула визначає, як оцінюється довгостроковий ефект кожної дії, і є основою алгоритму Q-навчання, який використовується агентом для покрокового вдосконалення своїх рішень.

У моєму підході ця логіка використовується для динамічної кластеризації вузлів: вузли з подібними характеристиками навантаження об'єднуються в кластери з метою досягнення рівномірного розподілу задач, запобігання перевантаженню та підвищення ефективності використання ресурсів. В процесі навчання агент постійно оновлює свою Q-функцію на основі отриманих даних, покращуючи якість кластеризації з кожною ітерацією.

Функція ставки винагороди позначається як $r(t)$, що визначає винагороду,

яку агент отримує за виконання певної дії в конкретному стані на момент часу t . Ставка дисконту позначена як β і є важливим параметром у процесі оцінки майбутніх винагород. Це дозволяє агенту приймати рішення з урахуванням не лише поточної, але й майбутньої винагороди.

Підхід керованого подіями у часі застосовується в Q-навчанні, яке є адаптивною стратегією онлайн-навчання з підкріпленням. Це означає, що агент навчається у реальному часі, поступово покращуючи свою стратегію, знижуючи витрати, пов'язані з періодичними оновленнями дискретних значень часу. Q-навчання включає безперервне формулювання функції значення $Q(s, a)$, яка відображає вигоду або "цінність" виконання дії a в стані s .

Оновлення значення функції $Q(s, a)$ описується в рівнянні (16). Тут α – це швидкість навчання, яка контролює, наскільки швидко агент оновлює свої оцінки на основі нових даних, а β – це коефіцієнт дисконтування, що враховує зменшення значення майбутніх винагород.

Процес оновлення функції $Q(s, a)$ є результатом переходу в новий стан та вибору найкращої дії, що забезпечить максимальний виграш у майбутньому. Для кожного кроку, залежно від часу перебування агента в певному стані, відбувається корекція значення цієї функції.

Відмінною рисою цього підходу є використання змішаних стратегій, що дозволяє адаптувати методи розподілу ресурсів або інших задач до конкретних умов. В процесі застосування таких стратегій використовуються офлайн і онлайн фази, які включають процедури, що автоматично оцінюють та вдосконалюють стратегії на основі попередніх спостережень і даних.

Завдяки інтеграції реальних даних, зібраних на етапі створення автономної системи навчання з підкріпленням, можна розробити систему, що ефективно вирішує завдання, такі як розподіл ресурсів у хмарних обчисленнях. Зберігання профілів станів та оцінок $Q(s, a)$ дозволяє уникнути розбіжностей і сприяє подальшому вдосконаленню алгоритмів.

У контексті такого підкріплювального навчання Hybrid Lyrebird Falcon Optimization (HLFO) виступає як ефективний механізм для пошуку оптимального

набору параметрів кластеризації на кожному кроці прийняття рішення. HLFO забезпечує високу якість глобального та локального пошуку в просторі параметрів, що дозволяє швидко знаходити найкращі комбінації розподілу вузлів, вибору кластерних голів і маршрутизації даних. Комбінування HLFO з DRL створює синергію: агент DRL формує стратегічний план дій на основі довгострокових цілей, а HLFO виконує локальну оптимізацію для точного налаштування параметрів.

Такий підхід особливо ефективний у динамічних сценаріях, характерних для Інтернету речей (IoT), хмарних обчислювальних платформ і розподілених мереж, де умови швидко змінюються і потрібно оперативно реагувати на появу нових вузлів, зміну трафіку або погіршення стану мережі. Адаптивність системи дозволяє зменшити енергоспоживання, підвищити стійкість до відмов та забезпечити більш стабільну якість зв'язку.

Крім того, інтеграція HLFO з DRL сприяє зниженню обчислювальних затрат, оскільки HLFO здатний ефективно і швидко виконувати оптимізацію навіть у великих просторах параметрів, що важливо для систем з обмеженими ресурсами. Внаслідок цього, такі системи демонструють кращу масштабованість і більш тривалий час роботи у порівнянні з традиційними методами кластеризації.

3.2.2. Кластеризація на основі алгоритму оптимізації Hybrid Lyrebird Falcon

У цьому підрозділі розглядається вдосконалений алгоритм, який поєднує елементи оптимізації на основі мисливської поведінки соколів і лірохвостих птахів. Алгоритм орієнтований на вирішення складних завдань кластеризації та оптимізації, зокрема в умовах змінних або невизначених параметрів, де традиційні методи не завжди ефективні. Соколині тактики полювання, такі як використання чітко визначених прийомів і непередбачуваних маневрів, разом із поведінкою лірохвостих птахів, додають елемент непередбачуваності та гнучкості в алгоритм.

Кластеризація за допомогою HLFO дає змогу формувати більш оптимальні кластери з точки зору внутрішньої щільності та міжкластерної різниці. Це особливо актуально для складних динамічних середовищ, таких як хмарні обчислення,

розподілені системи, мережі IoT, де дані надходять потоково, і рішення мають прийматися в реальному часі. HLFO дозволяє гнучко адаптуватися до змінних умов, водночас зберігаючи високу точність кластеризації. Ось покрокове пояснення цього процесу

Крок 1: ініціалізація параметрів. У першій фазі алгоритму ініціалізуються фактори пошуку та межі регулювання. Це дозволяє визначити простір для руху та змін вектора швидкості для кожного елемента в популяції, а також встановлює межі для оптимальних значень. Основні параметри включають значення FOA (Factors of Attraction), які регулюють поведінку кожного птаха, а також межі швидкості – визначені значення, в межах яких швидкість птахів може варіюватися.

Крок 2: рух та позиціонування. На цьому етапі соколи змінюють своє розташування згідно з параметрами, які встановлюються для кожного виміру простору. Вектор швидкості для кожного елемента генерується випадковим чином у межах визначених значень. Швидкість для кожної позиції в просторі змінюється динамічно залежно від поточного стану, щоб дозволити птахам адаптуватися до змінних умов пошуку.

Крок 3: оцінка ефективності. Один з основних етапів — це оцінка ефективності (або придатності) кожної позиції птаха. Це досягається через визначення глобальних і персональних оптимальних значень, де кожен птах оцінюється за тим, як добре він виконує свою роль у пошуковому процесі. Значення здоров'я кожного птаха (або його придатність) визначається через генерування векторів оцінки. Птахи, які показують кращі результати, зберігають свої позиції або переміщуються ближче до оптимальних точок.

Крок 4: визначення ймовірності стрибка. Під час цього етапу визначається ймовірність, з якою птах здійснює стрибок або змінює свою позицію, ґрунтуючись на когнітивних та соціальних аспектах. Кожен птах створює два випадкові елементи для визначення ймовірності стрибка, що визначає його здатність змінювати своє місцезнаходження в залежності від стратегії пошуку та результатів попередніх кроків. Це додатково ускладнює передбачення рухів птахів для потенційних «хижаків» або інших учасників системи.

Крок 5: інтеграція поведінки лірохвостого птаха. Цей етап включає елементи полювання лірохвостого птаха, який використовує випадковість для втечі. Лірохвости, за своєю природою, втікають непередбачувано, що робить їх складними для передбачення хижаками. Цей елемент додається до алгоритму, щоб забезпечити додаткову непередбачуваність і адаптивність в процесі кластеризації. Алгоритм поєднує поведінку сокола, що орієнтується на «слабкі місця» цілі, з нестандартними стратегіями лірохвостого птаха, таким чином підвищуючи шанси на «втечу» або успішне завершення оптимізації. Використовуючи логарифмічні трансформації, алгоритм змінює стратегії птахів на основі того, чи знаходяться вони близько до своєї мети або здобичі.

Крок 6: оновлення позицій на основі стратегії втечі. З кожним кроком алгоритм визначає нове положення для кожного птаха, враховуючи зміни в його стратегії втечі. Алгоритм моделює ці рухи, щоб оптимізувати результат, а також зберігає можливість втечі від «хижаків». Завдяки цьому птахи можуть адаптувати свої стратегії в залежності від зміни умов пошуку, що дозволяє знаходити найбільш ефективні рішення для кластеризації.

Оскільки алгоритм поєднує елементи від декількох видів поведінки (соколині та лірохвості), він дозволяє моделювати різноманітні ситуації, зокрема такі, що включають змінні умови і нестабільність. Він використовує комбінацію випадкових рухів і детермінованих стратегій для знаходження оптимальних рішень, що робить його потужним інструментом для задач кластеризації в складних умовах.

Таким чином, алгоритм Hybrid Lyrebird Falcon дозволяє проводити кластеризацію на основі поєднання адаптивних стратегій і мисливської поведінки, що дозволяє забезпечити більш ефективний і непередбачуваний процес оптимізації. Це підвищує ймовірність знаходження глобальних оптимумів навіть в умовах високої невизначеності.

3.3 Вибір типу архітектури та зразків проектування

У процесі розробки методу розподілу навантаження у мультикомп'ютерних системах загального призначення згідно з кластеризацією вузлів за рівнем завантаженості ресурсів постає завдання забезпечити гнучке, масштабоване та надійне рішення. Для досягнення цих цілей необхідно визначити найбільш доцільний тип архітектури та застосувати відповідні зразки проектування (design patterns), що спростять подальшу реалізацію, супровід і розвиток системи.

З урахуванням особливостей проєкту — наявності кількох вузлів (комп'ютерів), які виконують обчислювальні завдання спільно, та необхідності постійного обміну даними про завантаженість ресурсів — доцільним є використання розподіленої (distributed) архітектури з елементами кластеризації.

Масштабованість - у розподіленій архітектурі можна додавати нові вузли або виводити їх з експлуатації динамічно, забезпечуючи гнучке масштабування системи відповідно до потреб у продуктивності.

Надійність та стійкість до відмов - у разі виходу з ладу одного з вузлів інші продовжують працювати, що зменшує ризик повної зупинки системи.

Ефективне використання ресурсів - кожен вузол виконує обчислення та обробку запитів відповідно до свого поточного стану завантаженості, що покращує загальну продуктивність та знижує час відгуку.

З огляду на потребу обробляти та аналізувати дані про завантаженість ресурсів (процесор, пам'ять тощо) у реальному чи квазіреальному часі, а також здійснювати періодичну кластеризацію, доцільно організувати роботу системи у вигляді багат шарової (layered) або мікросервісної структури. У даному дипломному проєкті найбільш виправданим буде багат шаровий підхід, оскільки він чітко відокремлює логіку збору даних, аналізу та ухвалення рішень, а також надає можливість поступового розширення функціональності. Приклад спрощеної тривірневої (3-layer) архітектури зображено на рисунку 3.2.

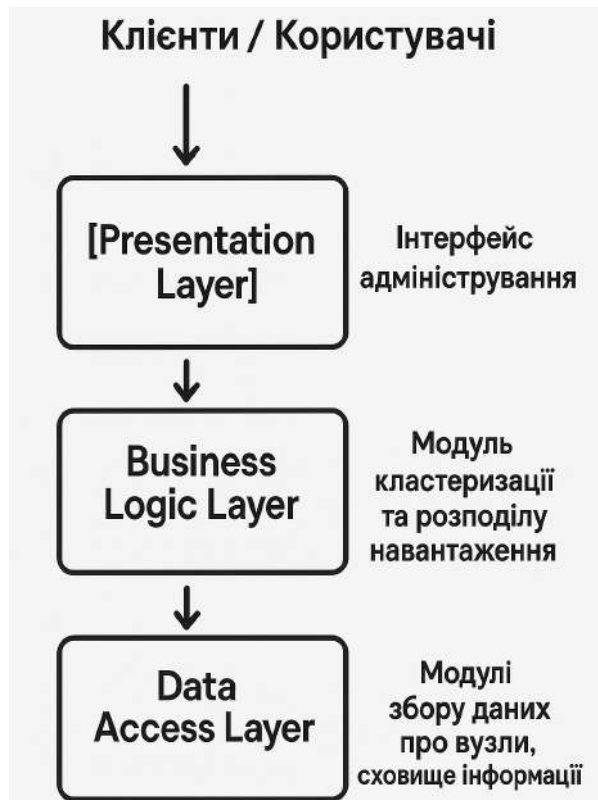


Рисунок 3.2 - Трирівнева (3-layer) архітектура

Рівень доступу до даних (Data Access Layer): збір інформації про стан кожного вузла (їхній поточний ступінь завантаженості, доступна пам'ять, пропускна здатність мережі, інші показники).

Рівень бізнес-логіки (Business Logic Layer): кластеризація вузлів за рівнем завантаженості, визначення оптимального вузла для розподілу нового або вже існуючого завдання, реалізація алгоритмів розподілу навантаження.

Рівень представлення (Presentation Layer): візуалізація даних (моніторинг стану вузлів, результатів кластеризації тощо), надання адміністративного інтерфейсу для налаштування політик розподілу та перегляду статистики.

Для ефективною реалізації обраного підходу використовують низку зразків проектування (design patterns), які допоможуть розв'язати типові завдання розподілених систем: гнучкість у зміні алгоритмів розподілу навантаження, динамічне додавання нових вузлів, організація збору та обробки даних про ресурси тощо.

Призначення: дозволяє інкапсулювати різні алгоритми розподілу

навантаження (наприклад, Random, Round Robin, Least Loaded, тощо) та динамічно обирати потрібну стратегію без зміни коду клієнта.

Використання в проєкті: для реалізації окремих методів розподілу навантаження залежно від завдань, які обробляються системою. Зокрема, можна швидко додати або модифікувати алгоритм, який опрацьовує завдання з урахуванням кластеризації (наприклад, обирає найменш завантажений кластер/вузол).

Призначення: організовує механізм сповіщення зацікавлених об'єктів про зміну стану іншого об'єкта.

Використання в проєкті: сервіс чи компонент, що здійснює моніторинг ресурсів вузлів, автоматично надсилатиме оновлення всім підписникам (наприклад, модулю кластеризації або бізнес-логіки). Це спростить облік актуальної інформації про навантаження і уникне постійних опитувань.

Singleton (Одинарний) - забезпечує існування лише одного екземпляра класу з глобальною точкою доступу.



Рисунок 3.3 - Спрощена діаграма послідовності

Використання в проєкті: для головного диспетчера або керуючого модуля розподілу навантаження, де потрібно зберігати глобальний стан (наприклад, статистику про всі вузли). Це спростить доступ до єдиної точки керування даними про стан системи.

Factory Method (Фабричний метод) або Abstract Factory (Абстрактна фабрика) - інкапсулює створення об'єктів, даючи змогу підкласам вирішувати, який саме клас буде створено.

Використання в проєкті: при додаванні нових типів вузлів або нових способів ініціалізації ресурсів. Наприклад, якщо з'являється необхідність підтримувати різні платформи або різні варіанти конфігурацій вузлів (з GPU, без GPU тощо), фабричний метод забезпечить розширюваність створення вузлів.

На рисунку 3.3 показано спрощену діаграму послідовності (Sequence Diagram), що ілюструє взаємодію основних компонентів (спостереження за вузлами, кластеризація, розподіл навантаження):

Node (кожен вузол) надсилає актуальні метрики своєму Monitor-компоненту.

Monitor отримує інформацію, обробляє її (застосовуючи, наприклад, патерн Observer, щоб повідомити зацікавлені підсистеми) та передає на LoadBalancer (модуль розподілу навантаження).

LoadBalancer виконує кластеризацію, обирає стратегію (патерн Strategy) та визначає, куди слід розподілити наступне завдання.

LoadBalancer дає команду розподілити конкретне завдання (distributeTask) на обраний вузол/кластер.

Вузол отримує завдання і виконує його (runTask), оновлюючи метрики, які знову потрапляють до Monitor-компонента.

3.4 Висновки до третього розділу

У даному розділі було обґрунтовано та реалізовано комплексний підхід до розподілу навантаження у мультикомп'ютерних системах загального призначення на основі кластеризації вузлів за рівнем їх ресурсної завантаженості.

Запропонована модель поєднує глибоке навчання (CNN та RNN) для точного аналізу просторово-часових характеристик навантаження, а також використовує підкріплювальне навчання (Reinforcement Learning) і гібридний алгоритм оптимізації (Hybrid Lyrebird Falcon) для динамічного формування кластерів і прийняття адаптивних рішень.

Побудована трирівнева архітектура системи із застосуванням шаблонів проєктування (Observer, Strategy, Singleton, Factory Method) дозволила забезпечити гнучкість, масштабованість та стійкість до відмов. Диспетчер та балансувальник виконують координацію дій між вузлами, гарантуючи рівномірний розподіл задач у реальному часі.

Класифікація вузлів за допомогою порогового методу дозволяє швидко виявляти перевантажені та недовантажені елементи системи й здійснювати їх ефективну взаємодію в рамках запропонованої системи.

Використання Python як основного інструменту реалізації надало можливість швидко і ефективно протестувати алгоритм розподілу навантаження, провести числовий аналіз результатів, реалізувати візуалізацію динаміки зміни навантаження та адаптувати систему до змінних умов без необхідності додаткових інструментів або середовищ розробки. Такий підхід демонструє практичну доцільність використання Python у задачах інтелектуального управління обчислювальними системами.

Таким чином, реалізований механізм дозволяє не лише оперативно реагувати на зміни навантаження, а й знижує енергоспоживання, покращує використання ресурсів, мінімізує затримки та забезпечує стабільну роботу мультикомп'ютерних систем у динамічному середовищі, зокрема в інфраструктурі розумного будинку.

4 ОЦІНКА ЕФЕКТИВНОСТІ МЕТОДУ РОЗПОДІЛУ НАВАНТАЖЕННЯ В МУЛЬТИКОМП'ЮТЕРНИХ СИСТЕМАХ

4.1 Визначення основних характеристик та експериментальні дослідження

Для оцінки ефективності запропонованого методу розподілу навантаження було здійснено серію експериментальних досліджень на тестовій мультикомп'ютерній системі. Система складалася з декількох автономних комп'ютерних вузлів, об'єднаних у єдину мережу. Експерименти проводилися з використанням спеціального програмного забезпечення, яке дозволяло імітувати різні ситуації навантаження та зміни в роботі вузлів.

Таким чином, проведені дослідження були спрямовані на оцінку потенціалу використання методу розподілу навантаження із метою:

1. Оптимізації використання ресурсів: забезпечення рівномірного навантаження на вузли для максимального використання доступних обчислювальних і енергетичних ресурсів.

2. Підвищення ефективності роботи системи: зниження часу відгуку та покращення загальної продуктивності через адаптивне розподілення навантаження.

3. Зменшення енергоспоживання: мінімізація надмірних витрат енергії шляхом рівномірного розподілу навантаження, що сприяє збереженню енергетичних ресурсів.

4. Забезпечення стабільності роботи системи: зменшення ризиків перевантаження вузлів і досягнення рівномірного навантаження на кожен з них, що підвищує загальну надійність системи.

Для оцінки ефективності застосування методу розподілу навантаження було використано набір метрик, які наведені у таблиці 4.1.

Таблиця 4.1 – Набір метрик, що були використані для оцінки методу розподілу навантаження

Метрика	Характеристика	Опис метрик
Час реакції системи	Продуктивність системи	Час, необхідний для обробки та виконання завдання після його надходження в систему. Визначає швидкодію системи.
Коефіцієнт балансування навантаження	Рівномірність розподілу навантаження	Показник, що характеризує рівень рівномірного розподілу навантаження між вузлами. Чим вищий коефіцієнт, тим ефективніше використання ресурсів.
Рівень використання ресурсів	Ефективність використання ресурсів	Співвідношення фактично використаних ресурсів до загальної кількості доступних ресурсів на кожному вузлі.
Кількість міграцій задач	Адаптивність системи	Кількість разів, коли завдання переміщуються між вузлами для оптимізації навантаження. Більш висока кількість міграцій може свідчити про неефективність методу.
Ступінь зменшення енергоспоживання	Економія енергії	Показник, що визначає зменшення енергоспоживання завдяки ефективному розподілу навантаження між вузлами.

4.2 Архітектура програмної реалізації інформаційної технології

Застосування методу розподілу навантаження в мультикомп'ютерних системах згідно кластеризації вузлів за рівнем завантаженості ресурсів може бути ефективно інтегроване у сферу розумних будинків. Дана інтеграція дозволяє підвищити енергоефективність, надійність та забезпечити стабільну роботу компонентів «розумного будинку».

Запропонована методика може бути реалізована шляхом створення єдиної мультикомп'ютерної мережі з вузлів (наприклад, мікрокомп'ютерів Raspberry Pi), кожен з яких відповідає за окремий блок задач системи розумного будинку:

- управління освітленням;
- системи безпеки та відеонагляду;
- регулювання клімату (опалення, кондиціонування);
- датчики руху, вологості та температури;
- управління побутовими пристроями та мультимедіа.

Інтелектуальний розподіл навантажень дає можливість ефективно управляти піковими ситуаціями, коли відбувається одночасне залучення великої кількості пристроїв (наприклад, ввечері або вранці). Метод кластеризації забезпечує адаптивне перерозподілення завдань, запобігаючи можливим зниженням продуктивності чи аварійним зупинкам обладнання.

Запропонована архітектура є кластеризованою мультикомп'ютерною системою зображена на рисунку 4.1, яка забезпечує розподілене управління ресурсами компонентів розумного будинку.

Центральним елементом архітектури виступає інтелектуальна система управління, яка реалізує розроблений автором алгоритм розподілу навантаження між вузлами згідно кластеризації за рівнем їхньої завантаженості ресурсів.

У системі «розумного» будинку організовано дві черги повідомлень: оперативну, призначену для подій, що використовуються для формування нових моделей штучних нейронних мереж, та історичну, яка зберігає дані подій у

документоорієнтованій базі для їхнього подальшого аналізу та виявлення взаємозв'язків із застосуванням штучного інтелекту.



Рисунок 4.1 - Кластеризована мультикомп'ютерна архітектура розподіленого управління ресурсами системи розумного будинку

Система складається з декількох автономних вузлів, кожен із яких утворює окремий кластер і відповідає за управління певною групою пристроїв будинку. Пристрої поділяються на функціональні категорії, такі як освітлення, відеонагляд, кліматичні системи (опалення, вентиляція, кондиціонування повітря), системи безпеки та контролю доступу, побутова техніка і сенсорні пристрої (датчики руху, температури, вологості тощо).

Інтелектуальна система управління здійснює постійний моніторинг поточного стану вузлів, оцінюючи їх завантаженість і визначаючи необхідність перерозподілу завдань між ними. При перевищенні допустимого рівня навантаження на окремому вузлі система автоматично виконує динамічну міграцію задач на менш завантажені вузли. Це дозволяє підтримувати стабільну роботу всієї системи, уникати перевантажень і аварійних ситуацій, підвищувати енергоефективність та забезпечувати надійність роботи компонентів будинку.

Мережева архітектура розробленої системи «розумного» будинку охоплює сенсори, побутові пристрої, клієнтські контролери та віддалені центри зберігання і обробки даних. Сенсори та побутова техніка підключаються до контролерів, які розміщені безпосередньо в приміщеннях житлового будинку. З'єднання пристроїв із контролерами реалізується за допомогою додаткових модулів, що підтримують віддалене керування та забезпечують комунікацію через бездротові протоколи передачі інформації.

Керування побутовими пристроями й отримання даних із сенсорів здійснюється шляхом передачі команд через інфрачервоне випромінювання, а також із використанням популярних бездротових стандартів зв'язку, таких як Wi-Fi, Bluetooth та ZigBee. Ці протоколи широко застосовуються у сучасних сенсорах та підтримуються більшістю вбудованих і додаткових модулів мікроконтролерів та мікрокомп'ютерів для бездротового зв'язку й інфрачервоної передачі інформації. Для реалізації запропонованого рішення використано сервіс Amazon API Gateway, який дозволяє створювати RESTful API на базі HTTP або REST API. Клієнтська частина взаємодіє з сервером шляхом надсилання HTTP-запитів (методів GET і POST) на попередньо визначені URL-адреси, отримуючи інформацію про новостворені моделі штучних нейронних мереж або передаючи дані про події, що відбуваються у приміщеннях.

Сервіс Amazon API Gateway забезпечує ефективний розподіл запитів між серверами, обираючи найменш завантажені або територіально найближчі вузли до клієнта. Це дозволяє підтримувати рівномірний розподіл обчислювального навантаження, мінімізувати затримки передачі інформації та максимально пришвидшити отримання відповіді від сервера. Таким чином, клієнтська частина швидко отримує результати обробки повідомлень про нові показники датчиків.

Крім того, API Gateway підтримує одночасне використання кількох версій одного API. Це дозволяє клієнтам продовжувати взаємодіяти зі старішими версіями навіть після виходу нових. Також API Gateway надає гнучке управління стадіями розгортання API (альфа, бета, продакшн). Кожну зі стадій можна налаштувати

окремо для взаємодії з різними серверними кінцевими точками відповідно до встановлених параметрів API.

Об'єкти даних передаються від клієнтських частин до серверної через мережу Інтернет. Особливістю систем «розумного» будинку є велика кількість подій, що надсилаються на сервер від численних клієнтів, які повинні бути опрацьовані в режимі реального часу. Для організації черги обробки подій та оптимізації навантаження на обчислювальні ресурси в системі використовується сервіс Amazon Simple Queue Service (SQS). Цей сервіс реалізує архітектурний патерн брокера повідомлень у розподілених системах, що дозволяє зберігати повідомлення в черзі на обробку під час пікових навантажень без необхідності додаткового виділення обчислювальних ресурсів для миттєвого опрацювання.

Кожна подія, що зберігається, буде оброблена при звільненні обчислювальних ресурсів відповідно до порядку її надходження в чергу. Сервіс Amazon SQS гарантує доставку повідомлень їх обробникам, підписаним на оновлення черги, згідно з політикою одноразової обробки повідомлень, тобто кожне повідомлення доставляється лише один раз і зберігається доступним до моменту його обробки та видалення підписником. Оскільки черги працюють за принципом «першим отримано — першим оброблено», забезпечується точний порядок надсилання та отримання повідомлень.

У системі «розумного» будинку створено дві черги повідомлень: оперативна черга для подій, які використовуються для створення нових моделей штучних нейронних мереж, і черга історичних даних для подій, які записуються в документоорієнтовану базу даних для подальшого аналізу і пошуку взаємозалежностей за допомогою штучного інтелекту. Розподіл подій між чергами здійснюється за типом події через окрему функцію AWS Lambda, яка активується після отримання нової події, не споживаючи обчислювальних ресурсів в періоди відсутності нових подій з боку клієнтів. Схема розподілу подій між чергами представлена на рисунку 4.2.

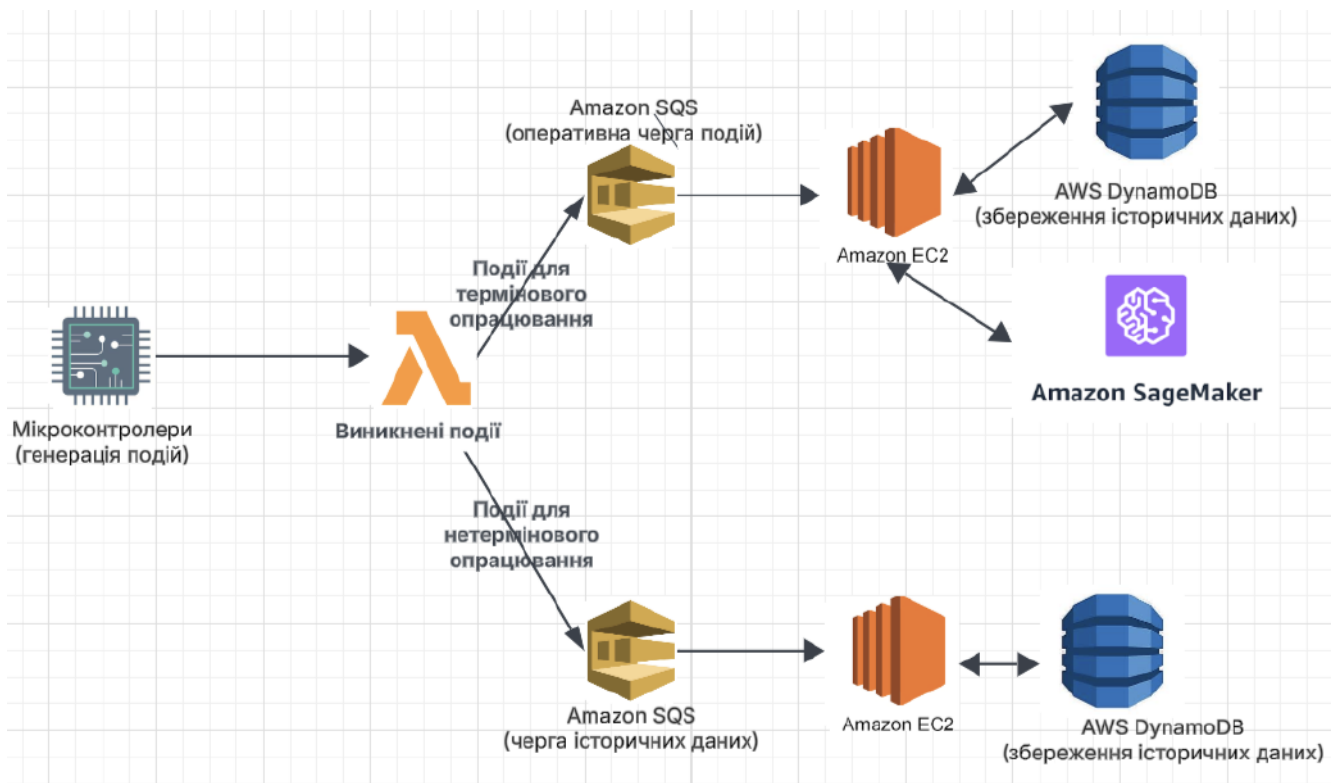


Рисунок 4.2 - Схема розподілу навантаження подій між чергами

У рамках реалізації методу оптимізації виконання завдань з використанням реплікації та цільової функції оптимізації у багатопроцесорній вбудованій системі було розроблено спеціалізовану програму на C++. Використання саме цієї мови обумовлено її високою продуктивністю, гнучкістю управління пам'яттю та можливостями тонкого контролю за багатопотоковими обчисленнями, що є критично важливим для мультикомп'ютерних систем.

Програма враховує:

- поточне завантаження ресурсів вузлів (процесор, пам'ять, мережеві ресурси);
- можливість реплікації завдань для підвищення надійності обчислень та скорочення часу їх виконання;
- цільову функцію оптимізації, яка мінімізує сумарний час виконання завдань при збалансованому використанні ресурсів;
- динамічну кластеризацію вузлів на основі рівня їх завантаженості;

– взаємодію між вузлами через стандартизовані інтерфейси для обміну інформацією про стан ресурсів.

Програма реалізує:

– алгоритм кластеризації вузлів за критеріями завантаженості (з використанням методів машинного навчання або евристичних правил);

– розподіл завдань між кластерами та всередині кластерів на основі заданої функції оптимізації;

– механізм реплікації для завдань з високою критичністю або пріоритетом.

– моніторинг ресурсів у реальному часі, що забезпечує адаптивність розподілу навантаження.

– систему логування і аварійного відновлення у випадку збоїв окремих вузлів.

Розроблений алгоритм кластерного розподілу навантаження (зображений на рисунку 4.3) має на меті забезпечення ефективного використання обчислювальних ресурсів у мультикомп'ютерній системі шляхом перенаправлення задач від перевантажених вузлів до менш завантажених. Його робота базується на поетапному аналізі поточного стану ресурсів системи та виконанні обчислень, що дозволяють визначити оптимальну стратегію розподілу.

Початковим етапом є ініціалізація даних. Збираються вихідні параметри кожного вузла, а саме: ступінь завантаження центрального процесора (CPU) та оперативної пам'яті (RAM). Ці два показники у подальшому використовуються для визначення інтегрального рівня навантаження на кожному вузлі за допомогою формули, яка враховує вагу кожного ресурсу. Формально, інтегральне навантаження L_i обчислюється як зважена сума відносного навантаження на CPU та RAM, що дозволяє одержати уніфікований показник для порівняння між вузлами.

Після розрахунку L_i система здійснює класифікацію вузлів. Вузли, навантаження яких перевищує заданий пороговий рівень, потрапляють до групи донорів — джерел перевантаження, які потребують зниження обчислювального навантаження. Водночас вузли з низьким значенням L_i формують групу приймачів, що мають достатній резерв ресурсів для прийняття додаткових задач.



Рисунок 4.3 - Блок-схема алгоритму кластерного розподілу навантаження

Далі відбувається пропорційний розподіл задач: кожен приймач отримує ту частину задач, яка відповідає його частці доступного ресурсу відносно загального резерву всіх приймачів. Таким чином, балансування здійснюється справедливо й оптимально відповідно до поточного стану всієї системи.

Після завершення розподілу, обчислювальні навантаження вузлів коригуються: у донорах значення L_i зменшуються відповідно до кількості переданих задач, а в приймачів, навпаки, збільшуються згідно з обсягом

прийнятого навантаження. Це дозволяє отримати оновлену картину розподілу ресурсів, яка відображає ефект від балансування.

Для підтвердження ефективності алгоритму проводиться візуалізація результатів. Зазвичай використовуються графіки, що демонструють рівень навантаження до та після розподілу. Це дає змогу оцінити результат як візуально, так і кількісно. Оцінка виконується шляхом аналізу середнього значення навантаження та його дисперсії, що дозволяє визначити ступінь вирівнювання системи.

Завершальним етапом є прийняття рішення про завершення процесу або його повторне виконання в наступному циклі. У разі реалізації системи в режимі реального часу, алгоритм працює циклічно — з певною періодичністю або на основі подієвих тригерів, що дозволяє постійно адаптуватися до змін у навантаженні.

4.3 Експериментальне дослідження ефективності методу розподілу навантаження в умовах інтенсивної роботи розумного будинку

У сучасних розумних будинках, що функціонують як мультикомп'ютерні системи, використовується широкий спектр інтегрованих пристроїв — від сенсорних модулів та систем контролю клімату до відеоспостереження з елементами штучного інтелекту. Така інфраструктура передбачає наявність кількох обчислювальних вузлів, які обробляють дані локально або передають їх на периферійні чи хмарні ресурси.

Під час пікових періодів активності, наприклад, у вечірній час, коли всі мешканці перебувають вдома, інтелектуальна система будинку зазнає підвищеного обчислювального навантаження. У цей період активуються одночасно декілька підсистем.

Системи відеоспостереження виконують обробку зображень у реальному часі, використовуючи алгоритми розпізнавання облич та виявлення аномалій. Голосові інтерфейси та цифрові помічники відповідають на запити користувачів.

Системи автоматичного кліматичного регулювання зчитують показники температури, вологості, рівня CO₂ та здійснюють відповідні регулювання. Обчислювальні вузли запускають моделі машинного навчання для адаптації поведінки системи до потреб мешканців (наприклад, прогнозування дій, розкладів тощо). Працюють модулі локального збереження та обробки даних, зокрема для потреб приватності.

Внаслідок цього формується нерівномірний розподіл обчислювального навантаження між вузлами системи: деякі вузли (наприклад, ті, що відповідають за відеоаналітику чи машинне навчання) функціонують на межі своїх ресурсних можливостей, тоді як інші вузли (що обробляють менш критичні або періодичні дані) залишаються частково або повністю недовантаженими.

Подібна ситуація не лише призводить до зниження ефективності функціонування системи, але й може викликати затримки в обробці інформації, погіршення якості сервісу, а в окремих випадках — часткову деградацію функціональності (наприклад, втрата кадрів відеопотоку або запізніла реакція системи на критичні сигнали від сенсорів).

Кожен із вузлів має локальний процесор, оперативну пам'ять, програмне середовище та можливість обмінюватися даними з іншими вузлами через локальну мережу (наприклад, Wi-Fi або Ethernet) (таблиця 4.2).

На початковому етапі дослідження необхідно зафіксувати поточний стан обчислювальних ресурсів кожного вузла мультикомп'ютерної системи розумного будинку. Це дозволяє об'єктивно оцінити рівень завантаженості ресурсів, визначити ступінь нерівномірності навантаження між вузлами та виявити критичні зони, що потребують балансування.

Кожен обчислювальний вузол характеризується двома основними показниками навантаження:

- процентне завантаження центрального процесора (CPU, %);
- процентне використання оперативної пам'яті (RAM, %).

Таблиця 4.2 - Основні елементи архітектури

Вузол	Назва модуля	Основне навантаження
A	Відеоаналітика	Аналіз відеопотоків у реальному часі
B	Голосовий інтерфейс	Розпізнавання мови, TTS, NLP
C	Модуль сенсорного моніторингу	Збір і передача показників сенсорів
D	Клімат-контроль	Обробка даних про температуру, CO ₂
E	Обчислювальний модуль (ML)	Виконання моделей штучного інтелекту

Дані показники збираються шляхом моніторингу системних ресурсів кожного вузла у реальному часі. Для моделювання сценарію максимального навантаження використовуються зразкові дані, що відображають пікові вечірні навантаження на розумний будинок.

$$L_i = \alpha \cdot \frac{CPU_i}{100} + \beta \cdot \frac{RAM_i}{100}, \quad (4.1)$$

де L_i – інтегральний коефіцієнт завантаження i -го вузла;

CPU_i – рівень використання центрального процесора вузлом i у відсотках;

RAM_i – рівень використання оперативної пам'яті вузлом i у відсотках;

α – вага процесорного навантаження (наприклад, $\alpha = 0.6$);

β – вага навантаження пам'яті (наприклад, $\beta=0.4$).

Ось інтегральний графік на рисунку 4.4 показує навантаження вузлів L_i системи до балансування. Він враховує як завантаження процесора (CPU), так і оперативної пам'яті (RAM) з відповідними вагами: $\alpha=0.6$ для CPU і $\beta=0.4$ для RAM. Чим вище стовпчик — тим більше загальне навантаження на вузол. Добре видно,

що вузли А та Е мають критично високі значення L_i , тоді як В, С та D завантажені мінімально. Блакитний стовпчик — завантаження CPU, помаранчевий стовпчик — завантаження оперативної пам'яті (RAM).

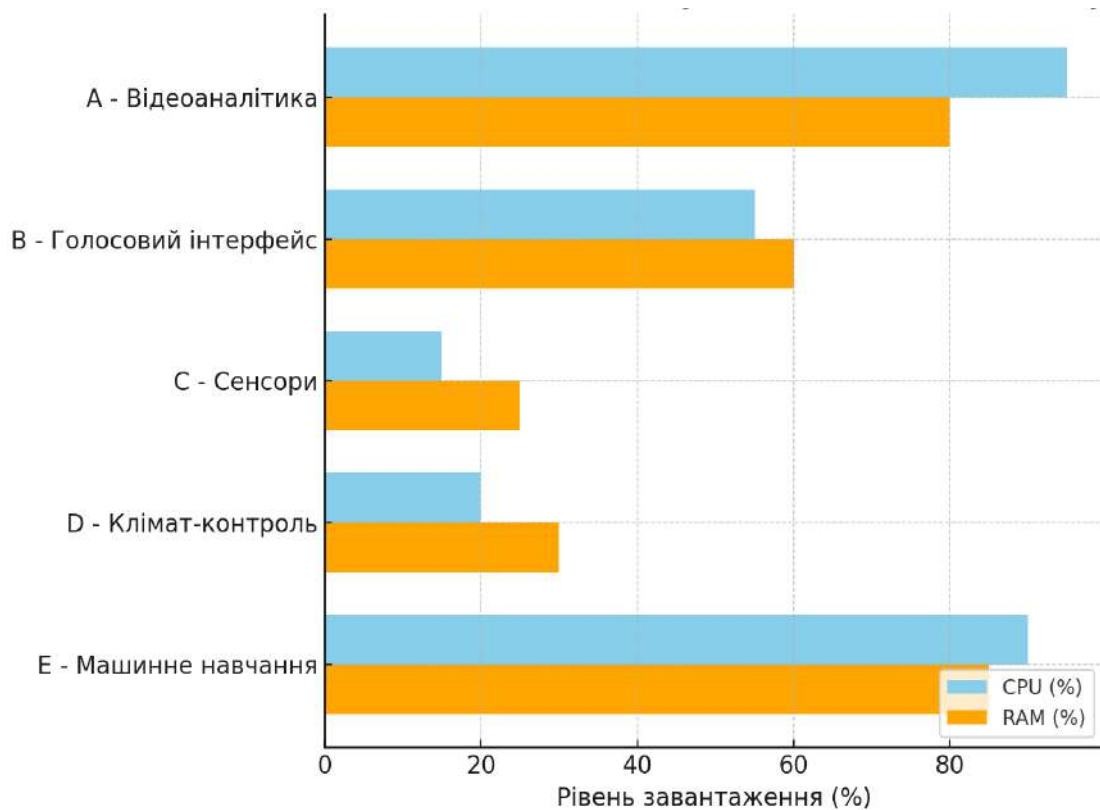


Рисунок 4.4 - Горизонтальна порівняльна діаграма навантаження вузлів

Таблиця 4.3 – Показники завантаження вузлів системи до балансування

Вузол	CPU завантаження (%)	RAM завантаження (%)
А	95	80
В	55	60
С	15	25
D	20	30
Е	90	85

Оскільки процесорне навантаження та використання пам'яті можуть мати різну вагу впливу на продуктивність системи, вводиться інтегральний показник завантаження вузла — коефіцієнт завантаження L_i (рис. 4.5).

Розглянемо приклад розрахунку інтегрального навантаження для вузла А:

$$L_A = 0.6 \cdot 0.95 + 0.4 \cdot 0.80 = 0.57 + 0.32 = 0.89. \quad (4.2)$$

Таким чином, вузол А має коефіцієнт завантаження $L_A = 0.89$, що свідчить про його критичне навантаження.

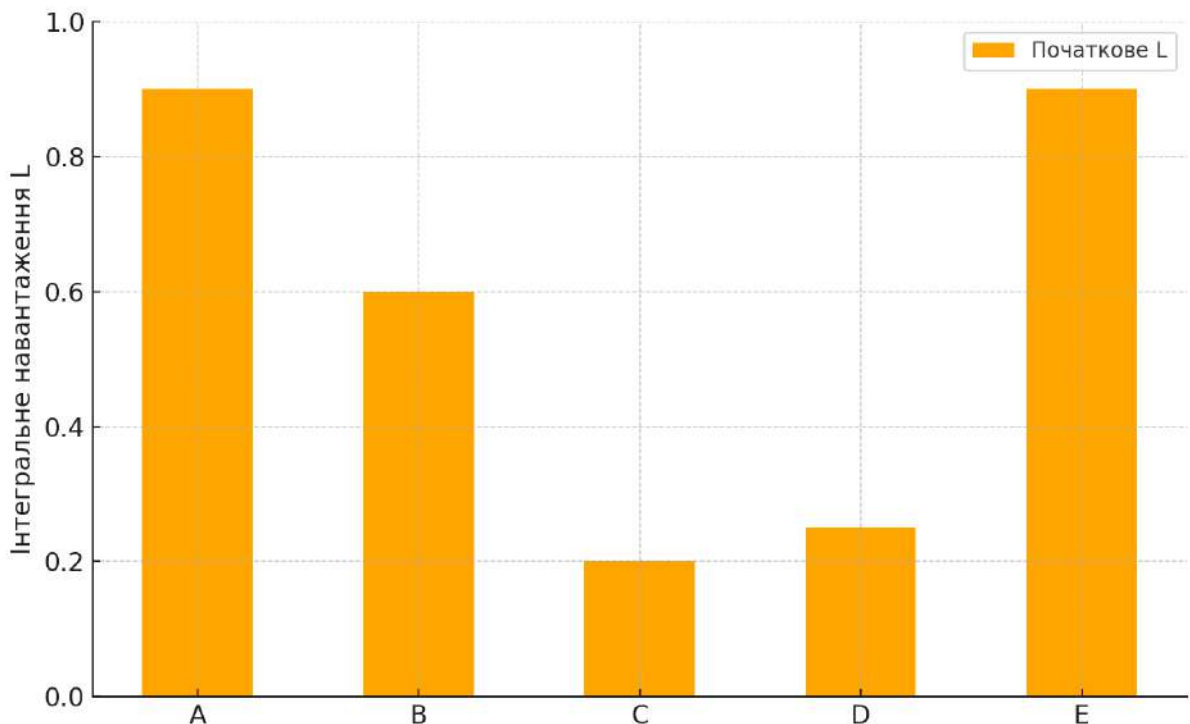


Рисунок 4.5 - Значення початкового інтегрального навантаження L для вузлів

З аналізу обчислених значень видно, що вузли А та Е мають найбільші коефіцієнти навантаження ($L_i > 0.85$), що підтверджує їх критичний стан. Вузли В, С та D мають мінімальні значення інтегрального навантаження, що свідчить про можливість залучення їх ресурсів для перерозподілу обчислювальних задач.

Аналогічні обчислення виконуються для всіх інших вузлів системи. Отримані результати представлені в таблиці 4.4

Таблиця 4.4 - Інтегральні коефіцієнти завантаження вузлів

Вузол	Назва модуля	L_i (ступінь завантаження)
A	Відеоаналітика	0.89
B	Голосовий інтерфейс	0.57
C	Модуль сенсорного моніторингу	0.19
D	Клімат-контроль	0.26
E	Обчислювальний модуль (ML)	0.90

З аналізу обчислених значень видно, що вузли А та Е мають найбільші коефіцієнти навантаження ($L_i > 0.85$), що підтверджує їх критичний стан. Вузли В, С та D мають мінімальні значення інтегрального навантаження, що свідчить про можливість залучення їх ресурсів для перерозподілу обчислювальних задач.

Для балансування навантаження у системі було запропоновано виконати міграцію задач наступним чином:

- Частину задач відеоаналітики, які виконуються на вузлі А, перенести на вузли С та D.
- Частину обчислень моделі машинного навчання з вузла Е перенести на вузол С.

Таким чином, менш завантажені вузли отримають додаткові задачі, що дасть змогу знизити навантаження на перевантажені вузли і вирівняти загальне навантаження в системі.

У нашому випадку розподіл завдань між вузлами не просто випадковий — ми застосували балансувальний принцип на основі пропорційного перерозподілу задач, враховуючи:

- поточний рівень завантаження вузлів;
- наявний вільний ресурс на кожному вузлі,
- відсоток задач, які можна перенести без втрати продуктивності.

Основною задачею є перерозподілити надлишок навантаження з вузлів з найвищим L_i на вузли з найнижчим L_i , пропорційно до їхньої доступної потужності.

Для задачі, яку потрібно мігрувати, обсяг міграції на вузол j визначається за формулою:

$$\Delta W_{i \rightarrow j} = W_{\text{надлишок}} \times \frac{R_j}{\sum_{k \in S} R_k}, \quad (4.3)$$

де $W_{\text{надлишок}}$ — кількість навантаження, яку потрібно зняти з вузла i ;

R_j — ресурсна ємність (запас потужності) вузла j ;

S — множина вузлів-кандидатів на прийом задач;

$\sum_{k \in S} R_k$ — сумарний ресурс вузлів-кандидатів.

Ідея - чим більше у вузла вільних ресурсів, тим більшу частку задач йому можна передати.

Наприклад, вузол А мав $L_A = 0.89$, а бажаний рівень — близько 0.6. Тобто потрібно було знизити навантаження на:

$$W_{\text{надлишок}(A)} = (0.89 - 0.6) = 0.29. \quad (4.4)$$

Аналогічно для вузла Е:

$$W_{\text{надлишок}(E)} = (0.90 - 0.66) = 0.24. \quad (4.5)$$

Перед балансуванням вузли з найменшим L_i були С і D. Вузол С: доступний запас ресурсів $\approx 1 - 0.19 = 0.81$. Вузол D: доступний запас ресурсів $\approx 1 - 0.26 = 0.74$

Скільки задач отримує кожен вузол – пропорційно їх ресурсним запасам:

$$\Delta W_{A \rightarrow C} = 0.29 \times \frac{0.81}{0.81 + 0.74} = 0.1515, \quad (4.6)$$

$$\Delta W_{A \rightarrow D} = 0.29 \times \frac{0.74}{0.81 + 0.74} = 0.1385. \quad (4.7)$$

Вузол С отримає приблизно 15.15% одиниць навантаження. Вузол D отримає приблизно 13.85% одиниць навантаження. Разом вони повністю покриють надлишок у 0.29 одиниці, який був на вузлі А.

Для вузла Е інтегральне навантаження вузла становило $L_E = 0.9$. Таким чином, надлишок навантаження, який потрібно зняти з вузла Е, дорівнює:

$$W_{\text{надлишок}(E)} = L_{\text{початкове}(E)} - L_{\text{цільове}(E)} = 0.90 - 0.66 = 0.24. \quad (4.8)$$

Отже, потрібно зменшити навантаження вузла Е на 0.24 одиниці. Нове навантаження вузла С:

$$L_{\text{нове}(C)} = L_{\text{початкове}(C)} + W_{\text{надлишок}(E)} = 0.19 + 0.24 = 0.43. \quad (4.9)$$

Ось візуальне резюме міграції навантаження всіх вузлів - показано в таблиці 4.5 початковий рівень інтегрального навантаження L_i .

Таблиця 4.5 - навантаження всіх вузлів після розподілу

Вузол	L_i (поч.)	Надлишок	L_i (після)	CPU завантаження (%)	RAM завантаження (%)
А	0.89	-0.29	0.58	60	50
В	0.57	0.0	0.57	55	60
С	0.19	0.39149	0.5	50	50
D	0.26	0.1385	0.43	40	25
Е	0.90	-0.24	0.66	70	60

Відображено (рис. 4.6), який надлишок задач було передано або прийнято кожним вузлом. І фінальні значення L_i після балансування системи.

Жовті стовпчики - початковий стан системи, де вузли А і Е були перевантажені, помаранчеві стовпчики - стан після балансування, де навантаження розподілено більш рівномірно.

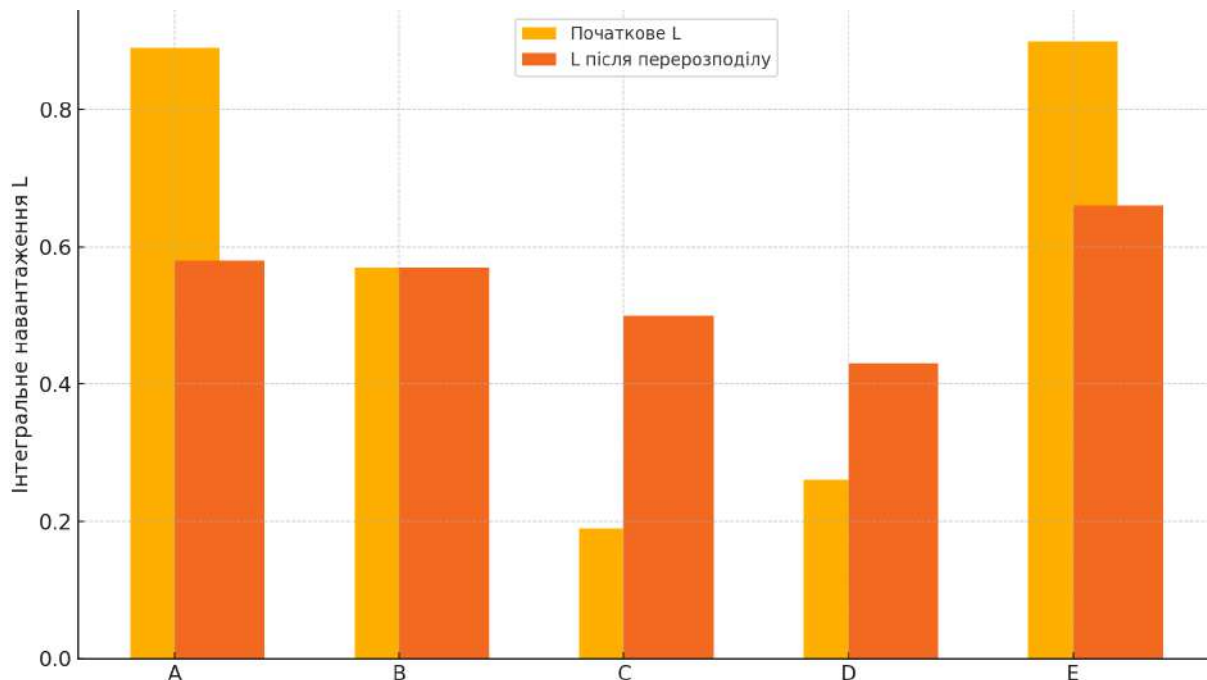


Рисунок 4.6 - порівняння інтегрального навантаження вузлів до і після розподілу

В результаті впровадження запропонованого алгоритму міграції обчислювальних задач між вузлами мультикомп'ютерної системи розумного будинку досягаються такі ключові ефекти.

Зниження навантаження на критичні вузли. Після часткового перенесення завдань обробки відео з вузла А та обчислень машинного навчання з вузла Е, їх інтегральні коефіцієнти навантаження L_A та L_E суттєво зменшуються. Це дозволяє уникнути перевищення критичних порогів використання ресурсів і запобігає ризикам деградації їх продуктивності.

Більш ефективне використання ресурсів малонавантажених вузлів. Вузли С (сенсорний модуль) та D (кліматичний контроль), які до цього моменту працювали із мінімальним навантаженням, після міграції приймають на себе додаткові задачі. Це забезпечує більш збалансоване використання доступних обчислювальних потужностей системи.

Досягнення рівномірнішого розподілу навантаження в системі. Завдяки оптимізації розподілу задач серед вузлів вдалося значно знизити дисперсію інтегральних коефіцієнтів навантаження L_i . Це сприяє стабільнішій і передбачуваній роботі розумного будинку навіть у пікові періоди активності.

Підвищення надійності функціонування системи. Вирівнювання обчислювального навантаження суттєво зменшує ймовірність виникнення відмов або деградації продуктивності через локальні перевантаження окремих вузлів. Система набуває здатності до більш ефективного самовідновлення та продовження роботи навіть у випадку збільшення кількості запитів або неочікуваного росту обсягу оброблюваних даних.

Таким чином, впровадження механізму динамічного перерозподілу навантаження на основі кластеризації вузлів дозволяє підвищити ефективність використання ресурсів, зменшити ризики критичних збоїв та забезпечити стабільну роботу розумного будинку в умовах змінного навантаження.

4.4 Особливості практичної реалізації

Для реалізації математичної моделі кластерного аналізу та ефективного розподілу обчислювального навантаження між вузлами мультикомп'ютерної системи в середовищі розумного будинку було обрано мову програмування Python. Такий вибір продиктований поєднанням технічних, практичних і методологічних факторів, які в сукупності забезпечують високу гнучкість, масштабованість і простоту інтеграції з сучасними бібліотеками машинного навчання, обробки даних та візуалізації результатів.

Python є високорівневою мовою з виразним синтаксисом, що значно прискорює розробку прототипів та полегшує супровід коду. Завдяки наявності потужних бібліотек, таких як NumPy, Pandas, scikit-learn, Dask, Joblib та Matplotlib, Python дозволяє ефективно реалізовувати складні математичні обчислення, алгоритми машинного навчання, а також графічне представлення отриманих результатів. Більше того, його активне використання у галузях Data Science та HPC

(High Performance Computing) робить Python ідеальним інструментом для дослідження та експериментального моделювання в межах інтелектуальних систем.

Етап 1 - збір і формалізація початкових даних. На цьому етапі здійснювалося очищення, нормалізація та трансформація вхідної інформації з сенсорів або журналів системи управління розумним будинком. Було реалізовано попередню фільтрацію шуму, відновлення втрачених даних і приведення їх до уніфікованого формату. Також проводиться фіксація вхідних параметрів - поточне завантаження процесора (CPU) кожного вузла у відсотках, поточне завантаження оперативної пам'яті (RAM) у відсотках.

```
import numpy as np
# Завантаження ресурсів (у відсотках)
cpu = np.array([95, 55, 15, 20, 90]) # A, B, C, D, E
ram = np.array([80, 60, 25, 30, 85])
```

Етап 2 - розрахунок інтегрального навантаження L_i . Результат — масив інтегрального навантаження L_i для кожного вузла. Важливим кроком було виділення релевантних ознак, які найкраще характеризують навантаження на кожен вузол. Це могли бути такі показники, як інтенсивність запитів, частота активації пристроїв, кількість підключених користувачів.

```
alpha, beta = 0.6, 0.4
L_i = alpha * (cpu / 100) + beta * (ram / 100)
print("Інтегральне навантаження L:", np.round(L_i, 2))
```

Етап 3 - визначення перевантажених та резервних вузлів. Вузли, у яких $L_i > 0.6$, визначаються як перевантажені. Інші — як вузли з потенціалом для прийому задач. За допомогою алгоритмів кластерного аналізу, вузли були об'єднані у кластери за подібністю навантаження або функціонального призначення. Це дозволило логічно згрупувати обчислювальні ресурси та оптимізувати подальший розподіл задач.

```
L_target = 0.6
excess = np.maximum(0, L_i - L_target)
print("Надлишок задач:", np.round(excess, 2))
```

Етап 4 - розрахунок доступного ресурсу у приймаючих вузлів. Ми знаходимо доступний ресурс кожного вузла. Після формування кластерів проводився аналіз їхньої внутрішньої структури, щільності та гомогенності. Визначались потенційно перевантажені або малоефективні вузли, що дозволяло заздалегідь виявити "вузькі місця" у системі.

```
available_capacity = np.where(L_i < L_target, 1 - L_i, 0)
print("Доступний ресурс вузлів:", np.round(available_capacity,
2))
```

Етап 5 - пропорційний розподіл задач. Згідно з отриманими кластерними структурами, було запрограмовано логіку балансування навантаження. Вона забезпечувала динамічне перенаправлення задач від перевантажених вузлів до менш завантажених у межах одного кластеру або між кластерами. Розподіляємо надлишкове навантаження вузла А (0.29) між вузлами С і D згідно їхньої доступної ємності.

```
excess_A = 0.29
res_C, res_D = 0.81, 0.74
total_res = res_C + res_D

delta_C = excess_A * (res_C / total_res)
delta_D = excess_A * (res_D / total_res)
print(f"З вузла А: до С -> {delta_C:.4f}, до D ->
{delta_D:.4f}")
```

Аналогічно — повне перенесення задач з вузла Е на вузол С:

```
excess_E = 0.24
delta_E_C = excess_E # усе йде на С
```

Етап 6 - оновлення навантаження. Було створено графіки, теплові карти та інтерактивні діаграми для наочного представлення структури кластерів, змін навантаження та ефективності перерозподілу. Після міграції задач оновлюємо значення навантаження для кожного вузла:

```
L_i_new = L_i.copy()
L_i_new[0] -= (delta_C + delta_D) # А
L_i_new[2] += (delta_C + delta_E_C) # С
```

```
L_i_new[3] += delta_D # D
L_i_new[4] -= delta_E_C # E
```

Етап 7 - оцінка ефективності. Розраховується середнє значення та дисперсія навантаження до і після:

```
print(f"Середнє L (до): {np.mean(L_i):.3f}, дисперсія:
{np.var(L_i):.3f}")
print(f"Середнє L (після): {np.mean(L_i_new):.3f}, дисперсія:
{np.var(L_i_new):.3f}")
```

4.4 Висновки до четвертого розділу

У цьому розділі було проведено експериментальну перевірку ефективності запропонованого методу кластерного розподілу навантаження в умовах мультикомп'ютерної системи розумного будинку. Було змодельовано ситуацію пікового навантаження, за якої деякі вузли працювали у критичних режимах, тоді як інші залишалися недовантаженими.

Запропонований алгоритм дозволив виявити перевантажені вузли та ефективно перерозподілити задачі на менш навантажені компоненти. В результаті вдалося знизити інтегральне навантаження критичних вузлів (А та Е), покращити рівномірність розподілу задач і підвищити загальну ефективність використання ресурсів системи.

Використання Python для реалізації методу забезпечило зручність моніторингу, аналізу та візуалізації даних. Запропонована система показала високу адаптивність, стабільність та енергоефективність, що підтверджує практичну доцільність її впровадження в інфраструктуру розумного будинку.

ВИСНОВКИ

У кваліфікаційній роботі магістра, за результатами виконаних теоретичних та практичних досліджень розроблено новий метод балансування навантаження який на відміну від відомих використовує кластеризацію вузлів за рівнем їхнього завантаження, що дозволяє знизити енергоспоживання, підвищити рівномірність розподілу задач та зменшити затримки в мультикомп'ютерних системах.

1. Проведено всебічний аналіз існуючих методів балансування навантаження в мультикомп'ютерних системах, що дозволило виокремити їх сильні та слабкі сторони, а також обґрунтувати необхідність пошуку нових адаптивних рішень.

2. Розроблено метод балансування навантаження який на відміну від відомих використовує кластеризацію вузлів за рівнем їхнього завантаження, що дозволяє знизити енергоспоживання, підвищити рівномірність розподілу задач та зменшити затримки в мультикомп'ютерних системах.

3. Проведено оцінку ефективності балансування навантаження, підвищено ефективність розподілу задач між вузлами, зменшено ймовірність затримок окремих компонентів і покращити загальну стабільність роботи системи;

4. Досліджено результати моделювання та експериментального тестування, підтвердилось доцільність використання розробленого методу, продемонструвавши покращення показників продуктивності, рівномірності розподілу задач та зменшення затримок в мультикомп'ютерних системах.

ПЕРЕЛІК ДЖЕРЕЛ ПОСИЛАНЬ

1. Tanenbaum A. S., Van Steen M. *Distributed Systems: Principles and Paradigms*. Pearson, 2021. 672 p.
2. Smith J. *High-Performance Distributed Systems*. Springer, 2020. 498 p.
3. Foster I., Kesselman C. *The Grid: Blueprint for a New Computing Infrastructure*. Elsevier, 2021. 696 p.
4. Ghosh S. *Applications of Distributed Computing Systems*. Wiley, 2019. 432 p.
5. Microsoft Azure. *What is cloud computing*. URL: <https://azure.microsoft.com/en-in/resources/cloud-computing-dictionary/what-is-cloud-computing/#uses> (дата звернення: 17.11.2024).
6. OnBiz. *Cloud computing models*. URL: <https://onbiz.biz/cloud-computing-models/> (дата звернення: 23.11.2024).
7. Сімоненко В. П., Дехтярук М. Т., Забара С. С. Програмне забезпечення комп'ютерних мереж: підручник. Київ : Університет «Україна», 2014. 220 с.
8. Харріс Д. М. Цифрова схемотехніка та архітектура комп'ютера. Морган Кауфман, 2013.
9. Alon N. Distributed Computing and Scalability Challenges. *ACM Transactions*. 2020. Vol. 38, No. 4. P. 101–115.
10. Dean J., Ghemawat S. MapReduce: Simplified Data Processing on Large Clusters // *Communications of the ACM*. 2020. Vol. 63, No. 1. P. 22–30.
11. Simplilearn. *What is Hadoop*. URL: <https://www.simplilearn.com/tutorials/hadoop-tutorial/what-is-hadoop> (дата звернення: 30.11.2024).
12. Amazon Web Services. *What is Apache Spark?* URL: https://aws.amazon.com/what-is/apache-spark/?nc1=h_ls (дата звернення: 07.12.2024).
13. Jassy A. Building Scalable Applications with AWS. *AWS Whitepapers*, 2022. 36 p.

14. Lee G. Titan Supercomputer: Applications and Architecture. Oak Ridge National Laboratory, 2020. 112 p.
15. Fox G. C. Challenges in Managing Distributed Resources. *Journal of Parallel and Distributed Computing*. 2020. Vol. 140. P. 1–10.
16. Fujitsu Ltd. Fugaku: Revolutionizing High-Performance Computing. Fujitsu Whitepapers, 2021. 22 p.
17. Coulouris G., Dollimore J., Kindberg T., Blair G. Distributed Systems: Concepts and Design. 5th ed. Pearson Education, 2020. 1064 p.
18. Ghosh S. Static Load Balancing in Distributed Computing. *Journal of Parallel and Distributed Computing*. 2020. Vol. 138. P. 85–93.
19. Kruger J. Analysis of Round-Robin Algorithms for Load Balancing. *ACM Transactions on Computing Systems*. 2021. Vol. 39, No. 2. P. 47–55.
20. Nguyen T. Random Assignment Strategies in High-Performance Computing. *IEEE Computing Journal*. 2020. Vol. 58, No. 4. P. 35–42.
21. Brownlee N. Simplified Static Methods for Resource Distribution. *Computing Journal*. 2021. Vol. 49, No. 1. P. 77–84.
22. Foster I. Challenges in Static Load Balancing. *Elsevier Distributed Systems Review*. 2020. Vol. 18, No. 2. P. 12–19.
23. Smith J. Dynamic Load Balancing in High-Performance Systems. Springer, 2021. 354 p.
24. Kalita R. Real-Time Monitoring for Dynamic Resource Allocation. Springer Distributed Systems, 2020. 218 p.
25. Alon N. Least Connection Algorithms for Efficient Load Distribution. *ACM Transactions*. 2020. Vol. 38, No. 5. P. 117–123.
26. Ghosh S. Load-Aware Balancing Techniques in Distributed Systems. Wiley High-Performance Computing Series, 2021. 376 p.
27. Fox G. C. Adaptive Balancing Methods in Modern Computing. Pearson Distributed Systems Analysis, 2020. 295 p.
28. Johnson M. Challenges and Costs of Dynamic Load Distribution. *IEEE Distributed Systems Journal*. 2020. Vol. 32, No. 3. P. 22–29.

29. McCarthy J. Comparative Analysis of Load Balancing Techniques. Springer Distributed Computing Systems, 2021. 288 p.
30. Kesselman C., Foster I. Hybrid Load Balancing: Static Meets Dynamic Approaches. Wiley, 2021. 305 p.
31. Coulouris G., Dollimore J., Kindberg T., Blair G. Distributed Systems: Concepts and Design. 5th ed. Pearson Education, 2020. 1064 p.
32. Tanenbaum A. S., Van Steen M. Distributed Systems: Principles and Paradigms. 2nd ed. Pearson, 2021. 672 p.
33. Dynamic Resource Management in Distributed Systems. Springer. URL: <https://link.springer.com> (дата звернення: 17.02.2025).
34. Ghosh S. Geographical Node Clustering for Resource Optimization. *Journal of Parallel and Distributed Computing*. 2021. Vol. 139. P. 120–128.
35. Kruger J. Priority-Based Clustering in High-Performance Systems. *ACM Transactions on Computing Systems*. 2021. Vol. 40, No. 1. P. 99–106.
36. Foster I., Kesselman C. The Grid: Blueprint for a New Computing Infrastructure. Elsevier, 2021. 696 p.
37. Jain A. K., Murty M. N., Flynn P. J. Data Clustering: A Review. *ACM Computing Surveys*. Vol. 31, No. 3. P. 264–323.
38. Lloyd S. Least Squares Quantization in PCM. *IEEE Transactions on Information Theory*. Vol. 28, No. 2. P. 129–137.
39. Hartigan J. A., Wong M. A. Algorithm AS 136: A K-Means Clustering Algorithm. *Journal of the Royal Statistical Society. Series C (Applied Statistics)*. Vol. 28, No. 1. P. 100–108.
40. Ester M., Kriegel H. P., Sander J., Xu X. A Density-Based Algorithm for Discovering Clusters in Large Spatial Databases with Noise. *Proceedings of the 2nd International Conference on Knowledge Discovery and Data Mining*. P. 226–231.
41. Schubert E., Sander J., Ester M., Kriegel H. P., Xu X. DBSCAN Algorithm: Revisited. *ACM Digital Library*. URL: <https://dl.acm.org> (дата звернення: 28.02.2025).
42. Bezdek J. C., Ehrlich R., Full W. FCM: The Fuzzy C-Means Clustering Algorithm. *Computers & Geosciences*. Vol. 10, No. 2–3. P. 191–203.

43. Pal N. R., Pal K. A Review on Image Segmentation Techniques. *Pattern Recognition Letters*. Vol. 14, No. 9. P. 839–852.
44. Brownlee N. Load Balancing and Node Clustering. *Wiley Distributed Systems Journal*. 2021. Vol. 35, No. 2. P. 45–53.
45. Kalita R. Cluster-Based Optimization for Cloud Systems. Springer, 2020. 228 p.
46. Куліш О. М., Грищенко В. В. Ресурси комп'ютерних систем: класифікація та методи оптимізації. Київ : Вид-во КПІ, 2022. 304 с.
47. Захарченко О. В. Процесорні ресурси: особливості використання в мультикомп'ютерних системах. *Журнал обчислювальних технологій*. 2020. № 4. С. 34–42.
48. Малишев К. П. Оцінка продуктивності багатоядерних процесорів: сучасні підходи. *Збірник наукових праць з прикладної інформатики*. 2021. № 5. С. 78–85.
49. Горбач В. І., Палійчук С. В. Методика оцінки завантаженості процесорних ядер у мультипроцесорних системах. *Комп'ютерні науки і кібернетика*. 2019. № 3. С. 59–66.
50. Кулиняк І. Я. Основи аналізу пам'яті у мультикомп'ютерних системах. *Інформаційні технології в освіті і науці*. 2020. № 7. С. 112–118.
51. Воробйова О. А. Нові підходи до використання кешу в багатоядерних процесорах. *Журнал сучасних досліджень*. 2021. № 6. С. 123–130.
52. Мельниченко О. А. Методика аналізу ефективності оперативної пам'яті. *Аналітичний вісник*. 2023. № 2. С. 44–51.
53. Борисенко Л. В. Мережеві ресурси мультикомп'ютерних систем: виклики і перспективи. *Журнал інформаційних технологій*. 2021. № 5. С. 89–95.
54. Семенюк П. Т. Методика вимірювання пропускнуої здатності мереж у розподілених системах. *Збірник статей науково-практичного семінару*. 2022. С. 98–104.

55. Назаров А. О. Вплив мережеских затримок на продуктивність мультимп'ютерних систем. *Праці Міжнародної конференції з інформаційних технологій*. 2023. С. 135–142.
56. Поліщук О. М. Системний аналіз взаємодії процесорів, пам'яті та мережеских компонентів. *Журнал інженерії програмного забезпечення*. 2022. № 4. С. 54–63.
57. Сидоренко К. В. Балансування навантаження в мультимп'ютерних системах: теорія та практика. *Збірник наукових праць з прикладної математики*. 2023. № 3. С. 67–74.
58. Лазаренко І. Г. Хмарні обчислення: аналіз продуктивності системи. *Інформаційні системи і технології*. 2021. № 6. С. 101–109.
59. Гладкий Ю. П. Грід-системи: основи оптимізації передачі даних. *Міжнародний журнал інформаційних технологій*. 2020. № 7. С. 115–121.
60. Сучасні підходи до аналізу обчислювальних систем. НТОУ. 2023. URL: <https://nto.ua/assets/files/ntou-library-war-%D0%86mpact-analysis.pdf> (дата звернення: 17.11.2024).
61. Процесорні, пам'ятні та мережескі ресурси в мультимп'ютерних системах. Gartner, 2021. URL: <https://www.gartner.com/resources/computing-report> (дата звернення: 15.03.2025).
62. Гнатенко А. П. Проблеми нерівномірного розподілу обчислювальних ресурсів у мультимп'ютерних системах. *Журнал системних досліджень*. 2021. № 3. С. 45–52.
63. Назаренко В. О., Малишев К. П. Вплив гетерогенності ресурсів на продуктивність мультимп'ютерних систем. *Інформаційні технології та системний аналіз*. 2020. № 5. С. 33–39.
64. Борисенко Л. В., Ткаченко С. М. Динамічність завдань у мультимп'ютерних системах: проблеми та підходи до вирішення. *Комп'ютерні науки і кібернетика*. 2021. № 7. С. 65–73.

65. Петров І. М., Сидоренко О. В. Проблеми масштабування мультикомп'ютерних систем. *Журнал обчислювальних технологій*. 2022. № 2. С. 28–34.
66. Звіт: Балансування навантаження в розподілених обчислювальних системах. НТОУ. 2023. URL: <https://nto.ua/resources/load-balancing-analysis.pdf> (дата звернення: 26.03.2025).
67. Мельниченко О. А. Час обробки завдань у мультикомп'ютерних системах: фактори впливу. *Журнал сучасних досліджень*. 2022. № 6. С. 78–85.
68. Тараненко К. Г. Неврахування гетерогенності як основна причина зниження ефективності систем. *Аналітичний вісник*. 2021. № 3. С. 112–118.
69. Сидоров А. В. Проблема перевантажених вузлів у мультикомп'ютерних системах. *Праці конференції «Інформаційні технології у системному аналізі»*. 2023. С. 56–63.
70. Захарченко О. В., Горбач В. І. Аналіз статичних алгоритмів розподілу завдань у мультикомп'ютерних системах. *Інформаційні системи*. 2022. № 4. С. 89–96.
71. Поліщук О. М., Назаренко В. О. Недовантажені вузли: причини та наслідки. *Журнал інформаційних технологій*. 2020. № 8. С. 44–51.
72. Палійчук С. В. Недовантаження вузлів у системах: аналіз гнучкості алгоритмів. *Комп'ютерні науки*. 2023. № 5. С. 98–104.
73. Гординюк І. О. Витрати на комунікацію в мультикомп'ютерних системах. *Журнал інженерії програмного забезпечення*. 2021. № 9. С. 101–109.
74. Семенюк П. Т., Мельниченко О. А. Кластеризація вузлів для зниження комунікаційних витрат. *Праці конференції «Інновації в комп'ютерних системах»*. 2023. С. 135–142.
75. Лазаренко І. Г. Масштабування мультикомп'ютерних систем: виклики та перспективи. *Інформаційні технології*. 2021. № 10. С. 112–119.
76. Гладкий Ю. П. Масштабування алгоритмів розподілу завдань. *Міжнародний журнал обчислювальних технологій*. 2022. № 6. С. 115–121.

77. Збереження даних у мультикомп'ютерних системах. НТОУ, 2023.
URL: <https://nto.ua/resources/data-preservation-report.pdf> (дата звернення: 02.04.2025).
78. Назаров А. О., Ткаченко С. М. *Вплив перевантаження вузлів на надійність систем*. Праці конференції «Інформаційні технології та надійність систем». 2023. С. 67–74.
79. Поліщук О. М. Енергоспоживання у мультикомп'ютерних системах. *Журнал сучасних технологій*. 2021. № 7. С. 101–109.
80. Оптимізація енергоспоживання в розподілених системах. Gartner, 2023.
URL: <https://www.gartner.com/resources/energy-optimization.pdf> (дата звернення: 15.04.2025).

ДОДАТОК А (обов'язковий)

ЛІСТИНГ ПРОГРАМНОГО ЗАБЕЗПЕЧЕННЯ

```

import numpy as np
import matplotlib.pyplot as plt

# Етап 1 - Вхідні дані: CPU і RAM у відсотках
cpu = np.array([95, 55, 15, 20, 90]) # Вузли: А, В, С, D, Е
ram = np.array([80, 60, 25, 30, 85])
nodes = ['A', 'B', 'C', 'D', 'E']

# Етап 2 - Розрахунок інтегрального навантаження
alpha, beta = 0.6, 0.4
L_i = alpha * (cpu / 100) + beta * (ram / 100)
print("Інтегральне навантаження:", np.round(L_i, 3))

# Етап 3 - Визначення вузлів-донорів і приймачів
threshold = 0.6
excess = np.maximum(0, L_i - threshold)
available = np.where(L_i < threshold, 1 - L_i, 0)

# Етап 4 - Розподіл надлишку з вузла А на вузли С і D
res_C, res_D = 0.81, 0.74
total_res = res_C + res_D
excess_A = 0.29
delta_C = excess_A * (res_C / total_res)
delta_D = excess_A * (res_D / total_res)

# Етап 5 - Перенесення з вузла Е на вузол С
excess_E = 0.24
delta_E_C = excess_E

# Етап 6 - Оновлення навантаження після міграції
L_i_new = L_i.copy()
L_i_new[0] -= (delta_C + delta_D) # А
L_i_new[2] += (delta_C + delta_E_C) # С
L_i_new[3] += delta_D # D
L_i_new[4] -= delta_E_C # Е

# Етап 7 - Візуалізація результатів
x = np.arange(len(nodes))
bar_width = 0.35

plt.bar(x - bar_width/2, L_i, width=bar_width, label='До розподілу')
plt.bar(x + bar_width/2, L_i_new, width=bar_width, label='Після
розподілу')
plt.xticks(x, nodes)
plt.ylabel('Інтегральне навантаження L')
plt.title('Розподіл навантаження до і після балансування')
plt.legend()

```

```
plt.grid(axis='y', linestyle='--')
plt.tight_layout()
plt.show()

# Етап 8 - Оцінка ефективності
print(f"Середнє L до: {np.mean(L_i):.3f}, дисперсія:
{np.var(L_i):.3f}")
print(f"Середнє L після: {np.mean(L_i_new):.3f}, дисперсія:
{np.var(L_i_new):.3f}")
```

ДОДАТОК Б
(обов'язковий)
НАУКОВА ПРАЦЯ ЗДОБУВАЧА

УДК 004.75

ГОБА КАРІНА

Хмельницький національний університет
<https://orcid.org/0009-0008-8553-1674>
e-mail: karina2002goba@gmail.com

ЖУКОВСЬКИЙ ПАВЛО

Хмельницький національний університет
<https://orcid.org/0009-0007-3461-9919>
e-mail: 777reiste777@gmail.com

КОВАЛЬЧУК ВАСИЛЬ

Хмельницький національний університет
<https://orcid.org/0009-0008-7013-5919>
e-mail: vasuli458@gmail.com

КЛЕЙН ОЛЕКСАНДР

Хмельницький національний університет
<https://orcid.org/0000-0002-1896-943X>
e-mail: olexandrkleyn@gmail.com

МЕТОД РОЗПОДІЛУ НАВАНТАЖЕННЯ В МУЛЬТИКОМП'ЮТЕРНИХ СИСТЕМАХ ЗАГАЛЬНОГО ПРИЗНАЧЕННЯ ЗГІДНО КЛАСТЕРИЗАЦІЇ ВУЗЛІВ ЗА РІВНЕМ ЗАВАНТАЖЕНОСТІ РЕСУРСІВ

У роботі розглянуто методи розподілу обчислювального навантаження в мультикомп'ютерних системах загального призначення з акцентом на кластеризацію вузлів за рівнем їх завантаженості. Дослідження спрямоване на підвищення ефективності використання обчислювальних ресурсів шляхом динамічного балансування навантаження між вузлами відповідно до їх поточного стану. У процесі виконання здійснено аналіз сучасних підходів до балансування навантаження в багатовузлових системах, розроблено модель кластеризації вузлів залежно від їхньої завантаженості, а також алгоритм, що дозволяє ефективно розподіляти навантаження між сформованими кластерами з метою мінімізації затримок і підвищення рівня використання доступних ресурсів. Запропонований підхід підлягає експериментальному тестуванню в умовах моделювання, що дозволяє порівняти його ефективність із базовими стратегіями розподілу навантаження.

Наукова новизна дослідження полягає у впровадженні нового методу балансування, який передбачає кластеризацію вузлів за ступенем їх завантаження, що забезпечує не лише більш рівномірний розподіл задач у системі, а й сприяє зниженню енергоспоживання та скороченню часу очікування. Практична значущість результатів полягає в можливості застосування запропонованого методу для оптимізації обчислювальних процесів у таких середовищах, як хмарні сервіси, дата-центри, суперкомп'ютери та корпоративні обчислювальні мережі, де ефективне використання ресурсів має критичне значення.

Для реалізації дослідження застосовуються методи кластеризації, що враховують поточну завантаженість обчислювальних вузлів. Це дає змогу адаптивно розподіляти навантаження відповідно до динамічного стану системи.

Ключові слова: розподіл навантаження, мультикомп'ютерні системи, кластеризація вузлів, завантаженість ресурсів, алгоритми розподілу навантаження, вузли системи, ресурсний моніторинг, оптимізація розподілу.

**KARINA GOBA, ZHUKOVSKYI PAVLO, VASILY KOVALCHUK, KLEIN
OLEXANDR**

Khmelnytskyi National University, Khmelnytskyi, Ukraine

LOAD DISTRIBUTION METHOD IN GENERAL-PURPOSE MULTI-COMPUTER SYSTEMS ACCORDING TO NODE CLUSTERING BY RESOURCE UTILIZATION LEVEL

The study examines methods for distributing computational load in general-purpose multicomputer systems, with a focus on clustering nodes based on their load levels. The research aims to improve the efficiency of resource utilization through dynamic load balancing across nodes, according to their current operational state. It involves an analysis of contemporary approaches to load balancing in multi-node systems, the development of a node clustering model based on load levels, and the creation of an algorithm that enables efficient load distribution among clusters to minimize delays and maximize resource utilization. The proposed approach is subject to experimental testing under simulation conditions, allowing for a comparative assessment of its effectiveness against baseline load distribution strategies.

The scientific novelty of the study lies in the implementation of a new load balancing method that involves clustering nodes according to their load levels, which not only ensures a more even distribution of tasks across the system but also contributes to reduced energy consumption and lower latency. The practical significance of the findings lies in the potential application of the proposed method for optimizing computational processes in environments such as cloud services, data centers, supercomputers, and corporate computing networks, where efficient resource management is of critical importance.

To carry out the research, clustering methods that take into account the current load of computing nodes are employed, enabling adaptive load distribution based on the dynamic state of the system.

Keywords: load balancing, multicomputer systems, node clustering, resource load, load distribution algorithms, system nodes, resource monitoring, distribution optimization.

ПОСТАНОВКА ПРОБЛЕМИ У ЗАГАЛЬНОМУ ВИГЛЯДІ

У сучасних розумних будинках, які виступають прикладом мультикомп'ютерних систем загального призначення, використовується широкий спектр інтегрованих пристроїв — від сенсорів і систем клімат-контролю до інтелектуальних систем відеоспостереження. В основі такої інфраструктури лежить координація численних обчислювальних вузлів, які або здійснюють локальну обробку даних, або передають їх на периферійні чи хмарні обчислювальні ресурси.

У періоди пікового навантаження, особливо у вечірні години, система розумного будинку зазнає значного зростання обчислювальної активності. Одночасно активуються кілька критично важливих підсистем: системи відеоаналітики працюють у режимі реального часу з використанням алгоритмів розпізнавання облич та виявлення аномалій; голосові асистенти обробляють запити мешканців; системи клімат-контролю регулюють мікроклімат на основі показань сенсорів; локальні модулі обробки даних забезпечують конфіденційність персональної інформації.

Унаслідок цього виникає нерівномірний розподіл обчислювального навантаження між вузлами. Одні вузли працюють на межі своїх ресурсних можливостей, тоді як інші залишаються недостатньо задіяними. Така ситуація призводить до зниження загальної ефективності системи, збільшення затримок в обробці даних, погіршення якості обслуговування, а в окремих випадках — до деградації окремих функціональних компонентів системи.

АНАЛІЗ ГОТОВИХ РІШЕНЬ ТА ДОСЛІДЖЕНЬ

Мультикомп'ютерні системи — це обчислювальні системи, що складаються з декількох автономних вузлів (комп'ютерів), об'єднаних в єдину інфраструктуру для виконання спільних завдань. Їх основною характеристикою є розподілений характер обчислень, що дозволяє виконувати задачі паралельно на різних вузлах, оптимізуючи час і використання ресурсів [1]. Мультикомп'ютерні системи можуть будуватися як із гомогенних (однакових за ресурсами), так і з гетерогенних (різних за можливостями) вузлів. Зазвичай для комунікації між вузлами використовується мережа передачі даних [2]. Масштабованість є однією з головних переваг таких систем, бо є можливість додавання нових вузлів для підвищення продуктивності без значних змін у загальній архітектурі [3]. Мультикомп'ютерні системи загального призначення розроблені для вирішення широкого спектра завдань [4]. Кластерні системи — це обчислювальні системи, у яких вузли з'єднані через локальну мережу (LAN). Вони зазвичай застосовуються для високопродуктивних обчислень і вирішення ресурсомістких задач [5]. Хмарні обчислення у таких системах вузли взаємодіють через Інтернет, надаючи доступ до ресурсів за моделями IaaS (інфраструктура як сервіс), PaaS (платформа як сервіс) і SaaS (програмне забезпечення як сервіс). Прикладами є Amazon Web Services (AWS) та Google Cloud Platform [6]. Grid-системи — це системи, що з'єднують вузли, розташовані у різних географічних місцях, для виконання масштабних обчислень. Їх відрізняє низький рівень інтеграції між вузлами [7]. Суперкомп'ютери — складаються з великої кількості вузлів, оптимізованих для високопродуктивних обчислень. Застосовуються для моделювання клімату, обробки даних з космосу тощо [8].

Прикладами мультикомп'ютерних систем є кластери для обробки даних, такі як Hadoop і Spark, хмарні платформи на зразок Amazon Web Services і Google Cloud Platform, а також суперкомп'ютери Titan і Fugaku, які використовуються в наукових дослідженнях [9,10]. Hadoop — це платформа для розподіленої обробки великих наборів даних. Вона базується на моделі MapReduce і дозволяє розподіляти завдання між численними вузлами системи, забезпечуючи паралельну обробку даних. Основні сфери застосування включають аналіз даних, обробку логів, побудову рекомендаційних систем і аналіз соціальних мереж [11]. Apache Spark — це платформа розподілених обчислень, призначена для швидкої обробки даних. Підтримує обробку даних у пам'яті, що суттєво прискорює

виконання завдань порівняно з Hadoop. Spark використовується в машинному навчанні, потоковій обробці даних, аналізі великих даних і побудові моделей прогнозування [12]. AWS — хмарна платформа, яка надає інфраструктуру, платформи та програмне забезпечення як послуги (IaaS, PaaS, SaaS). Забезпечує ресурси для обчислень, зберігання даних, машинного навчання, аналітики й розробки масштабованих і високо доступних систем. Широко використовується для розробки вебзастосунків, розміщення баз даних і побудови хмарних сервісів [13]. Titan — суперкомп'ютер, який використовується в наукових дослідженнях, таких як моделювання змін клімату, обробка великих обсягів даних у геноміці, моделювання поведінки матеріалів і аналіз фізичних явищ. Його висока продуктивність дозволяє виконувати складні обчислення з використанням паралельних алгоритмів [14,15]. Fugaku — один із найпотужніших суперкомп'ютерів у світі, застосовується у розробці ліків, кліматичному моделюванні, прогнозуванні землетрусів, фізичних дослідженнях і машинному навчанні. Його можливості дозволяють працювати з задачами, що потребують обробки надзвичайно великих обсягів даних [16].

Балансування навантаження це критично важливий аспект функціонування мультикомп'ютерних систем. Продуктивність, масштабованість і стабільність системи значною мірою залежать від ефективного використання обчислювальних ресурсів. Залежно від принципів організації методи балансування поділяються на статичні та динамічні [17]. Статичні методи передбачають заздалегідь визначений розподіл завдань між вузлами, який не змінюється під час виконання. Розподіл ґрунтується на фіксованих правилах, таких як кругове (Round-Robin) або випадкове (Random Assignment) призначення. Ці методи не враховують поточний стан системи, але відзначаються простотою реалізації [18].

Таким чином, актуальним завданням в контексті мультикомп'ютерних систем є розподіл навантаження з метою досягнення ефективного використання ресурсів.

ВИЗНАЧЕННЯ РІВНЯ НАВАНТАЖЕННЯ КОЖНОГО ВУЗЛА

Запропоноване рішення зосереджене на розподілі обчислювального навантаження в мультикомп'ютерних системах загального призначення шляхом кластеризації обчислювальних вузлів відповідно до поточного рівня використання їхніх ресурсів. Для реалізації цього підходу використовується гібридна модель машинного навчання, що поєднує методи контрольованого навчання, зокрема штучні нейронні мережі (ШНМ) та неконтрольовані методи, такі як кластеризація Best of K (БОК), з метою оптимального групування вузлів зі схожими характеристиками навантаження.

Система складається з множини вузлів обчислення $P = \{PM1, PM2, \dots, PMM\}$, кожен з яких містить набір програмно або апаратно ізольованих обчислювальних одиниць (віртуалізованих або фізичних), які в рамках моделі позначимо як $VM = \{VM1, VM2, \dots, VML\}$. Система оперує множиною завдань користувача $T = \{T1, T2, \dots, TS\}$, які надходять для обробки в режимі реального часу або пакетної обробки. Для балансування навантаження між вузлами та їх обчислювальними одиницями використовуються два ключові програмні компоненти. Диспетчер вузлів (Node Dispatcher) динамічно призначає завдання відповідним вузлам на основі їх поточного стану. Балансувальник навантаження (Load Balancer) координує рівновагу між кластерами, забезпечуючи стабільність та ефективність розподілу. Для точного оцінювання рівня навантаження кожного вузла застосовується гібридний підхід, заснований на глибокому навчанні, що включає два типи нейронних мереж. Згорткова нейронна мережа (CNN, Convolutional Neural Network) використовується для виділення просторових ознак із вхідних даних (наприклад, телеметрія, показники завантаженості процесора, пам'яті чи мережі). Такий тип мереж ефективно виявляє закономірності у структурованих сіткоподібних даних, таких як розподіл навантаження в масиві вузлів. Рекурентна нейронна мережа (RNN, Recurrent Neural Network) фіксує часові залежності в даних, що дозволяє системі враховувати динаміку змін навантаження з часом. Це є критично важливим для прогнозування стану системи в умовах динамічного середовища. Згорткові нейронні мережі є одними з найефективніших методів для аналізу структурованих та візуально закодованих даних. Архітектура CNN чергує згорткові шари, шари субдискретизації (pooling) та повноз'язані шари. У запропонованій системі модуль CNN виконує роль попереднього аналізатора даних, виявляючи типові шаблони використання ресурсів на рівні вузлів або окремих обчислювальних одиниць.

Для початкового аналізу та структурного виділення просторових ознак із даних про навантаження вузлів використовується згорткова нейронна мережа (CNN). Розглянемо ключові компоненти її архітектури та їх роль у розв'язанні поставленої задачі. Основою обробки в CNN є згортковий шар, який застосовує набір фільтрів (ядер згортки) до вхідних даних з метою виявлення локальних шаблонів. У нашому контексті дані відповідають часовим рядам показників використання ресурсів (наприклад, процесора, пам'яті, мережі) для кожного вузла. Кожне ядро згортки аналізує невелику ділянку вхідного простору, тобто поле сприйняття (receptive field) для виявлення локальних характеристик навантаження. Математично формалізуємо це як згорткову операцію між вхідним сигналом $i_c(x, y)$ та ядром $e_{kl}(u, v)$, що формує карту ознак $f_{kl}(p, q)$. Повну карту ознак задамо так:

$$F_{kl} = [f_{kl}(1,1), \dots, f_{kl}(p, q), \dots, f_{kl}(P, Q)]. \quad (1)$$

Шари субдискретизації (Pooling layers) зменшують розмірність просторових ознак шляхом агрегування найважливішої інформації в межах кожного поля сприйняття. Це робить мережу стійкою до незначних зсувів або спотворень у вхідних даних. Шари субдискретизації допомагають зменшити перенавчання та підвищити здатність моделі до узагальнення.

Процедура субдискретизації задамо функцією:

$$Z_{kl} = g_p(F_{kl}), \quad (2)$$

де g_p — тип pooling-операції (наприклад, max-pooling або average-pooling).

Для того щоб модель могла навчатися складним закономірностям використовується активаційна функція, яка вводить нелінійність у модель. На практиці це дозволяє вловлювати залежності між різними аспектами навантаження, які не можуть бути відображені виключно лінійними зв'язками. Задамо її так:

$$T_{kl} = g_a(F_{kl}), \quad (3)$$

де g_a — активаційна функція.

На завершальному етапі згортова нейронна мережа (CNN) використовує повнозв'язаний шар, який інтегрує всі раніше виділені ознаки для формування узагальненого представлення стану системи. Це дозволяє моделі враховувати взаємозв'язки між різними елементами інформації (наприклад, навантаженням процесора та активністю мережі), що є критично важливим для точної класифікації або регресійного аналізу. Вихідні дані цього шару зазвичай використовуються як основа для класифікації вузлів за рівнем навантаження (наприклад, для кластеризації), прогнозування майбутнього навантаження в наступні моменти часу (у поєднанні з RNN) передачі результатів модулю балансування навантаження для ухвалення рішень. Для моделювання динаміки змін навантаження в мультикомп'ютерній системі в режимі реального часу було використано рекурентну нейронну мережу (RNN), яка добре зарекомендувала себе у задачах обробки часових рядів. На відміну від класичних багатозарових перцептронів (MLP), RNN має внутрішні зворотні зв'язки між прихованими шарами, що дозволяє мережі зберігати та передавати інформацію про попередні стани. Завдяки цій властивості RNN здатна моделювати довгострокові часові залежності між подіями, що є критично важливим при оцінюванні навантаження вузлів, де історія попередніх станів може суттєво впливати на поточне навантаження.

У реалізованій моделі вхідні послідовності, сформовані з метрик використання ресурсів (процесор, оперативна пам'ять, мережа тощо), передаються через рекурентну структуру у межах ковзного часового вікна. Вихід прихованого стану в момент часу t обчислюється з використанням активаційної функції гіперболічного тангенса за наступною формулою:

$$h_t = \tanh(W_x x_t + W_h h_{t-1} + b), \quad (4)$$

де: h_t — поточний стан прихованого шару; x_t — вхідний вектор на момент часу t ; W_x — вагова матриця між входом і прихованим шаром; W_h — вагова матриця між попереднім і поточним станом прихованого шару; b — вектор зсуву (bias).

Вагові коефіцієнти мережі оптимізуються за допомогою зворотного поширення помилки в часі (Backpropagation Through Time, BPTT), із урахуванням похідних функцій втрат відносно кожного шару. Ключовою особливістю моделі є інтеграція архітектур CNN та RNN, де згортова нейронна мережа (CNN) виконує попереднє виділення просторових ознак із вхідних метрик, а рекурентна нейронна мережа (RNN) аналізує їх часову еволюцію. Після об'єднання вихідних даних обох мереж формується комплексна оцінка навантаження для кожного вузла.

Отримані значення навантаження використовуються для кластеризації вузлів системи відповідно до їх поточного рівня завантаженості, що дозволяє ефективно розподіляти ресурси згідно з реальним станом системи в режимі реального часу.

ВИБІР ТИПУ АРХІТЕКТУРИ ТА ШАБЛОНІВ ПРОЄКТУВАННЯ

Під час розробки методу розподілу навантаження для мультикомп'ютерних систем загального призначення, що базується на кластеризації вузлів за рівнем завантаженості ресурсів, надзвичайно важливо забезпечити гнучкість, масштабованість і надійність запропонованого рішення. Для досягнення цих цілей необхідно визначити найбільш доцільний тип архітектури та застосувати відповідні шаблони проектування, які сприятимуть подальшій реалізації, супроводу та розвитку системи. З урахуванням особливостей проекту, що полягає в наявності багатьох вузлів (комп'ютерів), що спільно виконують обчислювальні завдання, та необхідності безперервного обміну даними про завантаження ресурсів доцільно обрати розподілену архітектуру з елементами кластеризації. У розподіленій архітектурі можна динамічно додавати або вилучати вузли, що дозволяє системі гнучко масштабуватись відповідно до зростання вимог до продуктивності. У разі відмови одного з вузлів інші продовжують роботу, що мінімізує ризик повного припинення функціонування системи. Кожен вузол виконує обчислення та обробляє запити згідно зі своїм поточним станом завантаження, що покращує загальну продуктивність і зменшує час відгуку.

З огляду на потребу обробки та аналізу даних про завантаження ресурсів (ЦП, пам'ять тощо) у реальному або наближеному до реального часі, а також потребу у періодичній кластеризації, доцільно структурувати систему за принципами багатозарової (layered) або мікросервісної архітектури. Спрощена трирівнева (3-шарова) архітектура зображена на рис. 1.

Рівень доступу до даних (Data Access Layer) здійснює збір інформації про стан кожного вузла (поточний рівень завантаження, доступна пам'ять, пропускна здатність мережі та інші показники). Рівень бізнес-логіки (Business Logic Layer) виконує кластеризацію вузлів за рівнем навантаження, визначає оптимальний вузол для призначення нових або існуючих завдань та реалізує алгоритми розподілу навантаження. Рівень подання (Presentation Layer) візуалізує дані (наприклад, моніторинг стану вузлів, результати кластеризації) та забезпечує адміністративний інтерфейс для налаштування політик розподілу й перегляду статистики системи. Для ефективної реалізації обраного підходу використовується набір шаблонів проектування, що дозволяє вирішити типові задачі,

властиві розподілені системам: забезпечити гнучку зміну алгоритмів розподілу навантаження, динамічно додавати нові вузли, організувати збір та обробку даних тощо. Шаблон Strategy (Стратегія) інкапсулює різні алгоритми розподілу навантаження (наприклад, Random, Round Robin, Least Loaded) та дозволяє динамічно обирати стратегію без зміни клієнтського коду. Таке використання дає змогу реалізувати різні методи розподілу навантаження відповідно до потреб системи. Наприклад, дозволяє швидко додати або змінити алгоритм, який обирає найменш завантажений кластер чи вузол, з урахуванням результатів кластеризації. Шаблон Observer (Спостерігач) організовує механізм сповіщення зацікавлених об'єктів про зміну стану іншого об'єкта. Використання сервісу або компонент, що відстежує ресурси вузлів, може автоматично надсилати оновлення всім підписаним модулям (наприклад, модулю кластеризації чи рівню бізнес-логіки), що спрощує керування актуальною інформацією про навантаження та усуває потребу у постійному опитуванні. Шаблон Singleton (Одинак) гарантує існування лише одного екземпляра класу з глобальною точкою доступу. Використання головного модуля диспетчеризації або контролера балансування навантаження підтримує глобальний стан системи (наприклад, статистику вузлів). Це спрощує централізоване управління даними про стан системи. Шаблон Factory Method / Abstract Factory (Фабричний метод / Абстрактна фабрика) інкапсулює створення об'єктів, дозволяючи підкласам визначити, який клас необхідно створити.

Спрощена діаграма послідовності (Sequence Diagram), що ілюструє взаємодію основних компонентів (моніторинг вузлів, кластеризація, розподіл навантаження), наведена на рис. 2. Кожен вузол надсилає свої поточні метрики компоненту Monitor. Компонент Monitor приймає та обробляє цю інформацію (наприклад, застосовуючи шаблон Observer для сповіщення підписаних підсистем) і передає її модулю LoadBalancer. Модуль LoadBalancer виконує кластеризацію, обирає відповідну стратегію (через шаблон Strategy) та приймає рішення, куди направити наступне завдання. Далі LoadBalancer надсилає команду на розподіл конкретного завдання (distributeTask) до обраного вузла або кластера. Вузол приймає завдання та виконує його (runTask), оновлює свої метрики, які знову повертаються до компонента Monitor.

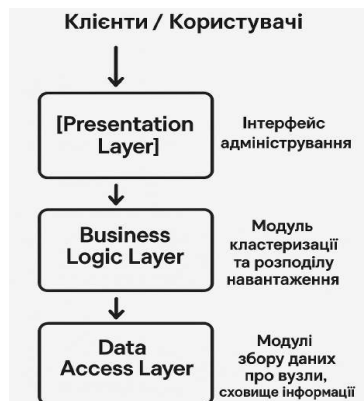


Рисунок 1 – Трирівнева (3-layer) архітектура

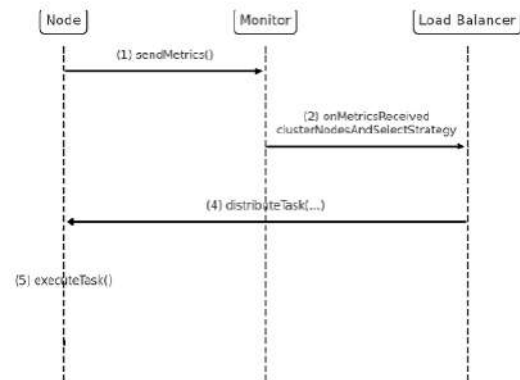


Рисунок 2 – Спрощена діаграма послідовності

Застосування методу балансування навантаження в мультикомп'ютерних системах, заснованого на кластеризації вузлів за рівнем використання ресурсів, може бути ефективно інтегроване в сферу розумного дому. Така інтеграція забезпечує підвищену енергоефективність, покращену надійність і стабільну роботу компонентів системи розумного дому. Реалізація запропонованого методу можлива шляхом створення єдиної мультикомп'ютерної мережі, що складається з вузлів (наприклад, мікрокомп'ютерів Raspberry Pi), кожен з яких відповідає за окрему групу задач розумного дому: керування освітленням, системи безпеки та відеоспостереження, клімат-контроль (опалення, кондиціонування), сенсори руху, вологості та температури, керування побутовою технікою та мультимедійними пристроями. Інтелектуальний розподіл навантаження дозволяє ефективно управляти піковими ситуаціями, коли велика кількість пристроїв працює одночасно (наприклад, у вечірній або ранковий час). Метод кластеризації забезпечує адаптивний перерозподіл задач, запобігаючи зниженню продуктивності чи відмовам обладнання.

Запропонована архітектура є кластеризованою мультикомп'ютерною системою, яка реалізує розподілене управління ресурсами компонентів розумного дому. Центральним елементом цієї архітектури виступає інтелектуальна система управління, що реалізує розроблений автором алгоритм розподілу навантаження на основі кластеризації вузлів за поточним рівнем використання ресурсів. Система складається з кількох автономних вузлів, кожен з яких формує окремий кластер і відповідає за управління певною групою побутових пристроїв. Пристрої поділяються функціонально на такі групи: освітлення; відеоспостереження; кліматичні системи (опалення, вентиляція, кондиціонування); системи безпеки та контролю доступу; побутова техніка; сенсорні пристрої (датчики руху, температури, вологості тощо) тощо. Інтелектуальна система управління постійно здійснює моніторинг поточного стану вузлів, оцінює рівень їхнього навантаження та визначає, чи потрібен перерозподіл задач між вузлами. У випадку перевищення допустимого навантаження на конкретному вузлі система автоматично виконує динамічну міграцію задач до менш завантажених вузлів. Це забезпечує стабільну роботу системи, запобігає перевантаженням та відмовам, підвищує енергоефективність і гарантує надійну роботу компонентів розумного дому. Мережева архітектура розробленої системи розумного дому включає сенсори, побутові прилади, клієнтські

контролери, а також віддалені центри зберігання та обробки даних. Сенсори та прилади з'єднані з контролерами, розташованими безпосередньо в приміщеннях житла. Зв'язок між пристроями та контролерами забезпечується за допомогою додаткових модулів, що підтримують дистанційне управління і комунікацію через бездротові протоколи передачі даних зображено на рис. 3.

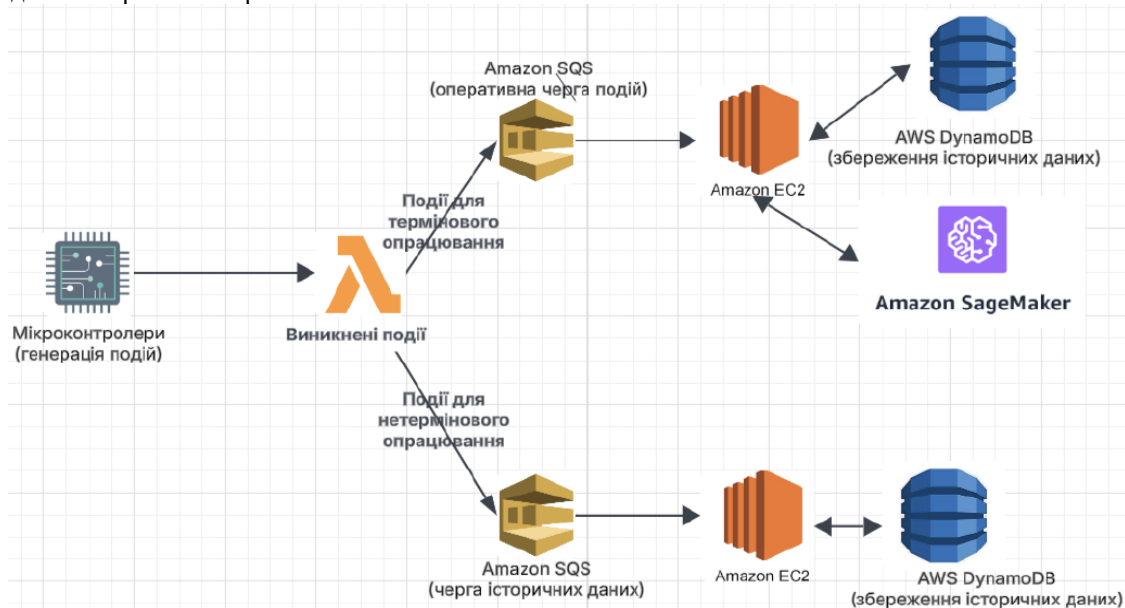


Рисунок 3 - Схема розподілу навантаження подій між чергами

Керування побутовими пристроями та зчитування даних із сенсорів здійснюється через інфрачервоні сигнали та популярні бездротові стандарти зв'язку, такі як Wi-Fi, Bluetooth і ZigBee. Ці протоколи широко використовуються в сучасних сенсорах і підтримуються більшістю вбудованих і додаткових модулів для мікроконтролерів і мікрокомп'ютерів. Використовуємо сервіс Amazon API Gateway для створення RESTful API на базі HTTP або REST-протоколів. Клієнтська сторона взаємодіє із сервером, надсилаючи HTTP-запити (методи GET і POST) до попередньо визначених URL-адрес для отримання інформації про новостворені моделі штучних нейронних мереж або для передачі подій, що відбуваються у приміщеннях. Amazon API Gateway ефективно розподіляє запити між серверами, обираючи найменш завантажений або географічно найближчий до клієнта вузол. Це забезпечує рівномірний розподіл навантаження, мінімізує затримки при передачі даних і прискорює час відгуку серверів. У результаті клієнтська сторона швидко отримує результати обробки нових даних із сенсорів. Крім того, API Gateway підтримує одночасне використання кількох версій API, що дозволяє клієнтам продовжувати роботу зі старими версіями навіть після випуску нових. Також сервіс забезпечує гнучке управління стадіями розгортання API (alpha, beta, production), кожна з яких може бути окремо налаштована для взаємодії з різними серверними кінцевими точками відповідно до конфігурації API. Об'єкти даних передаються з клієнтської сторони на серверну через Інтернет. Важливою особливістю систем розумного дому є велика кількість подій, що надсилаються на сервер від численних клієнтів і потребують обробки в режимі реального часу. Для організації черги обробки подій та оптимізації використання обчислювальних ресурсів застосовується Amazon Simple Queue Service (SQS). Цей сервіс реалізує архітектурний шаблон брокера повідомлень у розподілених системах, дозволяючи зберігати події в черзі під час пікового навантаження без потреби в негайному залученні додаткових обчислювальних ресурсів. Кожна збережена подія обробляється, щойно обчислювальні ресурси стають доступними, у порядку її надходження до черги. Сервіс Amazon SQS гарантує доставку повідомлень підписаним обробникам згідно з політикою "exactly-once" — кожне повідомлення доставляється лише один раз і залишається доступним до моменту обробки та видалення отримувачем. Оскільки черги працюють за принципом першим прийшов і першим обслуговується (FIFO), то забезпечується точна послідовність надсилання та отримання повідомлень. Система розумного дому використовує дві окремі черги повідомлень: операційну чергу для подій, що використовуються при створенні нових моделей штучних нейронних мереж, чергу історичних даних для подій, які зберігаються в документоорієнтованій базі даних з метою подальшого аналізу та виявлення залежностей за допомогою штучного інтелекту. Розподіл подій між чергами здійснюється за допомогою окремої функції AWS Lambda, яка активується при надходженні нової події, не використовуючи обчислювальні ресурси у періоди простою, коли нові події від клієнтів не надходять.

ЕКСПЕРИМЕНТАЛЬНЕ ДОСЛІДЖЕННЯ ЕФЕКТИВНОСТІ МЕТОДУ

У сучасних розумних будинках, що функціонують як мультикомп'ютерні системи, використовується широкий спектр інтегрованих пристроїв: від сенсорних модулів і систем клімат-контролю до систем відеоспостереження, підсилених елементами штучного інтелекту. Така інфраструктура включає декілька обчислювальних вузлів, які здійснюють локальну обробку даних або передають їх на периферійні чи хмарні ресурси. У періоди пікової активності, наприклад, у вечірній час, коли всі мешканці перебувають вдома, система розумного будинку зазнає значного зростання обчислювального навантаження. У цей час одночасно активуються декілька

підсистем. Системи відеоспостереження виконують обробку зображень у режимі реального часу із застосуванням алгоритмів розпізнавання облич та виявлення аномалій. Голосові інтерфейси та цифрові асистенти відповідають на запити користувачів. Автоматичні системи клімат-контролю зчитують показники температури, вологості та рівня CO₂, здійснюючи необхідні коригування. Обчислювальні вузли запускають моделі машинного навчання для адаптації поведінки системи до потреб мешканців (наприклад, передбачення дій, графіків та уподобань). Працюють модулі локального зберігання та обробки даних, зокрема для забезпечення вимог щодо конфіденційності. Унаслідок цього виникає нерівномірний розподіл обчислювального навантаження між вузлами системи: деякі вузли (наприклад, ті, що відповідають за відеоаналітику або машинне навчання) працюють на межі своїх ресурсних можливостей, тоді як інші (обробляють менш критичні або періодичні дані) залишаються частково або повністю недоавантаженими. Така ситуація не лише знижує загальну ефективність системи, але й може спричинити затримки в обробці інформації, погіршення якості обслуговування та, в окремих випадках, призвести до часткової втрати функціональності (наприклад, втрата відеокадрів або затримка реагування на критичні сигнали від сенсорів). Кожен вузол обладнаний локальним процесором, оперативною пам'яттю, програмним середовищем і здатністю обмінюватися даними з іншими вузлами через локальну мережу (наприклад, Wi-Fi або Ethernet). На початковому етапі дослідження необхідно зафіксувати поточний стан обчислювальних ресурсів кожного вузла в межах мультикомп'ютерної системи розумного дому. Кожен обчислювальний вузол характеризується двома основними показниками навантаження:

- відсоткове завантаження процесора (CPU, %);
- відсоткове завантаження оперативної пам'яті (RAM, %).

Ці показники збираються шляхом моніторингу системних ресурсів кожного вузла в режимі реального часу. Для моделювання сценарію максимального навантаження використовуються зразки даних, що відображають пікове вечірнє навантаження на систему розумного дому.

$$L_i = \alpha \cdot \frac{CPU_i}{100} + \beta \cdot \frac{RAM_i}{100}, \quad (4)$$

де: L_i — інтегральний коефіцієнт завантаження i -го вузла; CPU_i — рівень використання центрального процесора вузлом i у відсотках; RAM_i — рівень використання оперативної пам'яті вузлом i у відсотках; α — вага процесорного навантаження (наприклад, $\alpha = 0.6$); β — вага навантаження пам'яті (наприклад, $\beta = 0.4$).

Інтегральний графік на рис. 4 показує навантаження вузлів L_i системи до балансування. Він враховує як завантаження процесора (CPU), так і оперативної пам'яті (RAM) з відповідними вагами: $\alpha = 0.6$ для CPU і $\beta = 0.4$ для RAM. Чим вище стовпчик — тим більше загальне навантаження на вузол. Добре видно, що вузли А та Е мають критично високі значення L_i , тоді як В, С та D завантажені мінімально. Блакитний стовпчик — завантаження CPU, помаранчевий стовпчик — завантаження оперативної пам'яті (RAM).

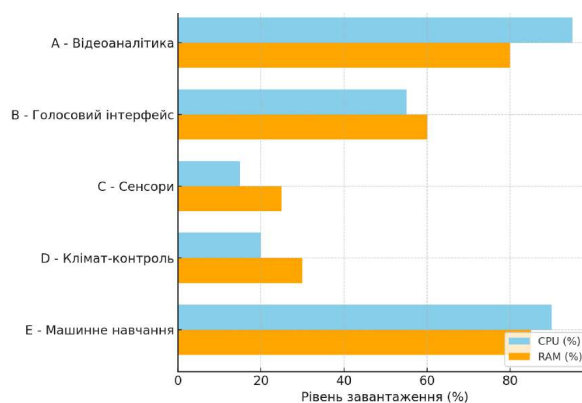


Рисунок 4 - Горизонтальна порівняльна діаграма навантаження вузлів

Оскільки процесорне навантаження та використання пам'яті можуть мати різну вагу впливу на продуктивність системи, вводиться інтегральний показник завантаження вузла — коефіцієнт завантаження L_i .

Розглянемо приклад розрахунку інтегрального навантаження для вузла А:

$$L_A = 0.6 \cdot 0.95 + 0.4 \cdot 0.80 = 0.57 + 0.32 = 0.89.$$

Таким чином, вузол А має коефіцієнт завантаження $L_A = 0.89$, що свідчить про його критичне навантаження. Аналогічні обчислення виконуються для всіх інших вузлів системи. З аналізу обчислених значень видно, що вузли А та Е мають найбільші коефіцієнти навантаження ($L_i > 0.85$), що підтверджує їх критичний стан. Вузели В, С та D мають мінімальні значення інтегрального навантаження, що свідчить про можливість залучення їх ресурсів для перерозподілу обчислювальних задач. Для балансування навантаження у системі було запропоновано виконати міграцію задач наступним чином: частину задач відеоаналітики, які виконуються на вузлі А, перенести на вузли С та D, частину обчислень моделі машинного навчання з вузла Е перенести на вузол С. Таким чином, менш завантажені вузли отримають додаткові задачі, що дасть змогу знизити навантаження на перевантажені вузли і вирівняти загальне навантаження в системі. У нашому випадку розподіл завдань між вузлами не просто випадковий. Було

застосовано балансувальний принцип на основі пропорційного перерозподілу задач, враховуючи поточний рівень завантаження вузлів, наявний вільний ресурс на кожному вузлі, відсоток задач, які можна перенести без втрати продуктивності.

Основною задачею є перерозподілити надлишок навантаження з вузлів з найвищим L_i на вузли з найнижчим L_i , пропорційно до їхньої доступної потужності. Для задачі, яку потрібно мігрувати, обсяг міграції на вузол j визначається за формулою:

$$\Delta W_{i \rightarrow j} = W_{\text{excess}} \times \frac{R_j}{\sum_{k \in S} R_k}, \quad (6)$$

де: $W_{\text{надлишок}}$ — кількість навантаження, яку потрібно зняти з вузла i , R_j — ресурсна ємність (запас потужності) вузла j , S — множина вузлів-кандидатів на прийом задач, $\sum_{k \in S} R_k$ — сумарний ресурс вузлів-кандидатів.

Чим більше у вузла вільних ресурсів, тим більшу частку задач йому можна передати. Наприклад, вузол А мав $L_A = 0.89$, а бажаний рівень — близько 0.6. Тобто, потрібно було знизити навантаження так:

$$W_{\text{excess}(A)} = (0.89 - 0.6) = 0.29. \quad (7)$$

Аналогічно для вузла Е:

$$W_{\text{excess}(E)} = (0.90 - 0.66) = 0.24. \quad (8)$$

Перед балансуванням вузли з найменшим L_i були С і D. Вузол С: доступний запас ресурсів $\approx 1 - 0.19 = 0.81$ Вузол D: доступний запас ресурсів $\approx 1 - 0.26 = 0.74$. Скільки задач отримує кожен вузол пропорційно їх ресурсним запасам визначаємо так:

$$\Delta W_{A \rightarrow C} = 0.29 \times \frac{0.81}{0.81 + 0.74} = 0.1515, \quad (9)$$

$$\Delta W_{A \rightarrow D} = 0.29 \times \frac{0.74}{0.81 + 0.74} = 0.1385. \quad (10)$$

Вузол С отримає приблизно 15.15% одиниць навантаження. Вузол D отримає приблизно 13.85% одиниць навантаження. Разом вони повністю покриють надлишок у 0.29 одиниці, який був на вузлі А. Для вузла Е інтегральне навантаження вузла становило $L_E = 0.9$. Таким чином, надлишок навантаження, який потрібно зняти з вузла Е, дорівнює 0.24. Отже, потрібно зменшити навантаження вузла Е на 0.24 одиниці. Нове навантаження вузла С 0.43. Візуальне порівняння інтегрального навантаження вузлів до та після перерозподілу: жовті стовпці — початковий стан системи, де вузли А і Е були перевантажені; помаранчеві стовпці — стан після балансування, де навантаження розподілено більш рівномірно.

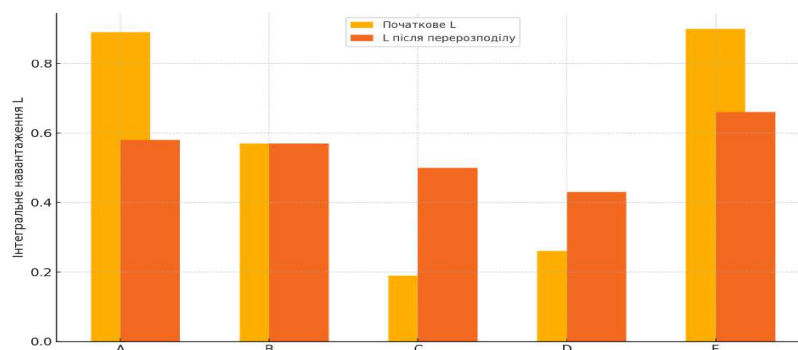


Рисунок 5 - Порівняння навантаження на вузли до та після балансування

ВИСНОВКИ

У результаті впровадження запропонованого алгоритму міграції обчислювальних задач між вузлами мультикомп'ютерної системи розумного дому було досягнуто таких ключових ефектів:

Зменшення навантаження на критично перевантажені вузли: після часткової міграції задач відеообробки з вузла А та обчислень машинного навчання з вузла Е інтегральні коефіцієнти навантаження L_a та L_e істотно знизилися. Це дозволило уникнути перевищення критичних порогів використання ресурсів і зменшити ризик погіршення продуктивності.

Більш ефективне використання недозавантажених вузлів: вузли С (сенсорний модуль) та D (клімат-контроль), які раніше працювали з мінімальним навантаженням, після міграції взяли на себе додаткові задачі. Це забезпечило більш рівномірне використання наявних обчислювальних ресурсів системи. Досягнення більш рівномірного розподілу навантаження в системі досягнуто завдяки оптимізації перерозподілу задач між вузлами, що значно зменшило дисперсію інтегральних коефіцієнтів навантаження L_i . Це сприяє стабільнішій і передбачуванішій роботі системи розумного дому навіть у періоди пікової активності. Балансування обчислювального навантаження суттєво знижує ймовірність відмов або зниження продуктивності, викликаних локальним перевантаженням вузлів.

Система отримує здатність ефективніше досягати самовідновлення та продовжувати функціонування навіть в умовах зростання кількості запитів або неочікуваних стрибків обсягу обробки даних.

Таким чином, впровадження механізму динамічного перерозподілу навантаження на основі кластеризації вузлів дозволило підвищити ефективність використання ресурсів, мінімізувати ризики критичних збоїв і забезпечило стабільну роботу системи розумного дому в умовах змінного навантаження.

Посилання

1. Tanenbaum A.S., Van Steen M. *Distributed Systems: Principles and Paradigms*. Pearson, 2021.
2. Smith J., "High-Performance Distributed Systems". Springer, 2020.
3. Foster I., Kesselman C. *The Grid: Blueprint for a New Computing Infrastructure*. Elsevier, 2021.
4. Ghosh S., "Applications of Distributed Computing Systems". Wiley, 2019.
5. Підручник під грифом МОН. Сімоненко В.П., Дехтярук М.Т., Забара С.С. Програмне забезпечення комп'ютерних мереж. – К.: Університет «Україна», 2014.-220 с.
6. Bedratyuk L., Savenko O. *MATCH Commun. Math. Comput. Chem.* 2018, 79, 407–414.
7. Alon N., "Distributed Computing and Scalability Challenges". *ACM Transactions*, 2020.
8. Dean J., Ghemawat S. "MapReduce: Simplified Data Processing on Large Clusters". *Communications of the ACM*, 2020.
9. Jassy A., "Building Scalable Applications with AWS". AWS Whitepapers, 2022.
10. Lee G. *Titan Supercomputer: Applications and Architecture*. Oak Ridge National Laboratory, 2020.
11. Fox G.C., "Challenges in Managing Distributed Resources". *Journal of Parallel and Distributed Computing*, 2020.
12. Fujitsu Ltd., "Fugaku: Revolutionizing High-Performance Computing". Fujitsu Whitepapers, 2021.
13. Coulouris G., Dollimore J., Kindberg T., Blair G. *Distributed Systems: Concepts and Design*. Pearson Education, 2020.
14. Ghosh S., "Static Load Balancing in Distributed Computing". *Journal of Parallel and Distributed Computing*, 2020.
15. Kruger J., "Analysis of Round-Robin Algorithms for Load Balancing". *ACM Transactions on Computing Systems*, 2021.
16. Nguyen T., "Random Assignment Strategies in High-Performance Computing". *IEEE Computing Journal*, 2020.
17. Brownlee N., "Simplified Static Methods for Resource Distribution". *Computing Journal*, 2021.
18. Foster I., "Challenges in Static Load Balancing". *Elsevier Distributed Systems Review*, 2020.

References

1. Tanenbaum A.S., Van Steen M. *Distributed Systems: Principles and Paradigms*. Pearson, 2021.
2. Smith J., "High-Performance Distributed Systems". Springer, 2020.
3. Foster I., Kesselman C. *The Grid: Blueprint for a New Computing Infrastructure*. Elsevier, 2021.
4. Ghosh S., "Applications of Distributed Computing Systems". Wiley, 2019.
5. Підручник під грифом МОН. Сімоненко В.П., Дехтярук М.Т., Забара С.С. Програмне забезпечення комп'ютерних мереж. – К.: Університет «Україна», 2014.-220 с.
6. Bedratyuk L., Savenko O. *MATCH Commun. Math. Comput. Chem.* 2018, 79, 407–414.
7. Alon N., "Distributed Computing and Scalability Challenges". *ACM Transactions*, 2020.
8. Dean J., Ghemawat S. "MapReduce: Simplified Data Processing on Large Clusters". *Communications of the ACM*, 2020.
9. Jassy A., "Building Scalable Applications with AWS". AWS Whitepapers, 2022.
10. Lee G. *Titan Supercomputer: Applications and Architecture*. Oak Ridge National Laboratory, 2020.
11. Fox G.C., "Challenges in Managing Distributed Resources". *Journal of Parallel and Distributed Computing*, 2020.
12. Fujitsu Ltd., "Fugaku: Revolutionizing High-Performance Computing". Fujitsu Whitepapers, 2021.
13. Coulouris G., Dollimore J., Kindberg T., Blair G. *Distributed Systems: Concepts and Design*. Pearson Education, 2020.
14. Ghosh S., "Static Load Balancing in Distributed Computing". *Journal of Parallel and Distributed Computing*, 2020.
15. Kruger J., "Analysis of Round-Robin Algorithms for Load Balancing". *ACM Transactions on Computing Systems*, 2021.
16. Nguyen T., "Random Assignment Strategies in High-Performance Computing". *IEEE Computing Journal*, 2020.
17. Brownlee N., "Simplified Static Methods for Resource Distribution". *Computing Journal*, 2021.
18. Foster I., "Challenges in Static Load Balancing". *Elsevier Distributed Systems Review*, 2020.

ДОДАТОК В
(обов'язковий)
ПРЕЗЕНТАЦІЯ

Метод розподілу навантаження в
мультимедійних системах загального
призначення згідно кластеризації вузлів за рівнем
завантаженості ресурсів

Хмельницький національний університет

Кваліфікаційна робота магістра

Факультет ІТ - Кафедра КІС

Каріна Гоба

Керівник: проф. Олег Савенко

2025

АКТУАЛЬНІСТЬ ТЕМИ

У сучасних умовах стрімкого зростання обсягів даних та вимог до обчислювальних потужностей, ефективне управління ресурсами мультимедійних систем набуває критичного значення. Особливо це стосується систем загального призначення, де навантаження є динамічним і нерівномірно розподіленим між вузлами. Традиційні методи балансування не завжди забезпечують належну адаптивність і продуктивність.

Запропонований підхід, що базується на кластеризації вузлів за рівнем завантаженості ресурсів, дозволяє забезпечити більш точне й ефективне розподілення задач, покращити використання системних потужностей, зменшити затримки та підвищити стабільність роботи системи.

Актуальність теми зумовлена потребою в розробці інтелектуальних, масштабованих рішень для балансування навантаження в умовах динамічного середовища, що особливо важливо для хмарних сервісів, грід-платформ, кластерних і розумних обчислювальних систем.

МЕТА

- **Об'єктом дослідження** є мультикомп'ютерні системи загального призначення, зокрема процеси розподілу навантаження між вузлами в таких системах.
- **Предметом дослідження** є методи та алгоритми розподілу навантаження в мультикомп'ютерних системах, засновані на кластеризації вузлів за рівнем завантаженості ресурсів.
- **Метою кваліфікаційної роботи** є забезпечення ефективності розподілу навантаження в мультикомп'ютерних системах загального призначення з використанням кластеризації вузлів за рівнем їхнього завантаження.

ПОСТАНОВКА ЗАВДАННЯ

1. Аналіз існуючих методів балансування навантаження в мультикомп'ютерних системах.
2. Розробка методу кластеризації вузлів на основі їхнього поточного завантаження.
3. Оцінка ефективності розподілу ресурсів при використанні запропонованого методу.
4. Моделювання та експериментальне тестування алгоритму на реальних або емуляційних даних.

НАУКОВА НОВИЗНА

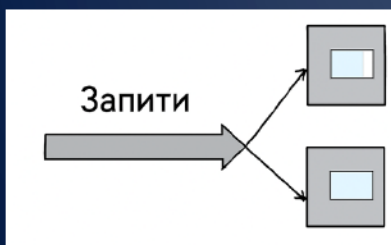
Наукова новизна дослідження полягає в розробці нового методу балансування навантаження, який на відміну від відомих використовує класифікацію вузлів на рівнем їхнього завантаження. Що дозволяє знизити енергоспоживання, підвищити рівномірність розподілу задач та зменшити затримки в мультикомп'ютерних системах.

Практична значущість - запропонований метод може бути застосований для оптимізації обчислювальних процесів у хмарних сервісах, дата-центрах, суперкомп'ютерах та корпоративних мережах, що дозволить значно покращити ефективність використання апаратних ресурсів.

ГОТОВІ МЕТОДИ РОЗПОДІЛУ

Round Robin - статичні

- Розподіляють задачі послідовно по вузлах за колом, незалежно від їх поточного стану.
- Простий у реалізації, але не враховує реального навантаження вузлів.
- Підходить для рівномірних задач.



Least Connection - динамічні

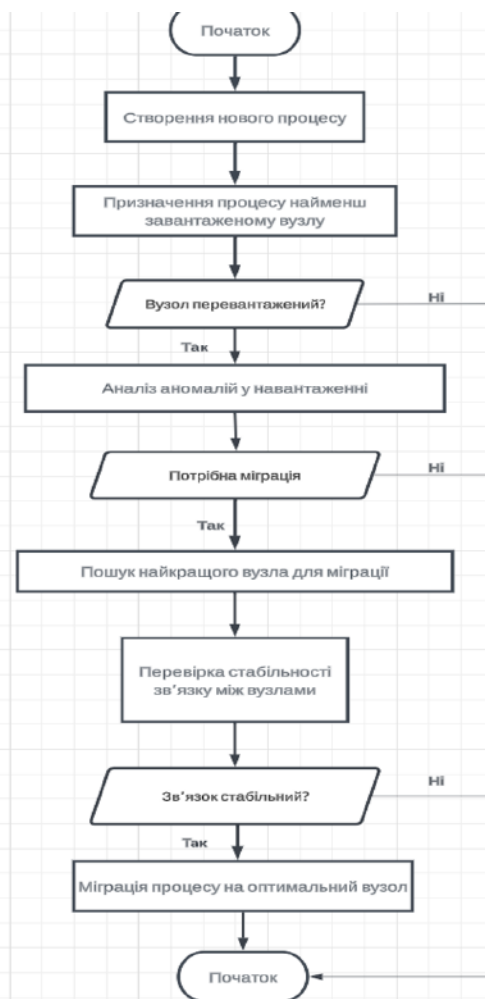
- Розподіляють задачі на основі поточного навантаження — вибирають вузол найменшою кількістю активних з'єднань або задач.
- Забезпечують адаптивність і балансування з урахуванням реальної завантаженості.
- Підходять для нерівномірних та змінних навантажень.

Комбіновані методи

- Поєднують елементи статичних та динамічних методів для оптимізації розподілу навантаження.
- Можуть враховувати як черговість, так і поточний стан вузлів.
- Забезпечують гнучкість і покращену продуктивність.

МОДЕЛЬ МЕТОДУ

1. Запуск нового завдання у системі.
2. Вибір вузла з мінімальним поточним навантаженням.
3. Вузол перевантажений? Так - аналіз аномалій у навантаженні. Ні - повернення до стабільного стану.
4. Перевірка рівномірності завантаженості вузлів.
5. Потрібна міграція? Так - пошук найкращого вузла для міграції. Ні - повернення до стабільного стану.
6. Вибір найменш завантаженого вузла.
7. Перевірка доступності мережевого каналу.
8. Зв'язок стабільний? Так - Міграція процесу. Ні - повернення до початкового стану.
9. Балансування навантаження шляхом перенесення процесу.
10. Фініш – завершення алгоритму.



1. Збираються вихідні параметри кожного вузла, а саме: ступінь завантаження центрального процесора (CPU) та оперативної пам'яті (RAM).

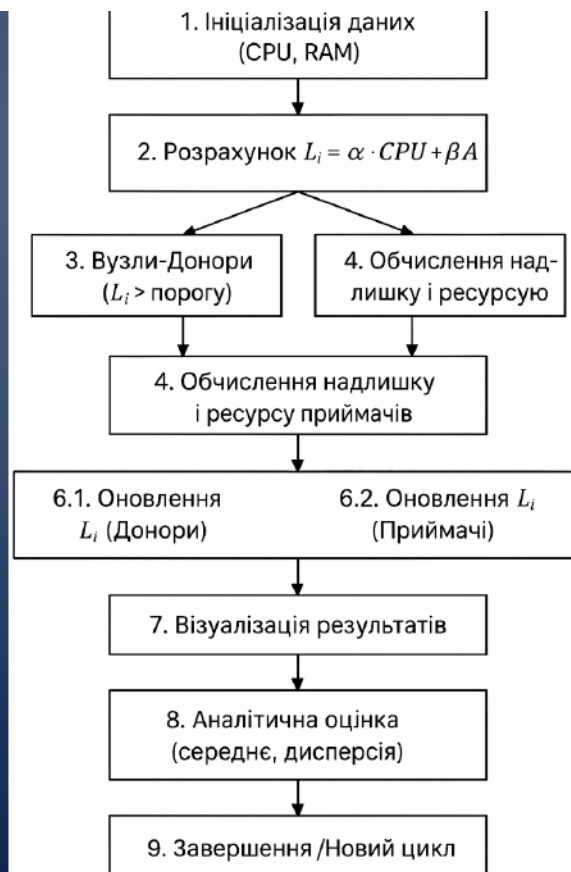
2. Інтегральне навантаження L_i обчислюється як зважена сума відносного навантаження на CPU та RAM

3. Після розрахунку L_i система здійснює класифікацію вузлів. Вузли, навантаження яких перевищує заданий пороговий рівень, потрапляють до групи донорів, вузли з низьким значенням L_i формують групу приймачів.

4. Обчислення обсягів надлишку на кожному з донорів, а також доступного ресурсу у кожному з вузлів-приймачів

5. Розподіл задач: кожен приймач отримує ту частину задач, яка відповідає його частці доступного ресурсу

6. Обчислювальні навантаження вузлів коригуються



Показники збираються шляхом моніторингу системних ресурсів кожного вузла у реальному часі. Для моделювання сценарію максимального навантаження використовуються зразкові дані, що відображають пікові вечірні навантаження на розумний будинок.

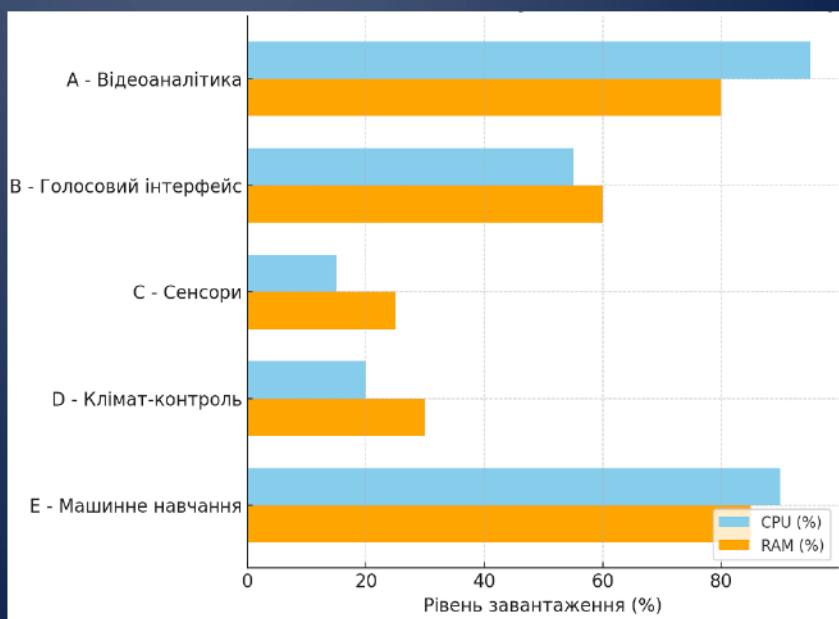
Вузол	Назва модуля	Основне навантаження
A	Відеоаналітика	Аналіз відеопотоків у реальному часі
B	Голосовий інтерфейс	Розпізнавання мови, TTS, NLP
C	Модуль сенсорного моніторингу	Збір і передача показників сенсорів
D	Клімат-контроль	Обробка даних про температуру, CO ₂
E	Обчислювальний модуль (ML)	Виконання моделей штучного інтелекту

ВХІДНІ ПАРАМЕТРИ

На рисунку видно показує навантаження вузлів L_i системи до балансування. Він враховує як завантаження процесора (CPU), так і оперативної пам'яті (RAM)

CPU (%)
95
55
15
20
90

RAM (%)
80
60
25
30
85



Розглянемо приклад розрахунку інтегрального навантаження для вузла А:

$$L_A = 0.6 \cdot 0.95 + 0.4 \cdot 0.80 = 0.57 + 0.32 = 0.89$$

Таким чином, вузол А має коефіцієнт завантаження $L_A = 0.89$, що свідчить про його критичне навантаження. Аналогічні обчислення виконуються для всіх інших вузлів системи.

Вузол	Назва модуля	L_i (ступінь завантаження)
А	Відеоаналітика	0.89
В	Голосовий інтерфейс	0.57
С	Модуль сенсорного моніторингу	0.19
Д	Клімат-контроль	0.26
Е	Обчислювальний модуль (ML)	0.90

Для балансування навантаження у системі було запропоновано виконати міграцію задач наступним чином:

- Частина задач відеоаналітики, які виконуються на вузлі А, перенести на вузли С та Д.
- Частина обчислень моделі машинного навчання з вузла Е перенести на вузол С.

Для задачі, яку потрібно мігрувати, обсяг міграції на вузол j визначається за формулою:

$$\Delta W_{i \rightarrow j} = W_{\text{надлишок}} \times \frac{R_j}{\sum_{k \in S} R_k},$$

Ідея - чим більше у вузла вільних ресурсів, тим більшу частку задач йому можна передати.

Наприклад, вузол А мав $L_A = 0.89$, а бажаний рівень — близько 0.6. Тобто потрібно було знизити навантаження на:

$$W_{\text{надлишок}(A)} = (0.89 - 0.6) = 0.29$$

Скільки задач отримує кожен вузол – пропорційно їх ресурсним запасам:

$$\Delta W_{A \rightarrow C} = 0.29 \times \frac{0.81}{0.81+0.74} = 0.1515$$

$$\Delta W_{A \rightarrow D} = 0.29 \times \frac{0.74}{0.81+0.74} = 0.1385$$

Для вузла Е інтегральне навантаження вузла становило $L_E = 0.9$. Таким чином, надлишок навантаження, який потрібно зняти з вузла Е, дорівнює:

$$W_{\text{надлишок}(E)} = L_{\text{початкове}(E)} - L_{\text{цільове}(E)} = 0.90 - 0.66 = 0.24$$

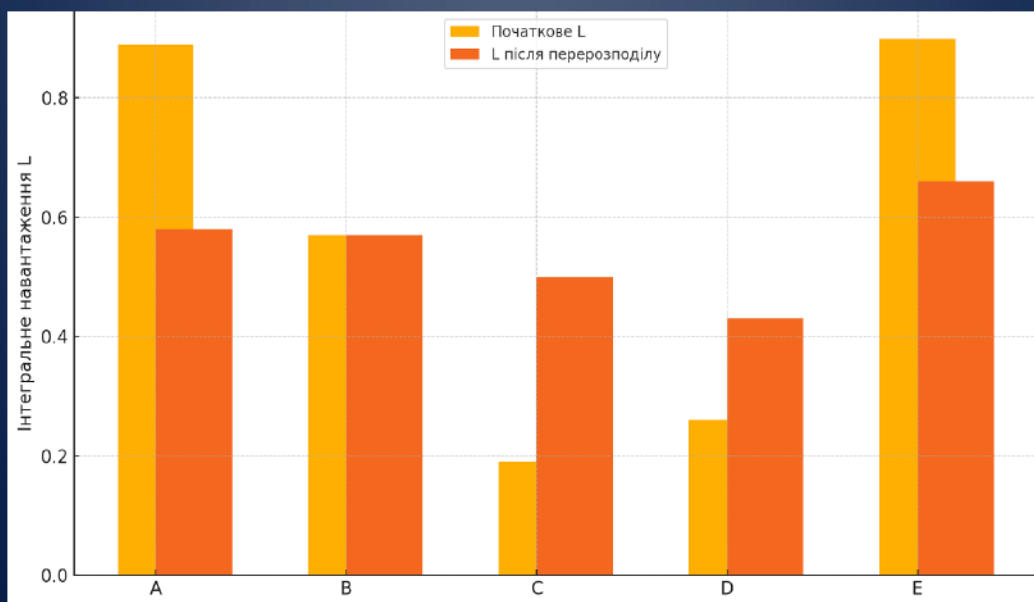
Отже, потрібно зменшити навантаження вузла Е на 0.24 одиниці. Нове навантаження вузла С:

$$L_{\text{нове}(C)} = L_{\text{початкове}(C)} + W_{\text{надлишок}(E)} = 0.19 + 0.24 = 0.43$$

Навантаження всіх вузлів після розподілу

Вузол	L_i (поч.)	Надлишок	L_i (після)	CPU завантаження (%)	RAM завантаження (%)
A	0.89	-0.29	0.58	60	50
B	0.57	0.0	0.57	55	60
C	0.19	0.39149	0.5	50	50
D	0.26	0.1385	0.43	40	25
E	0.90	-0.24	0.66	70	60

Порівняння інтегрального навантаження вузлів до і після розподілу



ВИСНОВКИ

У процесі дослідження було проведено:

1. Проведено всебічний аналіз існуючих методів балансування навантаження в мультимп'ютерних системах, що дозволило виокремити їх сильні та слабкі сторони, а також обґрунтувати необхідність пошуку нових адаптивних рішень.
2. Розроблено метод балансування навантаження який на відміну від відомих використовує кластеризацію вузлів за рівнем їхнього завантаження, що дозволяє знизити енергоспоживання, підвищити рівномірність розподілу задач та зменшити затримки в мультимп'ютерних системах.
3. Проведено оцінку ефективності балансування навантаження , підвищено ефективність розподілу задач між вузлами, зменшено ймовірність затримок окремих компонентів і покращити загальну стабільність роботи системи;
4. Досліджено результати моделювання та експериментального тестування, підтвердилось доцільність використання розробленого методу, продемонструвавши покращення показників продуктивності, рівномірності розподілу задач та зменшення затримок в мультимп'ютерних системах.

Завідувачу кафедри КІС
доктору філософії, доценту
Ользі ПАВЛОВІЙ

Гоби Каріни Сергіївни

ІІБ здобувача вищої освіти

ФІТ, 2 курсу, групи КІ2м-23-1

ЗАЯВА

З правилами чинного Положення «Про систему забезпечення академічної доброчесності у Хмельницькому національному університеті» від 01.07.2022, згідно з яким виявлення плагіату є підставою для відмови в допуску кваліфікаційної роботи до захисту та застосування заходів дисциплінарної та академічної відповідальності, ознайомлений(а). Про використання програмно-технічних засобів для перевірки кваліфікаційних робіт здобувачів вищої освіти на плагіат оповіщений(а) та надаю свою згоду на обробку та збереження університетом моєї роботи в інституційному репозитарії університету.

Також надаю університету право на передачу моєї роботи для обробки та збереження в базах даних програмно-технічних засобів (StrikePlagiarism та Anti-Plagiarism) та використання роботи для виявлення плагіату в інших роботах, які перевіряються програмно-технічними засобами та користувачами, що мають доступ до цих програмно-технічних засобів, виключно в обмежених цілях для виявлення плагіату в текстах робіт.

Робота для перевірки університетом надається в друкованому та електронному варіанті. Електронна версія моєї роботи збігається (ідентична) з друкованою.

21 травня 2025 року

Протокол аналізу звіту подібності експертом

Заявляю, що я ознайомився (-лась) з Повним звітом подібності, який був згенерований Системою виявлення і запобігання плагіату щодо роботи:

Автор: Каріна ГОБА

Співавтор:

Назва: Гоба_Метод розподілу навантаження в мультикомп'ютерних системах загального призначення згідно кластеризації вузлів за рівнем завантаженості ресурсів

Експерт:

Підрозділ: Кафедра комп'ютерної інженерії та інформаційних систем

Коефіцієнт подібності 1: 2.9%

Коефіцієнт подібності 2: 0.6%

Мікропробіли: 0

Заміна букв: 1

Інтервали: 0

Білі знаки: 1

Дата створення звіту: 2025-05-20 19:30:14.0

Після аналізу Звіту подібності констатую наступне:

- Запозичення, виявлені в роботі є законними і не є плагіатом. Рівень подібності не перевищує допустимої межі. Таким чином робота незалежна і приймається.
- Запозичення не є плагіатом, але перевищено граничне значення рівня подібностей. Таким чином робота повертається на доопрацювання.
- Виявлено запозичення і плагіат або навмисні текстові спотворення (маніпуляції), як передбачувані спроби укриття плагіату, які роблять роботу невідповідною вимогам законодавства (Ст. 32. ЗУ Про вищу освіту, пункт 3.1, Ст. 42. ЗУ Про освіту) та вимог НАЗЯВО (Критерій 5), а також кодексу етики і процедурам. Таким чином робота не приймається.

Обґрунтування:

2025-05-20

Дата



Доцент Андрій Нічепорук

експерт

РІШЕННЯ ЕКСПЕРТНОЇ КОМІСІЇ
КАФЕДРИ КОМП'ЮТЕРНОЇ ІНЖЕНЕРІЇ ТА ІНФОРМАЦІЙНИХ СИСТЕМ
ПРО ДОПУСК КВАЛІФІКАЦІЙНОЇ РОБОТИ ДО ЗАХИСТУ

Підтверджуємо ознайомлення з результатом звіту подібності щодо роботи, генерованою системою виявлення текстових збігів/ідентичності/схожості:

Назва: Метод розподілу навантаження в мультикомп'ютерних системах загального призначення згідно кластеризації вузлів за рівнем завантаженості ресурсів

Автор: Каріна Гоба

Спеціальність: 123 – Компютерна інженерія

Освітня програма: освітньо-наукова

Науковий керівник: Олег САВЕНКО, д.т.н., професор

Після аналізу звіту подібності зроблено такий висновок:

№	Висновок	Позначка про відповідність
1	Запозичення, виявлені в роботі, є законними і не є плагіатом. Робота приймається до захисту.	відповідає
2	Виявлені запозичення не є плагіатом, розміщені в розділах, які не описують безпосередньо авторське дослідження, але кількість цитат перевищує обсяг, виправданий поставленою метою роботи. Робота приймається до захисту, але має бути відкоригована. Відкоригований варіант має бути поданий на кафедру за 2 дні до захисту, разом із заявою щодо самостійності виконання письмової роботи та ідентичності друкованої та електронної версії роботи.	
3	Виявлені запозичення не є плагіатом, але частково розміщені в розділах, які описують безпосередньо авторське дослідження, а кількість цитат перевищує обсяг, виправданий поставленою метою роботи. В зв'язку з цим мета роботи та поставлені завдання не були досягнені. Робота може бути допущена до захисту (наступного року) після того як буде відкоригована та допрацьована і успішно пройде повторну перевірку на академічний плагіат.	
4	Робота містить навмисні текстові спотворення, передбачувані спроби укриття запозичень або інші прояви академічного плагіату. Робота містить фабрикацію або фальсифікацію даних. Робота не допускається до захисту.	

Підтвердження:

Запозичення, виявлені в роботі, є законними і не є плагіатом, оскільки:


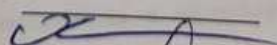
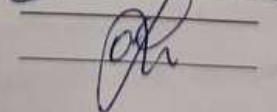
- окремі виявлені збіги є загальновживаними фразами або виразами, про що свідчить посилання системи на збіг з 10-40 джерелами на один фрагмент речення;
- всі зафіксовані системою ознаки модифікації тексту відносяться до комбінування латинських символів зі україномовними скороченнями індексів в формулах, що не є модифікацією тексту.

Сумарний обсяг всіх запозичень, визначений системою виявлення збігів/ідентичності/схожості, складає менше 4% і адресується до джерел з інтернету та бібліотеки, що, з урахуванням наведених обґрунтувань, відповідає характеру завдання і свідчить на користь кваліфікаційної роботи.

Керівник роботи

Гарант ОП

Завідувач кафедри КІС

Олег САВЕНКО

Олег САВЕНКО

Ольга ПАВЛОВА

МІНІСТЕРСТВО ОСВІТИ І НАУКИ УКРАЇНИ
ХМЕЛЬНИЦЬКИЙ НАЦІОНАЛЬНИЙ УНІВЕРСИТЕТ

РЕЦЕНЗІЯ НА ДИПЛОМНУ РОБОТУ

Дипломник: Каріна ГОБА

Тема: Метод розподілу навантаження в мультикомп'ютерних системах загального призначення згідно кластеризації вузлів за рівнем завантаженості ресурсів

Спеціальність: 123 «Комп'ютерна інженерія»

Обсяг кваліфікаційної роботи:

Кількість листів креслень -; кількість сторінок записки 73

1. Короткий зміст роботи та прийнятих рішень У роботі розроблено метод розподілу навантаження в мультикомп'ютерних системах загального призначення з використанням кластеризації вузлів за рівнем їхнього завантаження.

2. Висновок про відповідність роботи дипломному завданню _____
Кваліфікаційна робота відповідає виданому завданню

3. Характеристика виконання кожного розділу, ступінь використання останніх досягнень науки і техніки і передових методів роботи: У вступі подано об'єкт та предмет дослідження, мету, наукову новизну та практичну цінність роботи, а також характеристику структури роботи.

У першому розділі проведено аналіз існуючих методів балансування навантаження в мультикомп'ютерних системах.

У другому розділі здійснено розробку методу кластеризації вузлів на основі їхнього поточного завантаження.

У третьому розділі провели оцінку ефективності розподілу ресурсів при використанні запропонованого методу.

У четвертому розділі змоделивали та провели експериментальне тестування алгоритму на реальних або емуляційних даних.

У висновку підведено підсумки досягнення та ефективності методу та розв'язання задач дослідження

4. Позитивні сторони роботи: _____

5. Негативні сторони роботи: немає.

6. Оцінка графічного оформлення та пояснювальної записки роботи: -

7. Відгук про роботу в цілому: Робота виконана на належному рівні.

8. Інші зауваження: —

9. Оцінка дипломної роботи:

Розглянувши позитивні та негативні сторони представленої кваліфікаційної роботи вважаю, що робота заслуговує оцінки «відмінно» 4,75 (А)

Рецензент (прізвище, ім'я, по батькові, посада, місце роботи) Бедратюк Леонід Петрович, д.т.н., зав. кафедри ІІЗ ХНУ

“ 22 ” травня 2025р.



ХМЕЛЬНИЦЬКИЙ НАЦІОНАЛЬНИЙ УНІВЕРСИТЕТ

ГОЛОВІ ЕКЗАМЕНАЦІЙНОЇ КОМІСІЇ

Направляється студент Гоба Каріна Сергіївна на захист дипломного проєкту (роботи)
(прізвище, ім'я, по батькові)

за спеціальністю 123 - Комп'ютерна інженерія

На тему: Метод розподілу навантаження в мультикомп'ютерних системах загального призначення згідно кластеризації вузлів за рівнем завантаженості ресурсів

Дипломний проєкт (робота), рецензія з довідка про перевірку на плагіат додаються.

Декан факультету



Тетяна ГОВОРУЩЕНКО
(ім'я, прізвище)

ДОВІДКА УСПІШНОСТІ

Гоба К. С. за період навчання на факультеті інформаційних технологій з 2023 по 2025 роки повністю виконав навчальний план спеціальності з таким розподілом оцінок за: національною шкалою: відмінно 45,45 %, добре 18,18 %, задовільно 36,36 %. шкалою ЄКТС: А 36,84 %, В 5,26 %, С 15,79 %, D 5,26 %, E 36,84 %.

Методист факультету

[Signature]
(підпис)

[Signature]
(ім'я, прізвище)

ВИСНОВОК КЕРІВНИКА ДИПЛОМНОГО ПРОЄКТУ (РОБОТИ) ТА ОБГРУНТУВАННЯ ОЦІНКИ

Студентка Каріна Гоба виконала кваліфікаційну роботу на кавалетному науково-дослідницькому рівні

Оцінка дипломного проєкту (роботи) 4,75 / A / хороше

Керівник дипломного проєкту

[Signature]
(підпис)

[Signature]
(ім'я, прізвище)

" 01 " 05 2025 р.

ВИСНОВОК КАФЕДРИ ПРО ДИПЛОМНИЙ ПРОЄКТ (РОБОТУ)

Дипломний проєкт (роботу) розглянуто. Студент Гоба К. С. допускається до захисту цього проєкту (роботи) в екзаменаційній комісії.

Завідувач кафедри

КІІС

(назва)

[Signature]
(підпис, ім'я, прізвище)

" 01 " 05 2025 р.