

ГОБА КАРІНА

Хмельницький національний університет

<https://orcid.org/0009-0008-8553-1674>e-mail: karina2002goba@gmail.com**ЖУКОВСЬКИЙ ПAVЛЮ**

Хмельницький національний університет

<https://orcid.org/0009-0007-3461-9919>e-mail: 777reiste777@gmail.com**КОВАЛЬЧУК VАСИЛЬ**

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

Khmelnitskyi National University

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 validate the effectiveness of the proposed method, the study utilizes performance metrics such as task completion time, system throughput, and energy efficiency, comparing them across various simulated scenarios.

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.

Стаття надійшла до редакції / Received 27.04.2025

Прийнята до друку / Accepted 16.04.2025

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 простою реалізацією [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 $P = \{PM_1, PM_2, \dots, PM_M\}$, кожен з яких містить набір програмно або апаратно ізольованих обчислювальних одиниць (віртуалізованих або фізичних), які в рамках моделі позначимо як $VM = \{VM_1, VM_2, \dots, VM_L\}$. Система оперує множиною завдань користувача $T = \{T_1, T_2, \dots, T_S\}$, 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 даних, вільно закодованих даних. Архітектура CNN чергує згорткові шари, шари субдискретизації (pooling) та повнозв'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 згорткова 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

мультимедійній системі в режимі реального часу було використано рекурентну нейронну мережу (RNN), яка добре зарекомендувала себе у задачах обробки часових рядів. На відміну від класичних багатошарових 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 функції втрат відносно кожного шару. Ключовою особливістю моделі є інтеграція архітектур CNN та RNN, де згортовка 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 проектування

Під час розробки методу розподілу навантаження для мультимедійних 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 часу, а також потребу у періодичній кластеризації, доцільно структурувати систему за принципами багатошарової (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 ефективного реалізації обраного підходу використовується набір шаблонів проектування, що дозволяє вирішити типові задачі, властиві розподіленим 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 реалізовувати 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) організує 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.



Рис.1. Трирівнева (3-layer) архітектура

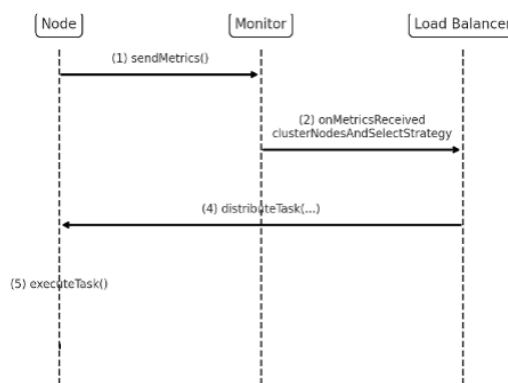


Рис.2. Спрощена діаграма послідовності

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, кондиціювання), 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, кондиціювання); 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 навантажених 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.

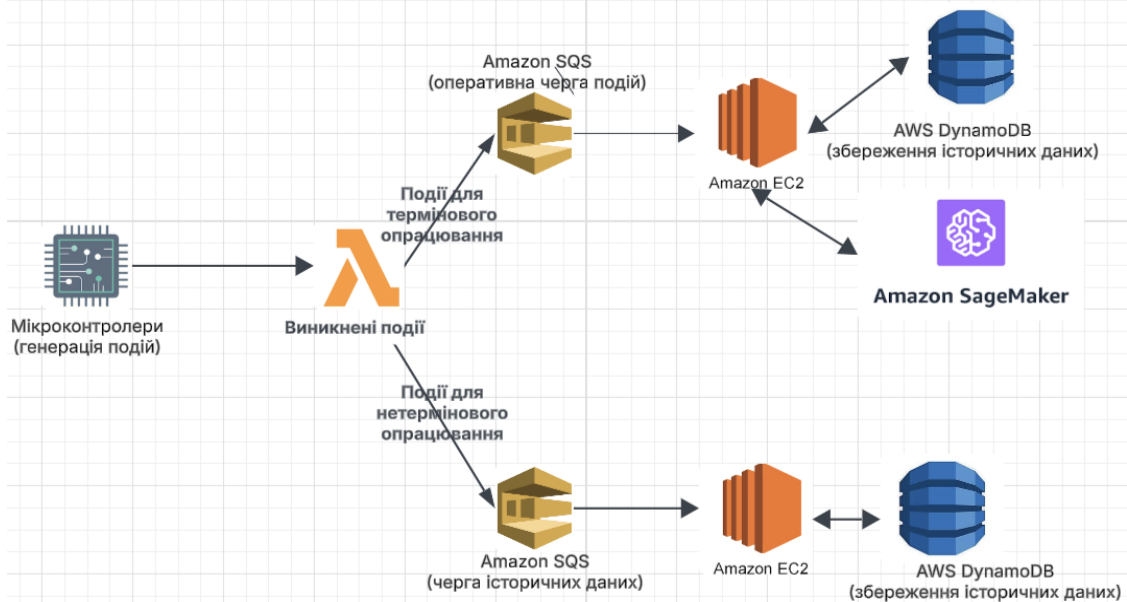


Рис.3. Схема розподілу навантаження подій між чергами

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. Використовуємомсервіс 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 недозавантаженими. 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).

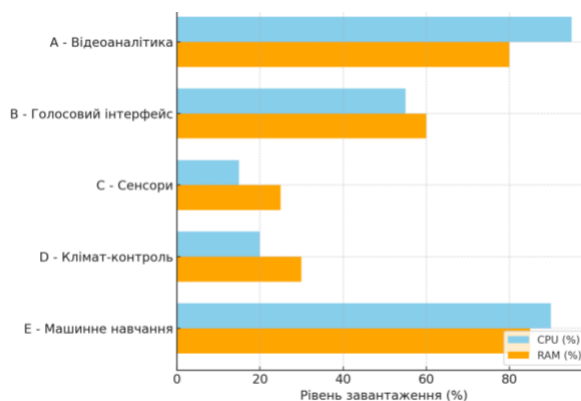


Рис.4. Горизонтальна порівняльна діаграма навантаження вузлів

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 отримують 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:

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

де: $W_{надлишок}$ — 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:

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

Аналогічно для вузла Е:

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

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:

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

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

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.

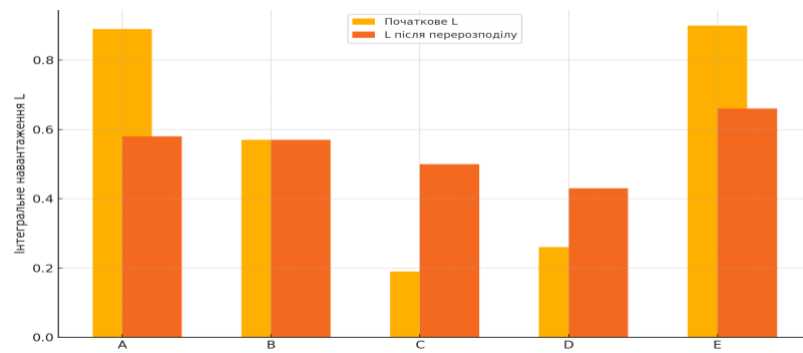


Рис. 5. Порівняння навантаження на вузли до та після балансування

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.

Література

1. Alon, N. (2020). *Distributed computing and scalability challenges*. *ACM Transactions*.
2. Bedratyuk, L., & Savenko, O. (2018). MATCH Commun. Math. Comput. Chem., 79, 407–414.
3. Brownlee, N. (2021). Simplified static methods for resource distribution. *Computing Journal*.
4. Coulouris, G., Dollimore, J., Kindberg, T., & Blair, G. (2020). *Distributed systems: Concepts and design* (5th ed.). Pearson Education.
5. Dean, J., & Ghemawat, S. (2020). MapReduce: Simplified data processing on large clusters. *Communications of the ACM*.
6. Foster, I. (2020). Challenges in static load balancing. *Elsevier Distributed Systems Review*.
7. Foster, I., & Kesselman, C. (2021). *The grid: Blueprint for a new computing infrastructure* (2nd ed.). Elsevier.
8. Fox, G. C. (2020). Challenges in managing distributed resources. *Journal of Parallel and Distributed Computing*.
9. Fujitsu Ltd. (2021). *Fugaku: Revolutionizing high-performance computing*. Fujitsu Whitepapers.
10. Ghosh, S. (2019). *Applications of distributed computing systems*. Wiley.
11. Ghosh, S. (2020). Static load balancing in distributed computing. *Journal of Parallel and Distributed Computing*.
12. Jassy, A. (2022). *Building scalable applications with AWS*. AWS Whitepapers.
13. Kruger, J. (2021). Analysis of round-robin algorithms for load balancing. *ACM Transactions on Computing Systems*.
14. Lee, G. (2020). *Titan supercomputer: Applications and architecture*. Oak Ridge National Laboratory.
15. Nguyen, T. (2020). Random assignment strategies in high-performance computing. *IEEE Computing Journal*.
16. Simonenko, V. P., Dekhtyaruk, M. T., & Zabara, S. S. (2014). *Programne zabezpechennia kompiuternykh merezh* [Програмне забезпечення комп'ютерних мереж]. Університет «Україна». (Підручник під грифом МОН України).
17. Smith, J. (2020). *High-performance distributed systems*. Springer.
18. Tanenbaum, A. S., & Van Steen, M. (2021). *Distributed systems: Principles and paradigms* (2nd ed.). Pearson.

References

1. Alon, N. (2020). *Distributed computing and scalability challenges*. *ACM Transactions*.
2. Bedratyuk, L., & Savenko, O. (2018). MATCH Commun. Math. Comput. Chem., 79, 407–414.
3. Brownlee, N. (2021). Simplified static methods for resource distribution. *Computing Journal*.
4. Coulouris, G., Dollimore, J., Kindberg, T., & Blair, G. (2020). *Distributed systems: Concepts and design* (5th ed.). Pearson Education.
5. Dean, J., & Ghemawat, S. (2020). MapReduce: Simplified data processing on large clusters. *Communications of the ACM*.
6. Foster, I. (2020). Challenges in static load balancing. *Elsevier Distributed Systems Review*.
7. Foster, I., & Kesselman, C. (2021). *The grid: Blueprint for a new computing infrastructure* (2nd ed.). Elsevier.
8. Fox, G. C. (2020). Challenges in managing distributed resources. *Journal of Parallel and Distributed Computing*.
9. Fujitsu Ltd. (2021). *Fugaku: Revolutionizing high-performance computing*. Fujitsu Whitepapers.
10. Ghosh, S. (2019). *Applications of distributed computing systems*. Wiley.
11. Ghosh, S. (2020). Static load balancing in distributed computing. *Journal of Parallel and Distributed Computing*.
12. Jassy, A. (2022). *Building scalable applications with AWS*. AWS Whitepapers.
13. Kruger, J. (2021). Analysis of round-robin algorithms for load balancing. *ACM Transactions on Computing Systems*.
14. Lee, G. (2020). *Titan supercomputer: Applications and architecture*. Oak Ridge National Laboratory.
15. Nguyen, T. (2020). Random assignment strategies in high-performance computing. *IEEE Computing Journal*.
16. Simonenko, V. P., Dekhtyaruk, M. T., & Zabara, S. S. (2014). *Programne zabezpechennia kompiuternykh merezh* [Програмне забезпечення комп'ютерних мереж]. Університет «Україна». (Підручник під грифом МОН України).
17. Smith, J. (2020). *High-performance distributed systems*. Springer.
18. Tanenbaum, A. S., & Van Steen, M. (2021). *Distributed systems: Principles and paradigms* (2nd ed.). Pearson.