

О. А. Двірна, С. В. Набока

Національний університет «Полтавська політехніка імені Юрія Кондратюка», Полтава, Україна

ОПТИМІЗАЦІЯ ПРОДУКТИВНОСТІ ХМАРНИХ СЕРВІСІВ: МЕТОДИ ТА ЇХ ЕФЕКТИВНІСТЬ

Анотація. Актуальність. Оптимізація продуктивності хмарних сервісів набуває критичного значення у 2025–2026 роках через експоненційний ріст обсягів даних, що обробляються в цифрових екосистемах, та необхідність швидкої адаптації до динамічних навантажень, спричинених впровадженням штучного інтелекту (ШІ) та IoT-пристроїв. Згідно з прогнозами, глобальний ринок хмарних обчислень зростає на 20–25% щорічно, а в Україні цей тренд посилюється державними ініціативами, такими як Національна стратегія розвитку ШІ до 2025 року, яка визначає хмарні сервіси як ключову інфраструктуру. Однак, часті інциденти з перевантаженнями, витратами та безпекою, особливо в умовах енергетичного дефіциту та геополітичних ризиків, свідчать про фрагментарність сучасних практик оптимізації. У бізнесі та державному управлінні України, де хмарні технології впроваджуються для цифровізації, неефективне управління ресурсами призводить до втрат до 30–40% бюджету на ІТ. Це робить тему стратегічно важливою для забезпечення конкурентоспроможності, сталого розвитку та стійкості до кіберзагроз у мультихмарних середовищах. **Об'єктом дослідження** є сучасні хмарні сервіси, включаючи IaaS, PaaS та SaaS-моделі провайдерів на кшталт AWS, Google Cloud, Azure та локальних українських платформ, з акцентом на їх продуктивність під змінними навантаженнями. Особлива увага приділяється українському ринку, де гібридні хмарні рішення поєднуються з локальною інфраструктурою для подолання обмежень інтернет-пропускної здатності та регуляторних бар'єрів. Дослідження охоплює ключові аспекти: динамічний розподіл ресурсів, балансування навантаження, автоматичне масштабування та інтеграцію ШІ для прогнозування попиту. **Мета** – систематизувати методи оптимізації, виявити виклики масштабованості, безпеки й енергоефективності, а також розробити рекомендації для впровадження в бізнесі (fintech, e-commerce) та державному секторі (електронне врядування, критична інфраструктура) з урахуванням специфіки України. **Методологія.** У статті застосовано комплексний підхід, що включає порівняльний аналіз сучасних методів оптимізації на основі даних реаналізу хмарних провайдерів (AWS Auto Scaling, Google Cloud Autoscaler, Azure Autoscale) та кейсів впровадження в 2025 році. Використано експертну оцінку ризиків з факторною моделлю, яка класифікує причини неефективності (людський фактор, перевантаження, слабе шифрування), а також математичне моделювання оптимізації через метрики SLO/SLA, лінійне програмування для балансу витрат і продуктивності. Додатково проведено аналіз міжнародних стандартів (ISO 27001 для безпеки, Green Cloud для енергоефективності) та регіональних даних України, включаючи чисельні експерименти з ШІ-алгоритмами прогнозування (машинне навчання для autoscaling). Кількісні оцінки базуються на статистиці: зменшення витрат на 30% при AI-автоматизації, моніторинг через Google Cloud Monitoring та моделювання енерговитрат дата-центрів. **Результати.** Дослідження підтвердило, що ключовими методами оптимізації є динамічний розподіл ресурсів з AI-прогнозуванням, який зменшує over-provisioning на 25–40%, та автоматичне масштабування, що реагує на пікові навантаження за секунди, як у випадку Pinterest з AWS (зниження витрат на 30%). Балансування навантаження та serverless-архітектури підвищують швидкість на 38–60%, усуваючи простой. Виклики включають безпеку (неправильне IAM, сліпі зони в ефемерних ресурсах), енергоефективність (зростання споживання на 40% через ШІ) та масштабування в Україні (обмежена інфраструктура, тіньова ІТ). Порівняння провайдерів показало перевагу гібридних моделей: Kubernetes-оркестрація з multi-cloud скорочує latency на 40%, а квантування нейромереж – обчислювальні витрати на 40%. В українському контексті виявлено нерівномірність впровадження: бізнес економить на адмінах, але держсектор страждає від енергодефіциту. **Висновки.** Запропоновано комплексний підхід до оптимізації: інтеграцію AI для predictive scaling, стандартизацію BRM-подібних процедур для хмар (адаптивні протоколи ризик-менеджменту), гармонізацію з локальними дата-центрами та відновлюваними джерелами енергії для зниження CO₂ на 30%. **Рекомендації для України:** впровадження autoscaling down у неробочий час, тестування безпеки (CSPM, SIEM), моніторинг SLO/SLA та тренінги для ШІ-управління ресурсами. Доведено доцільність гібридних рішень для сталого розвитку цифрової екосистеми, що забезпечує надійність, економічність і конкурентоспроможність у бізнесі та держуправлінні, з потенціалом скорочення витрат на 30–50% та підвищення продуктивності. Цей підхід створює основу для прогнозування навантажень і мінімізації ризиків у 2026 році та далі.

Ключові слова: динамічний розподіл ресурсів, балансування навантаження, автоматичне масштабування, машинне навчання, мультикритеріальна оптимізація, енергоефективність, масштабованість, надійність.

Вступ

Актуальність теми оптимізації продуктивності хмарних сервісів обумовлена стрімким розвитком цифрових технологій і все більшим обсягом даних та обчислювальних завдань, які переносяться у хмарні середовища. У 2025 році хмарні технології стали стратегічною платформою для бізнесу та організацій, оскільки дозволяють забезпечувати масштабованість, гнучкість та економічну ефективність ІТ-інфраструктури. Згідно [1] частка використання хмарних ресурсів у різних галузях на 2025 рік становить від 40 % до

90 %, усереднені показники подані на рис. 1. Офіційна статистика використання хмарних ресурсів в Україні у 2025 році свідчить про поступове зростання популярності хмарних технологій у різних галузях. Згідно з даними Державної служби статистики України та експертних досліджень, підприємства, зайняті у виробництві комп'ютерів, електроніки та програмуванні, збільшують частку використання хмарних послуг, прогнозуючи зростання долі користувачів із 12% у 2020 році до понад 20% у 2025 році у виробництві, хоча дещо спостерігається зниження у сфері програмування.

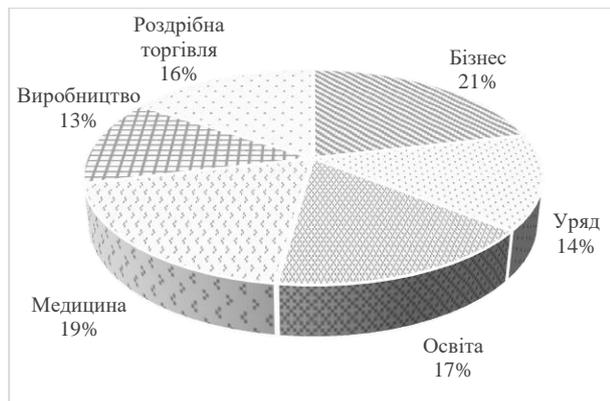


Рис. 1. Частка використання хмарних ресурсів у різних галузях на 2025 р. (%)

Урядові установи також активно впроваджують хмарні сервіси у свої процеси, зі зростаючою увагою до стандартизації та безпеки, що було законодавчо врегульовано у 2022 році. У сфері освіти та медицини хмарні технології використовуються для дистанційного навчання, зберігання великих обсягів даних і телемедицини, що особливо актуально в умовах цифрової трансформації країни [2].

Опитування IT-фахівців в Україні показують, що більше 90% використовують хмарні сервіси для зберігання і обробки даних, що впливає на робочі процеси і бізнес-моделі. Проте рівень обізнаності про хмарні рішення серед керівників поки що недостатній, що є бар'єром для швидшого впровадження. Водночас 94% респондентів планують розширювати використання хмар у майбутньому, визнаючи їх ключову роль в підвищенні продуктивності і гнучкості бізнесу.

Таким чином, ринок хмарних послуг в Україні перебуває у стадії активного розвитку з позитивною динамікою використання в промисловості, державному секторі, освіті і медицині. Проте для повного розкриття потенціалу необхідно підвищувати рівень професійної обізнаності, інвестувати у безпеку та стандартизацію, та сприяти інтеграції хмарних технологій в бізнес- і державні процеси.

Проте збільшення складності і масштабів хмарних систем породжує нові виклики у сфері управління ресурсами, продуктивністю та витратами. Саме тому оптимізація продуктивності стає ключовим фактором для підтримки високої якості сервісів, зменшення затримок, раціонального використання обчислювальних потужностей та зниження операційних витрат.

Важливо також відзначити, що інновації у сфері машинного навчання, автоматичного масштабування, динамічного розподілу ресурсів і балансування навантаження відкривають нові можливості для підвищення ефективності хмарних платформ. У зв'язку з цим, розробка і впровадження сучасних методів оптимізації мають велике значення для адаптації хмарних сервісів до постійно змінних умов, збільшення надійності та енергозбереження. Тема є актуальною для науковців і практиків, оскільки забезпечує фундамент для подальшого розвитку високопродуктивних, економічно вигідних та надійних хмарних тех-

нологій, які є критично важливими у сучасному цифровому світі.

Метою статті є аналіз та систематизація основних методів оптимізації продуктивності хмарних сервісів, а також оцінка їх ефективності з точки зору підвищення швидкодії, надійності та економії ресурсів у хмарних обчислювальних системах. Стаття спрямована на розробку рекомендацій щодо застосування сучасних алгоритмів та технологій, таких як автоматичне масштабування, балансування навантаження і оптимізація розподілу ресурсів, для досягнення максимального рівня продуктивності і ефективності у різних сценаріях використання хмарних сервісів.

Сучасний стан проблеми оптимізації продуктивності хмарних сервісів характеризується зростаючою складністю і масштабністю хмарних інфраструктур, що підштовхує до пошуку нових підходів і технологій для ефективного управління ресурсами. У 2025 році основними трендами стають перехід до гібридних і мультихмарних моделей, які поєднують публічні, приватні та виділені сервери для досягнення балансу між продуктивністю, безпекою і вартостями. Це дозволяє більш гнучко розподіляти робочі навантаження та уникати ризику завищених витрат або "vendor lock-in" (залежності від одного постачальника).

Також активно впроваджуються безсерверні архітектури (serverless), які автоматично масштабуються відповідно до навантаження, що сприяє оптимізації витрат і підвищенню швидкодії застосунків. Відбувається інтеграція машинного навчання та аналітики для точного передбачення навантажень та динамічного розподілу ресурсів, що дозволяє хмарним платформам адаптуватися до мінливих умов в реальному часі. Однак, потреба у балансі між масштабованістю, продуктивністю, надійністю та економією енергоресурсів залишається викликом, що визначає актуальність подальших досліджень у цій сфері.

Огляд наукових досліджень з оптимізації продуктивності хмарних сервісів у період 2020–2025 років вказує на інтенсивне зростання важливості цієї теми у зв'язку з поширенням цифрової трансформації і зростанням обсягів даних. Особливою увагою користуються дослідження, що спрямовані на адаптивне управління ресурсами, балансування навантаження і динамічне масштабування, що є ключовими факторами для збереження продуктивності та доступності хмарних сервісів [3].

Значний внесок у розвиток методів оптимізації внесли роботи, які розглядають застосування алгоритмів машинного навчання та штучного інтелекту для прогнозування навантажень та автоматичного регулювання розподілу ресурсів. Це дозволяє хмарним платформам самостійно адаптуватися у режимі реального часу, оптимізуючи витрати і підвищуючи надійність систем [4].

Дослідження також акцентують увагу на необхідності комплексного підходу до оптимізації, поєднуючи як технічні рішення, так і управлінські практики. Наприклад, у роботах вказується на важливість інтеграції мультихмарних і гібридних моделей, що

дозволяють комбінувати переваги різних типів хмарних середовищ та знижувати операційні ризики [5].

Окремо виділяються проблеми міграції у хмару, що часто супроводжуються значними технічними та організаційними викликами, включаючи безпеку, сумісність і витрати. Розробка стандартизованих планів міграції та використання досвіду провідних постачальників хмарних послуг є важливою складовою успішної оптимізації хмарних платформ [6].

Зростає роль аналітики й прогнозування у прийнятті рішень, що підвищує ефективність роботи хмарних систем. Впровадження розширених аналітичних та предиктивних моделей дозволяє компаніям краще розуміти поведінку своїх хмарних сервісів, оптимізувати алгоритми розподілу навантаження та зменшувати простій [7].

Крім того, огляди наукової літератури вказують на те, що постійне ускладнення і масштабування хмарних середовищ висуває нові вимоги до алгоритмів оптимізації, зокрема у сфері енергоефективності та екологічності. У сучасних системах увага приділяється не лише продуктивності, але і зменшенню споживання енергії, що стає важливим фактором сталого розвитку в інформаційних технологіях [8].

Таким чином, сучасний стан проблеми свідчить про активне впровадження інноваційних методів і технологій, котрі поєднують алгоритмічну оптимізацію, машинне навчання та управлінські методики для досягнення ефективного та надійного функціонування хмарних сервісів у складних і динамічних умовах сучасного цифрового середовища.

Теоретичні основи оптимізації продуктивності хмарних сервісів

Теоретичні основи оптимізації продуктивності хмарних сервісів базуються на визначенні комплексних критеріїв оцінювання якості роботи та ефективності ресурсів, що забезпечують виконання обчислювальних задач у хмарних середовищах. Основними показниками продуктивності є надійність, час відгуку, ефективність використання ресурсів, масштабованість, доступність і вартість обслуговування. Надійність визначає ймовірність безперервного функціонування сервісу без збоїв, що особливо важливо для критичних бізнес-додатків. Час відгуку характеризує швидкість реакції системи на запити користувачів і впливає на користувацький досвід. Ефективність ресурсів показує, наскільки оптимально використовуються обчислювальні потужності, пам'ять, пропускна здатність мережі та інші компоненти [9].

Відповідно до моделі якості хмарних сервісів NIST, оцінювання продуктивності враховує як технічні, так і організаційні аспекти: безпеку даних і додатків, доступність сервісу, масштабованість, а також можливість швидко адаптуватись до зміни навантаження. Методи системного оцінювання якості включають статистичний аналіз, SWOT-аналіз, моделювання і симуляції поведінки системи у різних сценаріях, що забезпечує гнучкий і глибокий аналіз сильних і слабких сторін хмарних сервісів [6].

Продуктивність хмарних сервісів тісно пов'язана з інфраструктурою віртуалізації, яка забезпечує

ізоляцію та ефективне розподілення ресурсів між користувачами. Важливим є підтримання балансу між перевантаженням ресурсів і нерівномірним використанням, що досягається за допомогою алгоритмів динамічного масштабування та балансування навантаження. Ці методи дозволяють забезпечити стабільне обслуговування навіть за умов різких коливань навантажень, що є характерним для сучасних хмарних середовищ [6, 7]. Масштабованість системи – це здатність швидко і ефективно збільшувати або зменшувати обчислювальні ресурси відповідно до поточної потреби. Вища масштабованість сприяє підвищенню загальної продуктивності та зниженню експлуатаційних витрат, оскільки ресурси надаються у потрібній кількості та у потрібний час. Сучасні підходи до оптимізації включають автоматичне масштабування, яке базується на аналізі метрик використання ресурсів у реальному часі [8].

Загалом, теоретичні основи оптимізації продуктивності хмарних сервісів ґрунтуються на комплексному підході до оцінки їх роботи, який поєднує математичне моделювання, статистичні методи, машинне навчання та інтелектуальне управління ресурсами для досягнення максимальної ефективності і надійності хмарних платформ.

Моделі та методи оптимізації продуктивності хмарних сервісів

Моделі та математичні методи оптимізації продуктивності хмарних сервісів спираються на мультикритеріальний підхід, який дозволяє одночасно враховувати кілька цілей, таких як максимізація надійності, мінімізація часу відгуку та вартості, а також оптимальне використання ресурсів. Одним із найбільш поширених підходів є створення багатокритеріальних моделей оптимізації, де завдання формалізується у вигляді функції кількох об'єктивних цілей, які потрібно мінімізувати або максимізувати при дотриманні певних обмежень. Прикладом такої моделі є мультиоб'єктивне програмування, у якому розглядаються взаємозалежні критерії та шукаються компромісні рішення з використанням технік, як-от метод вагових коефіцієнтів або метод Парето [9].

Математична задача оптимізації обчислень у хмарних середовищах зазвичай формулюється як багатокритеріальна оптимізація, мета якої – одночасно мінімізувати витрати ресурсів (CPU, пам'ять, пропускна здатність) і час виконання, при максимізації надійності та продуктивності.

Загальна формулювання має вигляд: знайти $x \in X$, при якому:

$$F(x) = (f_1, f_2, -f_3) \rightarrow \min$$

при обмеженнях:

$$\begin{cases} g_1(x) \leq G_1, \\ g_2(x) \leq G_2, \\ g_3(x) \leq G_3, \\ x \in X, \end{cases}$$

де $f_1(x)$ – функція вартості або споживання обчислювальних ресурсів за варіантом розподілу x ;

$f_2(x)$ – час виконання задачі; $f_3(x)$ – показник надійності (відмовостійкості); $g_1(x), g_2(x), g_3(x)$ – функції обмеження, а G_1, G_2, G_3 встановлені допустимі значення за максимальною вартістю, максимальним часом і мінімальною надійністю відповідно; x – вектор рішень, що визначає план розподілу ресурсу; X – допустимий простір розв'язків.

Часто простір розв'язків для таких задач є комбінаторною множиною перестановок чи розміщень, тоді мова йде про багатокритеріальну комбінаторну задачу. Методи розв'язування оптимізаційних проблем такого класу описані у роботах [10, 11].

Алгоритми планування у хмарних системах включають методи динамічного розподілу ресурсів, балансування навантаження та автоматичного масштабування. Планування задач часто базується на алгоритмах ранжування, евристичних методах та методах машинного навчання, які дозволяють прогнозувати навантаження на ресурси та відповідно адаптувати модель розподілу для підвищення ефективності. Крім того, алгоритми планування враховують специфіку хмарної архітектури, включаючи віртуалізацію, багатокористувацький режим і енергозбереження, що ускладнює класичні підходи і вимагає використання гібридних моделей із застосуванням аналітики даних [12].

Важливим теоретичним інструментом є закон Амдала, який встановлює межі масштабованості обчислювальних систем і впливає на оптимізацію багатоядерних платформ у хмарних середовищах. Таким чином, математичне моделювання використовується не лише для підвищення продуктивності, але й для оцінки ефективності паралельних обчислень і балансування навантаження між ресурсами [13].

Для точного прогнозування продуктивності у багатопотокових і розподілених системах застосовуються статистичні розподіли, зокрема розподіл Парето і розподіл Зіпфа, які моделюють характерні закономірності використання ресурсів. Це дозволяє розробляти більш точні моделі алгоритмів кешування, розподілу пам'яті та пропускної здатності, оптимізуючи при цьому обробку запитів у хмарі [14].

Сучасні математичні методи також включають використання систем моніторингу і аналізу великих даних (Big Data Analytics) для формалізації та автоматизації прийняття рішень у реальному часі. Інструменти на зразок Prometheus, Grafana і Elastic Stack дозволяють візуалізувати метрики продуктивності, виявляти потенційні вузькі місця у роботі системи та адаптувати алгоритми розподілу ресурсів відповідно до отриманих даних [15].

Загалом, мультикритеріальні моделі та алгоритми планування формують теоретичну основу для побудови адаптивних, масштабованих і високопродуктивних хмарних платформ, здатних задовольнити зростаючі вимоги користувачів і бізнесу у динамічних умовах сучасного ІТ-середовища. Моделі та математичні методи оптимізації продуктивності хмарних сервісів спираються на мультикритеріальний підхід, який дозволяє одночасно враховувати кілька цілей, таких як максимізація надійності, мінімізація

часу відгуку та вартості, а також оптимальне використання ресурсів. Одним із найпоширеніших є багатокритеріальне програмування, де шукаються компромісні рішення з урахуванням взаємозалежних об'єктивних функцій за допомогою методів вагових коефіцієнтів, методу Парето та інших технік [9- 11].

Алгоритми планування включають методи динамічного розподілу ресурсів, балансування навантаження, автоматичного масштабування і прогнозування навантажень із застосуванням евристичних методів і машинного навчання. Такі підходи враховують специфіку віртуалізації, багатокористувацький режим і енергозбереження, що потребує гібридних моделей та аналітики даних для підвищення ефективності [16].

Закон Амдала використовується для визначення меж масштабованості хмарних систем і оцінки продуктивності багатоядерних систем, що є важливою складовою при оптимізації паралельних обчислень і ресурсів [13].

Для аналізу використання ресурсів застосовують статистичні розподіли Парето і Зіпфа, що допомагають моделювати закономірності у споживанні пам'яті й пропускної здатності та оптимізувати керування кешуванням і розподілом даних [17].

Сучасні методи включають інструменти моніторингу (Prometheus, Grafana, ELK) для збору та аналізу метрик у реальному часі, що дозволяє автоматизувати прийняття рішень і своєчасно адаптувати алгоритми розподілу для забезпечення стабільності та продуктивності [18].

У сукупності мультикритеріальні моделі і алгоритми планування формують теоретичну базу для створення масштабованих, гнучких і високоефективних хмарних платформ, здатних відповідати сучасним викликам і потребам користувачів.

Порівняльний аналіз ефективності методів оптимізації хмарних сервісів

Порівняльний аналіз ефективності методів оптимізації хмарних сервісів базується на кількох ключових критеріях: продуктивності, надійності, собівартості та енергозбереженні. Продуктивність оцінюється через швидкість обробки запитів, час відгуку та стабільність роботи під навантаженням. Надійність включає безперервність надання послуг і стійкість до відмов. Собівартість враховує витрати на використання обчислювальних ресурсів, включно з оплатою за споживання і технічне обслуговування. Енергозбереження стає все більш актуальним критерієм через зростаючий вплив ІТ-інфраструктури на енергоспоживання і екологію систем.

Загалом, методи динамічного розподілу ресурсів, які адаптують обчислювальну потужність залежно від поточного навантаження, показують високу ефективність у підвищенні продуктивності та зменшенні енергоспоживання. Балансування навантаження дозволяє більш рівномірно використовувати ресурси, знижуючи ризик перевантаження окремих вузлів і тим самим підвищуючи надійність.

Практичним кейсом успішного застосування оптимізаційних методів є дослідження в Японії, де

впроваджуються методи динамічного розподілу, автоматичного масштабування і балансування навантаження з підтримкою машинного навчання для прогнозування майбутніх потреб.

Це дозволяє суттєво знизити собівартість і підвищити сталість роботи хмарних систем, водночас зберігаючи високі показники продуктивності [3, 7, 16].

Інший приклад – застосування гібридних методів управління ресурсами, які поєднують централізоване та децентралізоване управління, що забезпечує гнучкість і адаптивність хмарних платформ під час різких змін навантаження.

Експериментальні дослідження засвідчили, що ці підходи знижують затримки, покращують використання процесорів та суттєво оптимізують енергоспоживання [18].

Також сучасні платформи, такі як AWS, Azure і Google Cloud, впроваджують інструменти автоматичного масштабування, які зміцнюють адаптивність інфраструктури, підвищують продуктивність і одночасно контролюють операційні витрати. Проте вибір оптимального способу розподілу і масштабування залежить від специфіки задачі, обсягів обробки і вимог до надійності.

Порівняно з класичними підходами, новітні методи з використанням машинного навчання і предиктивної аналітики забезпечують більш тонке та ефективне регулювання ресурсів, дозволяючи підвищити продуктивність системи при зниженні витрат і збільшенні енергоефективності. Водночас це створює певні виклики у плані складності реалізації та потреби у великому обсязі даних для навчання моделей [15]. Результати аналізу узагальнено та подано у табл. 1.

Таблиця 1 – Порівняльна таблиця ефективності основних методів оптимізації хмарних сервісів

Метод оптимізації	Продуктивність	Надійність	Собівартість	Енергозбереження	Практичні кейси використання
Динамічний розподіл ресурсів	Висока адаптивність до навантажень	Висока завдяки рівномірному навантаженню	Зниження витрат через уникнення перевитрат	Середня, зниження пікових навантажень	Японія: впровадження з ML та прогнозування
Балансування навантаження	Покращення завдяки рівномірному навантаженню	Зниження ризиків відмов	Низькі додаткові витрати	Підвищення енергоефективності	Комерційні хмарні платформи, гібридні моделі
Автоматичне масштабування	Швидка адаптація до змін попиту	Вища надійність завдяки адаптивності	Оптимізація витрат за рахунок масштабування	Високий потенціал енергозбереження	AWS, Azure, Google Cloud
Машинне навчання та аналітика	Точне прогнозування навантаження	Покращення діагностики та моніторингу	Вища початкова собівартість впровадження	Значна економія за рахунок прогнозів	Практики Японії, США, застосування AI/ML

Отже, аналіз свідчить, що комплексне впровадження адаптивних алгоритмів, підтриманих аналітикою і штучним інтелектом, є найперспективнішим шляхом підвищення ефективності хмарних сервісів, що надалі зумовить їх надійність, продуктивність та економічну доцільність.

Обговорення та перспективи

Сучасні тенденції розвитку обчислювальних технологій свідчать, що хмарні обчислення залишатимуться ключовим елементом IT-стратегій у найближчі роки, зокрема завдяки впровадженню багатохмарних та гібридних рішень, які дозволяють підвищити гнучкість, стійкість систем та знизити залежність від окремих постачальників.

Перспективним напрямком розвитку є інтенсивна інтеграція штучного інтелекту і машинного навчання для оптимізації керування ресурсами, прогнозування навантаження та динамічного масштабування, що дозволить підвищити продуктивність і надійність хмарних платформ, одночасно знижуючи операційні витрати та енергоспоживання.

Рекомендується впроваджувати сучасні методи оптимізації, які поєднують динамічний розподіл ресурсів, балансування навантаження та інтелектуальний аналіз даних у промислових хмарних платформах. Акцент слід робити на автоматичному масштабуванні з використанням машинного навчання для адаптивного реагування на зміни в навантаженні та

потребах користувачів. Крім того, важливо інвестувати в енергоефективну інфраструктуру, впроваджувати політики сталого використання ресурсів та інтегрувати багатохмарні рішення для підвищення безперервності та безпеки обслуговування.

Проте масштабованість і складність сучасних хмарних систем створюють низку викликів. Збільшення обсягів обробки даних, різноманітність користувачів і додатків, а також потреба у високій доступності вимагають складних алгоритмів планування та управління ресурсами, що значно ускладнює адміністрування і підвищує ризики виникнення вузьких місць у системах. До того ж, зростає важливість забезпечення кібербезпеки та збереження конфіденційності, що вимагає постійного оновлення політик і технологій захисту даних. Ще однією проблемою є необхідність балансування між продуктивністю і енергоефективністю. Хмарні центри обробки даних потребують значної кількості енергії, тому інтеграція рішень для енергозбереження зберігає стратегічне значення. Використання периферійних обчислень (edge computing) частково знижує навантаження на центральні ресурси і зменшує затримки, покращуючи тим самим якість обслуговування.

Таким чином, подальший розвиток оптимізації хмарних обчислень має базуватися на комплексному поєднанні інтелектуальних алгоритмів, енергоефективних технологій і багатохмарних архітектур. Це дозволить зберегти баланс між продуктивністю, надійністю,

економічності і безпекою, що є критично важливим для сучасних ІТ-інфраструктур і бізнес-процесів.

Висновки

Узагальнення результатів дослідження свідчить, що оптимізація продуктивності хмарних сервісів є критично важливою складовою сучасних ІТ-інфраструктур, що дозволяє підвищити їх ефективність, надійність та економічність.

Запропоновані методи, що включають динамічний розподіл ресурсів, балансування навантаження, автоматичне масштабування та застосування алгоритмів машинного навчання, довели свою ефективність у зниженні часу відгуку, витрат на експлуатацію та енергоспоживання.

Практичні кейси, зокрема впровадження в Японії та на провідних світових платформах, підтверджують здатність цих підходів задовольняти зростаючі вимоги бізнесу та користувачів.

Оптимізація виконує ключову роль для підвищення продуктивності та ефективності хмарних сервісів, оскільки вона дозволяє більш раціонально використовувати ресурси, швидше реагувати на динамічні зміни навантаження та забезпечувати високу надійність обслуговування. Інтеграція штучного інтелекту та автоматизованих методів прогнозування дає змогу будувати самонавчальні системи, що адаптуються до складних умов експлуатації, що є важливою вимогою сучасних цифрових екосистем.

Перспективними напрямками подальших до-

сліджень є розвиток гібридних моделей оптимізації, які поєднують традиційні алгоритми з інтелектуальними підходами, розробка енергоефективних методів управління ресурсами та забезпечення безпеки у багатокористувацьких хмарних середовищах.

Важливим є також розширення використання квантових обчислень та периферійних (edge) технологій для підвищення масштабованості та мінімізації затримок. Подальша увага має бути приділена створенню стандартизованих фреймворків для адаптивної оптимізації, що забезпечують інтеграцію різномірних хмарних сервісів у комплексні екосистеми.

Таким чином, оптимізація хмарних обчислень залишається актуальною і динамічною сферою, розвиток якої визначає ефективність сучасних інформаційних технологій та становить основу цифрової трансформації підприємств і державних структур у 2026 році і надалі.

Конфлікт інтересів

Автори декларують, що не мають конфлікту інтересів стосовно даного дослідження, в тому числі фінансового, особистісного характеру, авторства чи іншого характеру, що міг би вплинути на дослідження та його результати, представлені в даній статті.

Використання засобів штучного інтелекту

Автори підтверджують, що не використовували технології штучного інтелекту при створенні представленої роботи.

СПИСОК ЛІТЕРАТУРИ

1. Deochake S. Cloud Cost Optimization: A Comprehensive Review of Strategies and Case Studies / S. Deochake // arXiv. 2023. URL: <https://arxiv.org/abs/2307.12479>
2. Sundaraperumal P. Cloud Profiling Techniques and Optimization Strategies for Cloud Computing. / P. Sundaraperumal, P. Kumar, A. Prabhakar, S.P. Chakravarthy // AIP Conf. Proc., 2025. 3279 (1). DOI: <https://doi.org/10.1063/5.0261982>
3. Поперешняк С. В. Хмарні технології як сервіси для оптимізації процесів адміністрування. / С. В. Поперешняк, А. С. Вечерковська, М. Ю. Хільченко, А. В. Антоненко // Таврійський науковий вісник. Серія: Технічні науки. – 2024. – № 6. – С. 54-63. doi: <https://doi.org/10.32782/tnv-tech.2023.6.7>
4. Sarkar S. An Investigation into the Performance Optimization of Cloud Computing Systems using Machine Learning Algorithms. SSRN. 2025. URL: https://papers.ssrn.com/sol3/papers.cfm?abstract_id=5317785 (дата звернення: 19.09.2025).
5. Nawrocki P., Smendowski M. Optimization of the Use of Cloud Computing Resources Using Exploratory Data Analysis and Machine Learning. Sciendo. 2024. URL: <https://sciendo.com/article/10.2478/jaiscr-2024-0016>
6. Song J. Improving Resource Efficiency in Cloud Computing [PhD dissertation]. Washington University in St. Louis, 2021. URL: https://openscholarship.wustl.edu/cgi/viewcontent.cgi?article=1788&context=eng_etds
7. Study on cloud resource optimization in Japan // AJCO. 2024. URL: <https://ajpojournals.org/journals/index.php/AJCE/article/view/2249>
8. A Review of Cloud Computing CPU Resource Optimization // ACM Digital Library. 2025. URL: <https://dl.acm.org/doi/10.1145/3745812.3745839>
9. Koliechkina L., Pichugina O., Dvirna O. Horizontal Method Application to Multiobjective Combinatorial Optimization over Permutations. (2022) 2022 IEEE 3rd International Conference on System Analysis and Intelligent Computing, SAIC 2022 - Proceedings. DOI: <https://doi.org/10.1109/SAIC57818.2022.9923018>
10. Koliechkina L., Dvirna O. Using Models of Combinatorial Optimization Problem to Estimate the Parameters of an Intelligent System. (2024) CEUR Workshop Proceedings, 3777, pp. 385 – 391. <https://www.scopus.com/inward/record.uri?eid=2-s2.0-85210084496&partnerID=40&md5=33c37ffcfbd74e2058973befcbdacd72> дата звернення: 19.09.2025).
11. Koliechkina L.N., Dvernaya O.A., Nagornaya A.N. Modified coordinate method to solve multicriteria optimization problems on combinatorial configurations. (2014) Cybernetics and Systems Analysis, 50 (4), pp. 620 - 626. DOI: <https://doi.org/10.1007/s10559-014-9650-4>
12. Волк М.О. Оптимізація ресурсів у хмарних обчисленнях: гібридний підхід до автоматизації операцій та енергозбереження. / М.О. Волк, А.М. Бугрій, С.І. Ковтун, Р.М. Брестовицький, Я.В. Лобач // Вчені записки ТНУ імені В.І. Вернадського. Серія: Технічні науки. 2024. Том 35 (74) № 5. С. 91–96. DOI <https://doi.org/10.32782/2663-5941/2024.5.1/15>
13. Yang Y. An Optimization Method for Reliable Cloud Service Composition // Atlantis Press. 2018. URL: <https://www.atlantispress.com/article/25888651.pdf>
14. Zeng R. Performance optimization for cloud computing systems in ... // Frontiers of Computer Science. 2022. URL: <https://journal.hep.com.cn/fcs/EN/10.1007/s11704-020-0072-3>

15. Pachipala Y. Optimizing Task Scheduling in Cloud Computing // ScienceDirect. 2024. URL: <https://www.sciencedirect.com/science/article/pii/S1877050924006094>
16. Aral A. Modeling and Optimization of Resource Allocation in Cloud Computing [PhD thesis]. 2014. URL: <https://www.slideshare.net/AtakanAral/proposal-37035819>
17. Singh A. Optimization of the Cloud-Native Infrastructure using AI/ML Techniques. 2023. URL: <https://www.diva-portal.org/smash/get/diva2:1830454/FULLTEXT01.pdf>
18. Cloud Optimization: 2025 Guide to Process, Tools & Best Practices // UmbrellaCost. 2025. URL: <https://umbrellacost.com/learning-center/cloud-optimization-why-its-important-6-critical-best-practices>

Received (Надійшла) 16.10.2025

Accepted for publication (Прийнята до друку) 21.01.2026

Publication date (Дата публікації) 27.02.2026

ВІДОМОСТІ ПРО АВТОРІВ / ABOUT THE AUTHORS

Двірна Олена Анатоліївна – кандидат фізико-математичних наук, завідувач кафедри комп'ютерних та інформаційних технологій і систем, Національний університет «Полтавська політехніка імені Юрія Кондратюка», Полтава, Україна;
Olena Dvirna – Candidate of Physical and Mathematical Sciences, Head of the Department of Computer and Information Technologies and Systems, National University “Yuri Kondratyuk Poltava Polytechnic”, Poltava, Ukraine;
e-mail: lenadvirna@gmail.com; ORCID Author ID: <https://orcid.org/0000-0002-0750-6958>;
Scopus Author ID: <https://www.scopus.com/authid/detail.uri?authorId=57195305100>.

Набока Сергій Володимирович – аспірант, Національний університет «Полтавська політехніка імені Юрія Кондратюка», Полтава, Україна;

Sergiy Naboka – PhD Student, National University “Yuri Kondratyuk Poltava Polytechnic”, Poltava, Ukraine;
e-mail: poltava-gambit@ukr.net; ORCID Author ID <https://orcid.org/0009-0007-4752-3705>.

Optimization of cloud services performance: methods and their efficiency

Olena Dvirna, Sergiy Naboka

Abstract. Relevance. Optimizing the performance of cloud services becomes critically important in 2025–2026 due to the exponential growth of data volumes processed within digital ecosystems and the need for rapid adaptation to dynamic workloads driven by the adoption of artificial intelligence (AI) and IoT devices. According to forecasts, the global cloud computing market is growing by 20–25% annually, and in Ukraine this trend is reinforced by government initiatives such as the National AI Development Strategy through 2025, which identifies cloud services as key infrastructure. However, frequent incidents related to overloads, costs, and security—especially under conditions of energy shortages and geopolitical risks—indicate the fragmented nature of current optimization practices. In Ukrainian business and public administration, where cloud technologies are deployed for digitalization (e.g., e-government and fintech systems), inefficient resource management leads to losses of up to 30–40% of IT budgets. This makes the topic strategically important for ensuring competitiveness, sustainable development, and resilience to cyber threats in multi-cloud environments. **The object of the study** is modern cloud services, including IaaS, PaaS, and SaaS models provided by platforms such as AWS, Google Cloud, Azure, and local Ukrainian providers, with a focus on performance under variable workloads. Particular attention is paid to the Ukrainian market, where hybrid cloud solutions are combined with local infrastructure to overcome limitations in internet bandwidth and regulatory barriers. The research covers key aspects such as dynamic resource allocation, load balancing, automatic scaling, and AI integration for demand forecasting. **The aim** is to systematize optimization methods, identify challenges related to scalability, security, and energy efficiency, and develop recommendations for implementation in business (fintech, e-commerce) and the public sector (e-governance, critical infrastructure), taking into account Ukraine's specific context. **Methodology.** The article applies a comprehensive approach that includes a comparative analysis of modern optimization methods based on reanalysis data from cloud providers (AWS Auto Scaling, Google Cloud Autoscaler, Azure Autoscale) and implementation case studies from 2025. An expert risk assessment using a factor model is employed to classify causes of inefficiency (human factors, overloads, weak encryption), along with mathematical optimization modeling using SLO/SLA metrics and linear programming to balance costs and performance. Additionally, the study analyzes international standards (ISO 27001 for security, Green Cloud for energy efficiency) and regional Ukrainian data, including numerical experiments with AI-based forecasting algorithms (machine learning for autoscaling). Quantitative assessments are based on statistics such as a 30% cost reduction through AI automation, monitoring via Google Cloud Monitoring, and data center energy consumption modeling. **Results.** The study confirms that key optimization methods include dynamic resource allocation with AI-based forecasting, which reduces over-provisioning by 25–40%, and automatic scaling that responds to peak loads within seconds, as demonstrated by the Pinterest case on AWS (30% cost reduction). Load balancing and serverless architectures increase performance by 38–60% by eliminating downtime. Challenges include security issues (misconfigured IAM, blind spots in ephemeral resources), energy efficiency (a 40% increase in consumption due to AI), and scaling constraints in Ukraine (limited infrastructure, shadow IT). Provider comparisons show the advantages of hybrid models: Kubernetes orchestration in multi-cloud environments reduces latency by 40%, while neural network quantization cuts computational costs by 40%. In the Ukrainian context, uneven adoption is observed: businesses save on administrative staff, while the public sector suffers from energy shortages. **Conclusions.** A comprehensive optimization approach is proposed, including AI integration for predictive scaling, standardization of BRM-like procedures for cloud environments (adaptive risk management protocols), and harmonization with local data centers and renewable energy sources to reduce CO₂ emissions by 30%. **Recommendations for Ukraine** include implementing autoscaling down during off-hours, conducting security testing (CSPM, SIEM), monitoring SLO/SLA metrics, and providing training for AI-driven resource management. The feasibility of hybrid solutions for the sustainable development of the digital ecosystem is demonstrated, ensuring reliability, cost efficiency, and competitiveness in business and public administration, with the potential to reduce costs by 30–50% and increase productivity. This approach lays the foundation for workload forecasting and risk minimization in 2026 and beyond.

Keywords: dynamic resource allocation, load balancing, auto-scaling, machine learning, multi-criteria optimization, energy efficiency, scalability, reliability.