

**НАЦІОНАЛЬНИЙ ТЕХНІЧНИЙ УНІВЕРСИТЕТ УКРАЇНИ  
«КИЇВСЬКИЙ ПОЛІТЕХНІЧНИЙ ІНСТИТУТ  
імені ІГОРЯ СІКОРСЬКОГО»**

**Навчально-науковий інститут телекомунікаційних систем  
Кафедра телекомунікацій**

До захисту допущено:

Завідувач кафедри

\_\_\_\_\_ Сергій КРАВЧУК

« \_\_\_\_ » \_\_\_\_\_ 2025 р.

**Дипломна робота**

**на здобуття ступеня бакалавра**

**за освітньо-професійною програмою «Інженерія та програмування  
інфокомунікацій»**

**спеціальності 172 «Телекомунікації та радіотехніка»**

**на тему: «Використання графового методу для пошуку оптимального  
шляху в телекомунікаційних системах соціальних платформ за допомогою  
LLM»**

Виконав:

студент IV курсу, групи ТЗ-12

Зозуля Владислав Сергійович \_\_\_\_\_

Керівник:

Старший викладач кафедри ТК НН ІТС,

Чуб Михайло Миколайович \_\_\_\_\_

Консультант з Практична реалізація та оцінка ефективності:

Старший викладач кафедри ТК НН ІТС, к.т.н.

Маньківський Володимир Броніславович \_\_\_\_\_

Рецензент:

Доцент кафедри ІТТ НН ІТС КПІ, д.т.н.

Астраханцев Андрій Анатолійович \_\_\_\_\_

Засвідчую, що у цій дипломній роботі  
немає запозичень з праць інших авторів  
без відповідних посилань.

Студент \_\_\_\_\_

Київ 2025

**Національний технічний університет України**  
**«Київський політехнічний інститут імені Ігоря Сікорського»**  
**Навчально-науковий інститут телекомунікаційних систем**  
**Кафедра телекомунікацій**

Рівень вищої освіти – перший (бакалаврський)

Спеціальність – 172 «Телекомунікації та радіотехніка»

Освітньо-професійна програма «Інженерія та програмування інфокомунікацій»

ЗАТВЕРДЖУЮ

Завідувач кафедри

\_\_\_\_\_ Сергій КРАВЧУК

«\_\_» \_\_\_\_\_ 2025 р.

### ЗАВДАННЯ

**на дипломну роботу студенту**

**Зозулі Владиславу Сергійовичу**

1. Тема роботи «Використання графового методу для пошуку оптимального шляху в телекомунікаційних системах соціальних платформ за допомогою LLM», керівник роботи Чуб Михайло Миколайович, затверджені наказом по університету від «26» травня 2025 р. № 1755-с.
2. Термін подання студентом роботи 9 червня 2025 р.
3. Вихідні дані до роботи: Теоретичні матеріали з теорії графів і класичних алгоритмів пошуку найкоротшого шляху. Праці та статті з використання GNN у задачах маршрутизації та мережевого аналізу. Огляд досліджень із застосування LLM для аналізу мережевих логів і прогнозування навантажень. Документація платформи (OpenAI API) щодо інтеграції LLM для семантичного аналізу тексту. Опис архітектури телекомунікаційних систем соціальних платформ (Facebook, Twitter, Instagram) з акцентом на CDN і розподілену маршрутизацію. Інструкції зі встановлення та використання бібліотек NetworkX, PyTorch Geometric для побудови та обчислення графів. Приклади симуляційних набір даних для тестування графових алгоритмів у мережевому середовищі
4. Зміст роботи:
  - Дослідити математичні основи графових алгоритмів маршрутизації (Dijkstra, A\*) та принципи побудови зважених графів для телекомунікаційних мереж.
  - Ознайомитися з існуючими інструментами та бібліотеками для обробки

графів і пошуку найкоротших шляхів (NetworkX, igragh тощо).

- Проаналізувати можливості використання GNN для прогнозування завантаження каналів і виявлення «вузьких місць».
- Вивчити, як LLM можуть аналізувати текстові логи мережевих пристроїв і повідомлення користувачів для семантичного передбачення пікових навантажень.
- Спроектувати архітектуру рішення: компоненти збору метрик, формування графа, виклик GNN/LLM та модуль пошуку маршруту.
- Реалізувати прототип: створити симульовану мережу, інтегрувати GNN для інференсу та звернення до LLM API.
- Розробити інтерфейс користувача для введення параметрів, запуску обчислень і візуалізації результатів (графів і маршруту).
- Провести тестування системи: перевірити коректність графових обчислень, точність прогнозів GNN/LLM, провести навантажувальні тести.
- Оцінити продуктивність та масштабованість рішення: порівняти час пошуку з класичними алгоритмами, проаналізувати вплив семантичного аналізу на швидкість і якість маршрутизації.
- Підготувати технічну документацію для встановлення, конфігурації та експлуатації розробленого модуля.
- Підготуватися до захисту дипломної роботи: створити презентацію з основними результатами й висновками.

5. Перелік ілюстративного матеріалу:

6. Консультанти розділів роботи □

| Розділ                                      | Прізвище, ініціали та посада консультанта                    | Підпис, дата   |                  |
|---|--|----------------|------------------|
|   |  | Завдання видав | Завдання прийняв |
| Практична реалізація та оцінка ефективності | Старший викладач кафедри ТК НН ІТС, к.т.н. Маньківський В.Б. |                |                  |

7. Дата видачі завдання «5» жовтня 2024 р.

#### Календарний план

| № з/п | Назва етапів виконання дипломної роботи  | Термін виконання етапів роботи | Примітка |
|-------|--|--------------------------------|----------|
| 1     | Створення, оформлення та затвердження технічного завдання з визначенням мети, завдань і меж дослідження: графова маршрутизація в соціальних мережах із залученням LLM. | 05.10.2024 - 25.12.2024        | Виконано |
| 2     | Аналіз науково-технічної літератури з теорії графів, алгоритмів пошуку найкоротшого шляху (Dijkstra, A*) та  | 05.01.2025 - 25.01.2025        | Виконано |

|   |  |                         |          |
|---|--|-------------------------|----------|
|   | сучасних підходів на базі GNN і LLM у мережевому моніторингу.  |                         |          |
| 3 | Вибір технологій і інструментів для реалізації: Python, NetworkX, PyTorch Geometric, OpenAI API (LLM), FastAPI.  | 26.01.2025 – 31.01.2025 | Виконано |
| 4 | Проектування архітектури рішення: компоненти збору метрик, формування зваженого графа, модуль GNN, виклики до LLM і механізм пошуку маршруту                                     | 01.02.2025 - 20.02.2025 | Виконано |
| 5 | Реалізація серверної частини: налаштування середовища, збір даних про затримки і пропускну здатність, побудова графової моделі, інтеграція GNN та LLM API.                       | 21.02.2025 - 24.03.2025 | Виконано |
| 6 | Розробка клієнтської частини: створення веб-інтерфейсу для подачі параметрів мережі, відображення графа, візуалізації оптимального шляху та рекомендацій LLM.                    | 26.03.2025 - 28.04.2025 | Виконано |
| 7 | Проведення тестування і налагодження: перевірка коректності побудови графа й обчислень маршруту, оцінка якості прогнозів GNN і семантичних висновків LLM, навантажувальні тести. | 29.04.2025 - 11.05.2025 | Виконано |
| 8 | Оформлення технічної документації: опис API, приклади роботи системи, а також звіт про виконані експерименти і результати.   | 12.05.2025 - 25.05.2025 | Виконано |
| 9 | Підготовка до захисту дипломної роботи: створення презентації, підготовка доповіді.  | 26.05.2025 - 05.06.2025 | Виконано |

Студент

Владислав ЗОЗУЛЯ

Керівник

Михайло ЧУБ

## РЕФЕРАТ

Пояснювальна записка: 77 сторінок, 16 рисунків, 19 джерел.

Мета роботи полягає у розробки системи для оптимізації маршрутизації в телекомунікаційних мережах соціальних платформ із використанням графових методів та великих мовних моделей.

У роботі розглянуто основні поняття мережевого графа, класичні алгоритми пошуку шляху (Dijkstra, A\*) та сучасні підходи на основі Graph Neural Networks. Досліджено можливості інтеграції LLM для семантичного аналізу мережевих лог-файлів і прогнозування пікових навантажень.

Для реалізації використано бібліотеки Python (NetworkX, PyTorch Geometric), FastAPI для бекенду. Реалізовано підсистему побудови та динамічного оновлення зваженого графа, модуль GNN-інференсу, виклики до OpenAI API і клієнтський інтерфейс із візуалізацією оптимального маршруту та рекомендацій.

Розроблена система дозволяє автоматично обирати найбільш ефективний шлях з урахуванням поточних метрик мережі й семантичних підказок від LLM, значно знижуючи латентність і ризик перевантаження каналів.

Ключові слова: ГРАФОВА МАРШРУТИЗАЦІЯ, DIJKSTRA, A\*, GNN, LLM, СОЦІАЛЬНА ПЛАТФОРМА.

## ABSTRACT

Explanatory note: 77 pages, 16 figures, 19 sources.

The purpose of this work is to develop a system for optimizing routing in telecommunications networks of social platforms using graph methods and large language models.

The work covers the basic concepts of network graphs, classic pathfinding algorithms (Dijkstra, A\*) and modern approaches based on Graph Neural Networks. The possibilities of integrating LLM for semantic analysis of network log files and peak load forecasting are explored.

Python libraries (NetworkX, PyTorch Geometric), FastAPI for the backend. A subsystem for building and dynamically updating a weighted graph, a GNN inference module, calls to the OpenAI API, and a client interface with visualization of the optimal route and recommendations were implemented.

The developed system allows you to automatically select the most efficient path, taking into account current network metrics and semantic hints from LLM, significantly reducing latency and the risk of channel overload.

Keywords: GRAPH ROUTING, DIJKSTRA, A\*, GNN, LLM, SOCIAL PLATFORM.

## ЗМІСТ

|  |    |
|--|----|
| ВСТУП.....   | 10 |
| РОЗДІЛ 1 АНАЛІЗ ПРЕДМЕТНОЇ ОБЛАСТІ .....   | 12 |
| 1.1. Особливості телекомунікаційних систем соціальних платформ.....                                    | 12 |
| 1.2. Виявлення вузьких місць у маршрутизації.....  | 15 |
| 1.3. Вплив затримок на якість обслуговування користувачів.....   | 15 |
| 1.4. Використання CDN для оптимізації передачі даних .....   | 16 |
| 1.5. Методи пошуку оптимального шляху в мережах .....  | 17 |
| 1.5.1. Алгоритм Dijkstra .....   | 17 |
| 1.5.2. Алгоритм A* .....   | 18 |
| 1.6. Графові нейронні мережі.....  | 20 |
| 1.7. Недоліки традиційних методів.....   | 22 |
| Висновки .....   | 23 |
| РОЗДІЛ 2 ЗАПРОПОНОВАНИЙ МЕТОД ОПТИМІЗАЦІЇ МАРШРУТИЗАЦІЇ ...  | 24 |
| 2.1. Графове уявлення мережі .....   | 24 |
| 2.1.1. Вузли та ребра: параметри каналів.....  | 25 |
| 2.1.2. Ідентифікація критичних вузлів і каналів.....   | 26 |
| 2.2. Пошук оптимального маршруту в реальному часі .....  | 28 |
| 2.3. Інтеграція LLM .....  | 30 |
| 2.3.1. Використання LLM для аналізу історичних даних трафіку та прогнозування пікових навантажень..... | 30 |
| 2.3.2. Визначення альтернативних маршрутів на основі семантичного аналізу ..                           | 32 |
| 2.4. Архітектура рішення.....  | 33 |
| 2.4.1. Компоненти та їх взаємодія.....   | 33 |
| 2.4.2. Розподілена база даних і моніторинг .....   | 35 |
| Висновки .....   | 38 |
| РОЗДІЛ 3 ПРАКТИЧНА РЕАЛІЗАЦІЯ ТА ОЦІНКА ЕФЕКТИВНОСТІ .....   | 39 |
| 3.1. Формування графової моделі мережі .....   | 39 |
| 3.2. Реалізація алгоритмів маршрутизації.....  | 41 |
| 3.2.1. Код для алгоритму Dijkstra.....   | 42 |
| 3.2.2. Код для алгоритму A* .....  | 45 |
| 3.2.3. Навчання та застосування GNN.....   | 48 |
| 3.3. Інтеграція LLM для прогнозування навантажень .....  | 54 |
| 3.4. Тестування та аналіз результатів .....  | 57 |

|  |    |
|--|----|
| 3.4.1. Методологія тестування (сценарії, метрики).....   | 57 |
| 3.4.2. Результати порівняння (Dijkstra, GNN).....        | 58 |
| Висновки .....   | 63 |
| РОЗДІЛ 4 ВИСНОВКИ ТА РЕКОМЕНДАЦІЇ .....                  | 65 |
| 4.1. Основні результати дослідження.....                 | 65 |
| 4.2. Порівняльний аналіз методів .....                   | 66 |
| 4.3. Практичні рекомендації для впровадження.....        | 69 |
| 4.4. Напрями подальших досліджень та удосконалення ..... | 72 |
| Висновки .....   | 73 |
| ЗАГАЛЬНІ ВИСНОВКИ ПО РОБОТІ.....                         | 74 |
| СПИСОК ВИКОРИСТАНИХ ДЖЕРЕЛ .....                         | 76 |

**ПЕРЕЛІК СКОРОЧЕНЬ**

|     |                          |                          |
|-----|--------------------------|--------------------------|
| CDN | Content Delivery Network | Мережа доставки контенту |
| LLM | Large Language Models    | Великі мовні моделі      |
| GNN | Graph Neural Networks    | Графові нейронні мережі  |

## ВСТУП

У сучасному світі соціальні платформи стали невід'ємною частиною повсякденного життя мільярдів людей. Такі сервіси, як Facebook, Instagram, Twitter, TikTok, Telegram та інші, забезпечують постійний обмін повідомленнями, мультимедійними файлами, відео- та аудіоконтентом. Вся ця комунікація здійснюється за допомогою складних телекомунікаційних систем, що складаються з розподілених серверів, маршрутизаторів, центрів обробки даних та численних логічних і фізичних мережевих з'єднань. Надійність, швидкість та ефективність передачі інформації в таких системах мають вирішальне значення для задоволення потреб користувачів.

З розвитком соціальних платформ обсяг оброблюваних і переданих даних зростає в геометричній прогресії. Зі збільшенням кількості активних користувачів, особливо в години пік, телекомунікаційні системи піддаються підвищеному навантаженню. Це призводить до збільшення затримок, перевантаження каналів зв'язку і, в деяких випадках, до тимчасових перебоїв у роботі сервісів. Традиційні методи маршрутизації, що використовуються для визначення шляхів передачі даних у таких мережах, не завжди забезпечують достатню продуктивність в умовах високодинамічних змін трафіку.

У зв'язку з цим існує потреба в інтелектуальних підходах до оптимізації маршрутизації мережі. Багатообіцяючим підходом є використання графічних методів, за допомогою яких телекомунікаційна мережа може бути представлена як сукупність вузлів (серверів, пристроїв, користувачів) і ребер (каналів зв'язку) і можна визначити оптимальні шляхи між ними з урахуванням певних критеріїв, таких як швидкість, пропускна здатність, затримка, надійність тощо.

Сучасний розвиток технологій штучного інтелекту, зокрема великих мовних моделей, відкриває нові можливості для аналізу мережевих даних і прийняття рішень у режимі реального часу. LLM можуть обробляти складні запити, прогнозувати майбутні навантаження, генерувати аналітичні висновки на основі великих обсягів даних і брати участь у виборі оптимального шляху

передачі, враховуючи не тільки технічні характеристики мережі, але й контекстуальні або поведінкові аспекти користувачів.

Таким чином, інтеграція графічних методів з LLM створює потужний інструмент для підвищення ефективності телекомунікаційних систем соціальних платформ. Це не тільки сприяє зменшенню затримок і поліпшенню якості обслуговування користувачів, але й створює основу для побудови адаптивних, самонавчальних мереж, здатних реагувати на зміни в режимі роботи.

### **Мета дослідження**

Основна ідея цієї роботи полягає в тому, щоб знайти спосіб зробити маршрутизацію в телеком-мережах соціальних платформ більш ефективною завдяки поєднанню графових алгоритмів із можливостями великих мовних моделей. Спочатку ми детально вивчаємо, як побудовані ці мережі, і визначаємо «вузькі місця» - ті ділянки, де найчастіше збираються затори чи виникають перевантаження. Далі створюємо графове представлення інфраструктури, у якому кожен канал зв'язку отримує вагу, що відображає його пропускну здатність, затримку та інші ключові характеристики. Уже над цим графом запускаємо перевірені алгоритми Dijkstra чи A\*, щоб знайти найкоротший або оптимальний шлях для передавання даних.

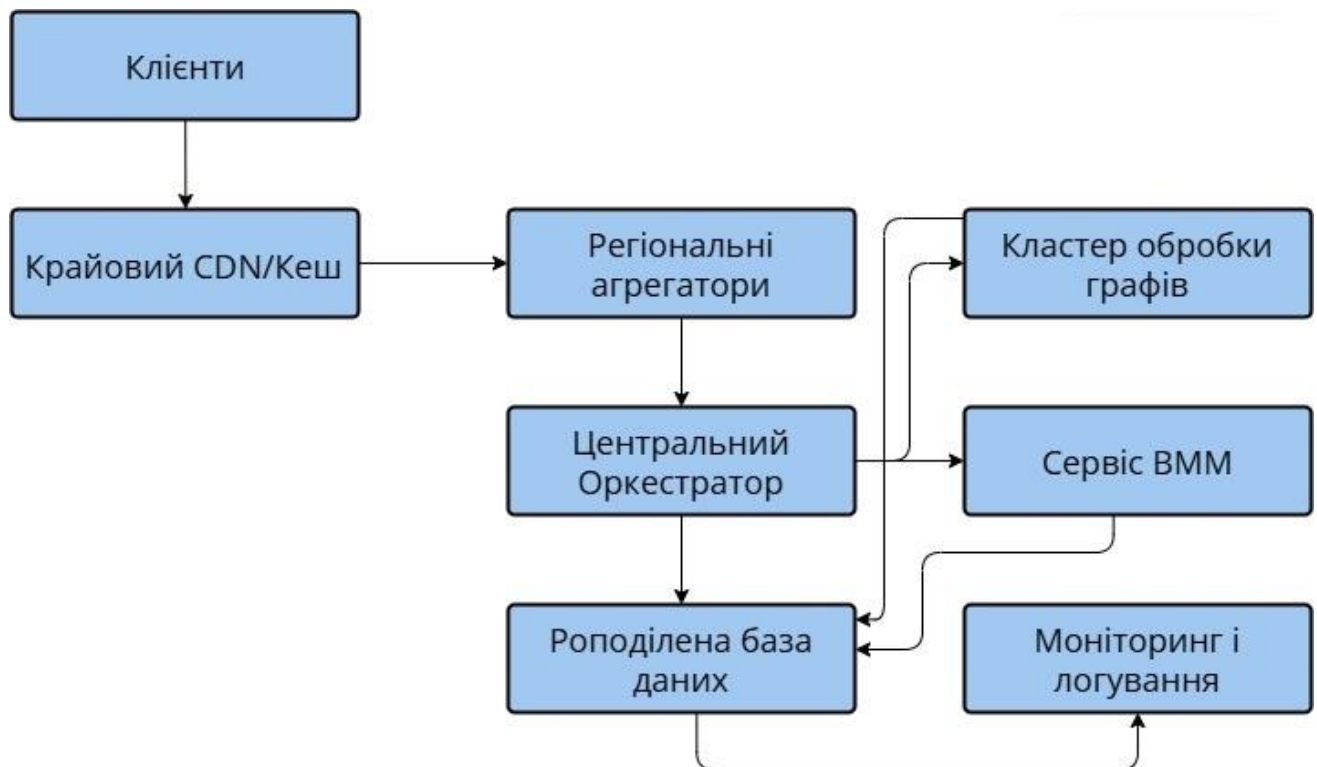
Щоб система працювала швидше і могла підлаштовуватися під змінні умови, ми додаємо шар із LLM-моделлю. Вона аналізує історичні дані про трафік, виявляє закономірності у поведінці користувачів і навіть прогнозує, коли трафік піде вгору, щоб заздалегідь підкоригувати маршрути. Після цього етапу йде розробка самого програмного рішення або прототипу, який реалізує описані підходи. У фіналі ми тестуємо його на реальних чи синтетичних даних, порівнюючи отримані результати з класичними методами. Завдяки цьому вдається побачити, наскільки наш підхід поліпшує швидкість доставки контенту та зменшує затримки. Важливо, що результати мають практичне значення: їх можна використовувати в компаніях, які створюють або підтримують соціальні платформи, а також у проєктах із розподілених обчислень, хмарних сервісів та інтелектуальних телекомунікацій.

## РОЗДІЛ 1

### АНАЛІЗ ПРЕДМЕТНОЇ ОБЛАСТІ

#### 1.1. Особливості телекомунікаційних систем соціальних платформ

Сучасні соціальні платформи функціонують на основі складної розподіленої архітектури, яка дозволяє обробляти гігантські обсяги трафіку, забезпечуючи при цьому мінімальні затримки та високу доступність контенту для мільйонів користувачів у реальному часі. Основу цієї архітектури складають дата-центри (ЦОДи), розміщені у стратегічних географічних регіонах світу. Вони взаємодіють між собою через потужні оптоволоконні канали зв'язку та використовують системи балансування навантаження, CDN і системи глобального кешування [10, 11].



*Рис. 1.1. Розподілена Архітектура Системи Оптимізації Маршрутизації*

На наведеній схемі відображені основні компоненти системи, що забезпечують збір, обробку та прийняття рішень щодо оптимальних маршрутів передачі даних у телекомунікаційних мережах соціальних платформ. Далі подано детальний опис послідовності взаємодії цих компонентів:

### **Клієнти (Clients)**

Кінцеві користувачі або їхні додатки надсилають HTTP/API-запит на отримання контенту (медіа, повідомлень, оновлень тощо).

### **Крайовий CDN/Кеш (Edge CDN/Cache)**

Запит спочатку потрапляє до крайового CDN або локального кешу, який надає збережену копію даних, якщо вона актуальна, щоб зменшити затримку. Якщо контент відсутній або застарілий, запит передається далі.

### **Регіональні Агрегатори (Regional Aggregators)**

Кеші одного регіону надсилають невирішені запити до агрегаційного вузла. Агрегатор перевіряє власний регіональний кеш і, за відсутності потрібних даних, спрямовує запит до центрального оркестратора. Це зменшує навантаження на центральні служби й оптимізує обробку запитів за регіонами.

### **Центральний Оркестратор (Central Orchestrator)**

Після надходження запиту центральний Оркестратор одночасно звертається до двох служб. Кластер обробки графів отримує актуальну топологію мережі разом із характеристиками каналів (пропускна здатність, затримка, завантаження) і виконує алгоритми пошуку найкоротшого чи оптимального шляху (Dijkstra, A\*, GNN), формуючи кілька варіантів маршрутів. Сервіс BMM (LLM Service) аналізує історичні метрики трафіку, прогнозує пікові навантаження та враховує контекст запиту (тип контенту, час доби, геолокацію), щоб надати рекомендації щодо перенаправлення трафіку з урахуванням мінімізації ризиків перевантаження. На основі результатів обох служб Оркестратор обирає остаточний маршрут (наприклад, звернення до певного кешу або первинного сервера).

### **Розподілена База Даних (Distributed Database)**

Усі дані про топологію мережі, історичні метрики навантаження та результати розрахунків маршрутів зберігаються у розподіленій БД. Це забезпечує збереження інформації для подальшого аналізу, корекції маршрутизації й тренування нових моделей (GNN).

## Моніторинг і Логування (Monitoring & Logging)

Після ухвалення рішення результати Оркестратора й метрики стану компонентів (завантаження ЦП, час відповіді, кількість запитів) надходять у модуль моніторингу. Логування фіксує вибрані маршрути, помилки та попередження. Це дає змогу оперативно виявляти збої, аналізувати тенденції та поліпшувати алгоритми.

Наприклад, **Facebook (Meta)** використовує архітектуру типу *Cold Start + Warm Cache*, де нові запити обробляються безпосередньо у дата-центрах, а популярний контент доставляється через кеш-сервери по всьому світу. Це дозволяє скоротити затримку під час масового доступу до вірусного контенту. У рамках своєї інфраструктури Meta також розгорнула власні підводні кабелі для з'єднання дата-центрів між континентами, що зменшує залежність від сторонніх операторів.

**Twitter** перейшов на мікросервісну архітектуру, де кожен тип контенту (твіти, медіа, лайки) обробляється окремим сервісом. Внутрішні API забезпечують маршрутизацію запитів між сервісами, а глобальна CDN мережа оптимізує доставку медіа та відео. При цьому особлива увага приділяється розподілу навантаження, оскільки платформа часто стикається з різкими стрибками трафіку під час глобальних подій.

**TikTok** значно підвищив вимоги до телекомунікаційної інфраструктури через активне використання алгоритмів персоналізації у реальному часі. Відео рекомендується користувачам буквально за частки секунди, що вимагає не лише високої пропускної здатності, але й мінімальної затримки між дата-центрами та кінцевими користувачами. Для цього TikTok використовує багаторівневу CDN-мережу, а також алгоритми динамічного маршрутування запитів.

**Instagram**, як частина екосистеми Meta, інтегрований із глобальною інфраструктурою Facebook, але має специфічні особливості, пов'язані з обробкою медіа. Платформа використовує системи балансування навантаження, які визначають оптимальні маршрути для доставки контенту на основі поточного

стану мережі, а також механізми edge-кешування для прискорення завантаження сторінок і відео.

## **1.2. Виявлення вузьких місць у маршрутизації**

Попри розвинену архітектуру, вузькі місця у маршрутизації залишаються актуальною проблемою. Вони можуть виникати через перевантаження окремих вузлів (серверів CDN або дата-центрів) під час пікових навантажень, обмежену пропускну здатність каналів зв'язку між регіональними ЦОДами, недосконалі алгоритми маршрутизації, які не враховують реальну завантаженість мережі в конкретний момент, а також через збої у системах балансування навантаження, що призводять до перенавантаження певних сегментів. Для усунення цих проблем сучасні платформи застосовують комбінований підхід: спочатку використовують графові алгоритми для пошуку оптимальних маршрутів передачі даних, потім залучають графові нейронні мережі для прогнозування перевантажень, а також інтегрують великі мовні моделі для аналізу поведінки користувачів і динамічного переналаштування маршрутів на основі семантичних патернів. В результаті такі системи забезпечують стабільну роботу соціальних мереж навіть у найгарячіші моменти — під час міжнародних спортивних подій, політичних виборів або запуску вірусного контенту.

## **1.3. Вплив затримок на якість обслуговування користувачів**

Затримка доставки даних є критичною для користувацького досвіду: навіть 50–100 мс можуть помітно знижувати відчуття “плавності” у соцмережах (лайки, коментарі, стріми). У глобальній мережі до цього додаються витрати часу на передачу сигналу через сотні чи тисячі кілометрів, обробку на проміжних маршрутизаторах та черги в пікові години.

Статичні алгоритми маршрутизації (оновлення таблиць раз на декілька секунд чи хвилин) не встигають реагувати на такі швидкі коливання трафіку. Інтелектуальні методи, навпаки, аналізують поточне й історичне навантаження («на льоту»), прогнозують піки й ураховують контекст (наприклад, “вірусні” теми або географічні сплески). Вони дозволяють миттєво перенаправляти трафік

на менш завантажені ланки, знижуючи затримки та підтримуючи добрий рівень сервісу навіть у пікові моменти.

#### **1.4. Використання CDN для оптимізації передачі даних**

Мережа доставки контенту - це розподілена мережа спеціалізованих серверів (edge-нодів), розташованих по всьому світу для максимально швидкої доставки контенту, зберігаючи його копії ближче до кінцевого користувача. Соціальні платформи активно використовують CDN для кешування найрізноманітніших ресурсів: від статичних файлів (зображень, відео, CSS (Cascading Style Sheets) та JavaScript) до динамічного вмісту, наприклад, персоналізованих фрагментів сторінок і мініатюр відео. Основною перевагою CDN є суттєве зниження затримки : edge-ноди розташовані у великих центрах обробки трафіку, тому дані доставляються за мілісекунди, а не сотні, якби їх надсилали безпосередньо з центрального дата-центру. Завдяки цьому користувачі отримують контент без відчутних пауз навіть у пікові години активності, коли одночасно обробляються тисячі чи мільйони запитів. Крім того, CDN дозволяє балансувати навантаження між різними фізичними вузлами, розвантажуючи центральні сервери та запобігаючи їхньому перевантаженню. Сучасні рішення пропонують інтегровані механізми захисту від DDoS-атак і бот-трафіку, що підвищує стійкість платформи та гарантує високу доступність сервісу. Великі гравці ринку, як-от Cloudflare і Akamai, мають тисячі edge-нодів, які автоматично оптимізують файли, наприклад, перетворюють зображення у сучасні формати WebP (WEB Pictures) чи AVIF (AV1 Image File Format) і підтримують протоколи HTTP/2 та HTTP/3 для одночасної передачі великої кількості об'єктів. Водночас такі компанії, як Meta (Facebook, Instagram) чи Google (YouTube), розробили власні CDN-інфраструктури, адаптовані під специфіку їхніх продуктів. Їхні рішення автоматично змінюють якість відео залежно від пропускної здатності клієнта та розподіляють контент між кількома дата-центрами з урахуванням географічного розташування, щоб забезпечити доступність навіть у разі виходу з ладу окремих вузлів. Завдяки впровадженню CDN соціальні платформи забезпечують швидкий доступ до контенту незалежно

від місцеперебування користувача, суттєво скорочують навантаження на основні сервери й підвищують загальну масштабованість і стабільність своєї інфраструктури [12, 13].

### **1.5. Методи пошуку оптимального шляху в мережах**

Класичні алгоритми, такі як Dijkstra та  $A^*$ , широко використовуються для пошуку найкоротшого шляху в графах, що робить їх невід'ємною частиною маршрутизації в телекомунікаційних мережах. У цьому контексті вузли графа відповідають маршрутизаторам, серверам або дата-центрам, а ребра — фізичним чи віртуальним каналам зв'язку, кожен з яких має задану вагу (наприклад, час затримки, пропускну здатність або вартість використання). Алгоритм Dijkstra ефективно знаходить мінімальний шлях у зваженому графі, гарантувавши оптимальне рішення, якщо всі ваги невід'ємні. Завдяки цьому Dijkstra став базовим інструментом у багатьох мережевих протоколах маршрутизації, забезпечуючи точність і надійність у підборі маршрутів для пакетів даних. З іншого боку, алгоритм  $A^*$  доповнює класичну ідею, вводячи евристичну оцінку відстані до цілі: це дозволяє скоротити число можливих шляхів, які потрібно дослідити, і віднайти потрібний маршрут швидше, за умови що евристика добре наближує реальну вартість. Умовний приклад — якщо в мережі відомі приблизні відстані між регіональними дата-центрами,  $A^*$  може сфокусуватися на тих сегментах, які найімовірніше приведуть до оптимального результату, пропускаючи малоімовірні варіанти. Обидва алгоритми мають свої обмеження: Dijkstra може бути повільним на великих мережах із сотнями тисяч вузлів, а  $A^*$  вимагає ретельно підібраної евристики, щоб не втратити у точності. Проте в умовах стабільної топології та достатнього обсягу обчислювальних ресурсів вони залишаються найпоширенішими методами для побудови таблиць маршрутизації, адже забезпечують баланс між точністю результату та швидкістю прийняття рішення.

#### **1.5.1. Алгоритм Dijkstra**

Алгоритм Dijkstra знаходить найкоротший шлях від одного джерела до всіх інших вузлів у зв'язаному неорієнтованому чи орієнтованому графі з

невід'ємними вагами ребер [1]. Основна ідея полягає у поступовому розширенні множини «оброблених» вузлів: спочатку обирається вузол із мінімальною відстанню від початку, потім оновлюються відстані до його сусідів, і цей процес повторюється, доки не опрацьовано всі вузли чи не досягнуто заданої цілі. Завдяки використанню пріоритетної черги складність алгоритму становить  $O((V + E) \log V)$ ,

де  $V$  – кількість вузлів,

$E$  – кількість ребер.

Dijkstra-гарантує оптимальний результат, але в дуже великих мережах може вимагати значних обчислювальних ресурсів при масивних  $V$ .

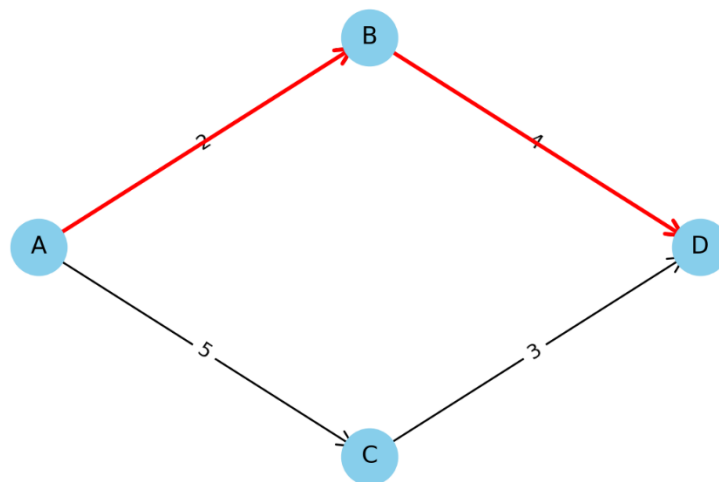


Рис. 1.1. Алгоритм Dijkstra

Ребра та їх ваги:

- $A \rightarrow B$ : 2
- $A \rightarrow C$ : 5
- $B \rightarrow D$ : 4
- $C \rightarrow D$ : 3

Виділений червоним ланцюжок  $A \rightarrow B \rightarrow D$  ( $2 + 4 = 6$ ) є найкоротшим, оскільки шлях  $A \rightarrow C \rightarrow D$  дає суму  $5 + 3 = 8$ .

### 1.5.2. Алгоритм $A^*$

Алгоритм  $A^*$  комбінує принципи пошуку в ширину з евристичним оцінюванням «відстані до мети». Для кожного розглянутого вузла він обчислює

суму  $f(n) = g(n) + h(n)$ , де  $g(n)$  - фактична відстань від початкового вузла до  $n$ , а  $h(n)$  - евристична оцінка відстані від  $n$  до цільового вузла [2]. Якщо  $h(n)$  завжди не перевищує реальної вартості (тобто є допустимою),  $A^*$  гарантовано знайде оптимальний шлях, але значно скоротить кількість перевічених вузлів порівняно з Dijkstra. Складність залежить від якості евристики: при «ідеальному» наближенні ( $h(n) \approx$  реальна відстань) час роботи може бути майже лінійним у кількості вузлів. У телекомунікаційних мережах для  $h(n)$  часто використовують прямі відстані (геодані) або приблизну кількість сегментів, що залишилися до цілі.

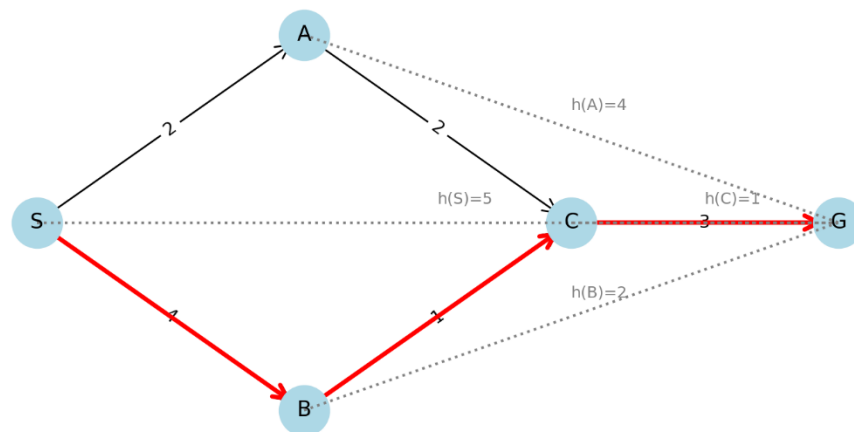


Рис. 1.2. Алгоритм  $A^*$

На наведеній схемі показано приклад графа для алгоритму  $A^*$ , де:

- Вузол S (старт) має двох сусідів: A (вага ребра 2) і B (вага ребра 4).
- Вузол A веде до C з вагою 2.
- Вузол B веде до C з вагою 1.
- Вузол C веде до G (ціль) з вагою 3.
- Сірим пунктирними лініями позначено евристику  $h(n)$  - оцінку відстані від кожного вузла до цілі G (для S: 5, A: 4, B: 2, C: 1, G: 0).

Алгоритм  $A^*$  обирає маршрут із найменшою сукупною функцією  $f(n) = g(n) + h(n)$ , де  $g(n)$  - вартість від початку S до поточного вузла, а  $h(n)$  - евристична оцінка до цілі. У прикладі зображений оптимальний шлях (червоними стрілками):  $S \rightarrow B \rightarrow C \rightarrow G$

- Спочатку з S до B:  $g(B) = 4$ ,  $h(B) = 2$ ,  $f(B) = 6$ .
- Далі  $B \rightarrow C$ :  $g(C) = 4 + 1 = 5$ ,  $h(C) = 1$ ,  $f(C) = 6$ .
- І, нарешті,  $C \rightarrow G$ :  $g(G) = 5 + 3 = 8$ ,  $h(G) = 0$ ,  $f(G) = 8$ .

Таким чином,  $A^*$  обходить вузол A, тому що його оцінка вища:

- Для A з S  $g(A) = 2$ ,  $h(A) = 4$ ,  $f(A) = 6$ . Обидва A і C мали однакове  $f = 6$ , але наступний крок перевагу має C через подальшу мінімальну суму.

## 1.6. Графові нейронні мережі

Графові нейронні мережі (GNN) — це сучасний клас моделей машинного навчання, який працює безпосередньо з графовими структурами. У контексті телекомунікацій GNN дозволяють прогнозувати майбутнє навантаження на мережу, виявляти критичні вузли, а також моделювати поведінку трафіку в реальному часі[3-4].

За допомогою GNN можна передбачити, які частини мережі можуть виявитися перевантаженими у найближчі хвилини, і завчасно скоригувати маршрути, щоб уникнути затримок або втрати даних. Для цього мережа описується у вигляді графа, де вузли (nodes) – це об’єкти, наприклад сервери, користувачі чи інші пристрої, а ребра (edges) – це мережеві канали між ними. Кожен вузол і ребро має власні ознаки (features), такі як поточне навантаження чи пропускна здатність, які служать вхідними даними для GNN. на основі цих фіч нейромережа навчається прогнозувати майбутню «вартість» або ризик перевантаження відповідних ділянок мережі[5, 9, 15].

Кожен вузол отримує інформацію від своїх сусідів. Наприклад:

$$h_v^{(k)} = \text{AGGREGATE}^{(k)}(\{h_u^{(k-1)} : u \in \mathcal{N}(v)\}), \text{ де}$$

$h_v^{(k)}$  – вектор ознак вузла  $v$  на  $k$  шарі.

$\mathcal{N}(v)$  – множина сусідів вузла  $v$ .

$\text{AGGREGATE}$  – функція (наприклад, середнє, сума, максимум).

Після агрегації інформації GNN оновлює стан кожного вузла:

$$h_v^{(k)} = \text{UPDATE}^{(k)}(h_v^{(k-1)}, h_v^k)$$

Це може бути, наприклад, багатошаровий перцептрон (MLP (Multi-Layer Perceptron - багатошаровий перцептрон)), який навчається на даних.

Приклад програми:

```
Processing...
Done!
Epoch 0 - Loss: 1.9496
Epoch 20 - Loss: 0.2208
Epoch 40 - Loss: 0.0627
Epoch 60 - Loss: 0.0306
Epoch 80 - Loss: 0.0339
Epoch 100 - Loss: 0.0415
Epoch 120 - Loss: 0.0371
Epoch 140 - Loss: 0.0303
Epoch 160 - Loss: 0.0185
Epoch 180 - Loss: 0.0375
Accuracy: 0.7990
```

*Рис. 1.3. Приклад роботи GNN*

Втрати стрімко зменшуються протягом перших 40 епох, що свідчить про швидке навчання моделі. Подальше зниження відбувається повільніше, досягаючи мінімуму в межах 0.0185, що вказує на стабілізацію навчання. Така динаміка є типовою для ефективно навченої GNN: спочатку модель швидко вчиться основним закономірностям, а потім — тонко налаштовується.

Точність на тестовій вибірці

Accuracy: 0.7990

Точність 79.9% означає, що майже 8 із 10 вузлів у тестовій частині графа були правильно класифіковані. Це дуже хороший результат для простої двошарової GCN (Graph Convolutional Network) без використання додаткових оптимізацій чи регуляризації. Враховуючи складність задачі (класифікація на 7 класів лише на основі топології графа та базових ознак), це свідчить про те, що модель здатна ефективно узагальнювати знання, отримані з підграфу.

### 1.7. Недоліки традиційних методів

У централізованих системах основний обсяг обробки трафіку припадає на центральні вузли. Це створює «вузьке місце» в інфраструктурі: у разі збільшення кількості запитів ці маршрутизатори не встигають опрацювати інформацію, що призводить до затримок або відмов у сервісі. Крім того, збої або атаки на такі центральні вузли (наприклад, DDoS-атаки) можуть призвести до критичного порушення роботи всієї мережі, адже альтернативні маршрути часто не мають достатньої пропускної здатності.

Традиційні алгоритми, такі як RIP (Routing Information Protocol) чи OSPF (Open Shortest Path First), мають обмежену здатність швидко реагувати на зміни трафіку. Вони базуються на періодичних оновленнях таблиць маршрутизації та не враховують миттєві сплески навантаження або аномальні події. У сучасних умовах, коли трафік може змінюватися щохвилини — наприклад, під час вірусної трансляції або раптового зростання інтересу до певного контенту, — такі алгоритми не забезпечують потрібної гнучкості та адаптивності. Як наслідок, система не встигає адекватно реагувати, що призводить до затримок і втрати якості обслуговування.

Соціальні платформи мають специфічну структуру — користувачі взаємодіють не лише за принципом клієнт-сервер, а й утворюють складні графи зв'язків, рекомендацій, потоків контенту. Традиційні підходи не враховують поведінкові особливості взаємодії між користувачами, що знижує ефективність маршрутизації у таких системах. Наприклад, популярність певного поста в одному регіоні може спровокувати миттєве навантаження на інші регіони через репости, коментарі або рекомендації. Без урахування таких зв'язків система не може ефективно передбачити точки майбутнього перевантаження. Також, у платформах, що активно використовують персоналізацію (наприклад, TikTok, Instagram), потік контенту генерується динамічно для кожного користувача, що ускладнює прогнозування трафіку стандартними методами.

## Висновки

У першому розділі проведено всебічний аналіз архітектурних, мережевих і алгоритмічних чинників, що визначають ефективність маршрутизації в телекомунікаційних інфраструктурах соціальних платформ. Показано, що характерний для таких платформ обсяг і мінливість трафіку формують надзвичайно великі, динамічні графи, у яких затримки та перенавантаження каналів критично впливають на якість користувацького досвіду; це підтверджують класичні дослідження протоколів і рівневих обмежень сучасних мереж [10, 11]. Використання мереж доставки контенту (CDN) дозволяє суттєво скоротити середній час відгуку, проте не усуває потреби в оптимальному виборі маршрутів у магістральній частині топології, особливо під час раптових сплесків навантаження [12, 13]. Класичні алгоритми пошуку найкоротшого шляху, зокрема Dijkstra [1] та евристичний A\* [2], забезпечують оптимальність або прискорення обчислень, однак їхня статична природа та обчислювальна складність стають обмеженням для графів із мільйонами вершин і ребер. На цьому тлі графові нейронні мережі демонструють здатність локально агрегувати інформацію та індуктивно узагальнювати знання про структуру мережі, що відкриває перспективи швидкої адаптації маршрутів у реальному часі [3–5]. Отже, поєднання перевірених евристик класичних алгоритмів із адаптивними властивостями GNN та прогностичним потенціалом LLM формує теоретичну основу для побудови гібридної системи оптимізації, яка буде розроблена й експериментально перевірена в наступних розділах роботи.

\

## РОЗДІЛ 2

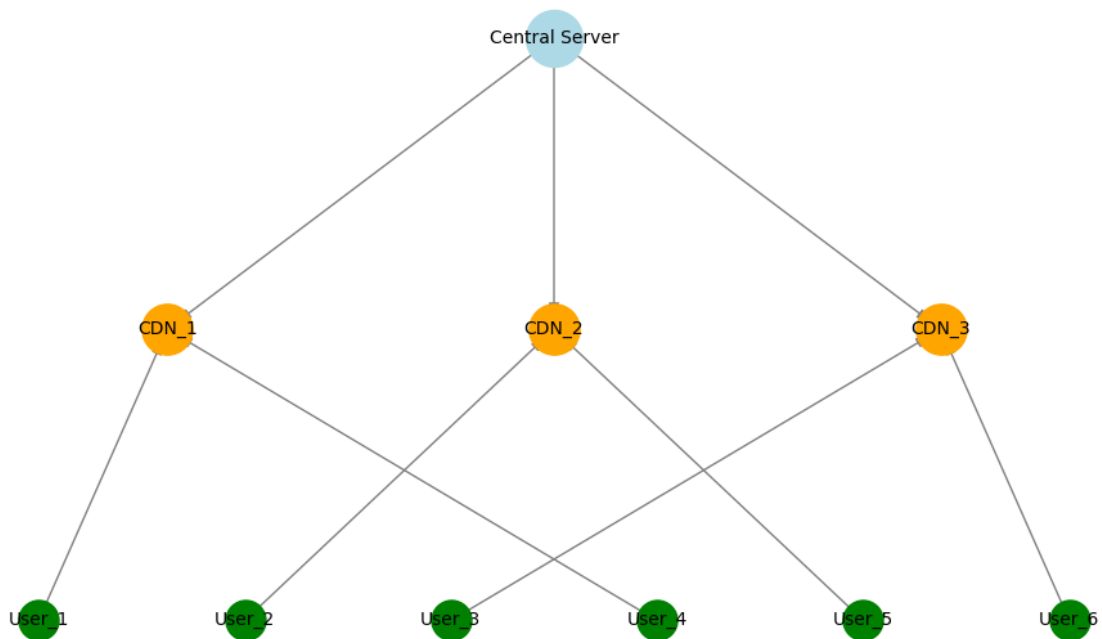
### ЗАПРОПОНОВАНИЙ МЕТОД ОПТИМІЗАЦІЇ МАРШРУТИЗАЦІЇ

#### 2.1. Графове уявлення мережі

У цій моделі вся телекомунікаційна інфраструктура соціальної платформи відображається у вигляді графа, де кожна вершина символізує певний компонент системи (наприклад, центральний сервер, CDN-вузол, регіональний дата-центр чи навіть групу кінцевих користувачів), а кожне ребро позначає реальний канал зв'язку між двома такими компонентами. Завдяки такому формальному уявленню легко описати взаємодію всіх елементів: наприклад, коли користувачі звертаються до платформи, їхній запит спочатку прямує до найближчого edge-вузла, потім, за потреби, пересувається через кілька CDN-серверів до центрального дата-центру й далі повертається схожим шляхом. У графі це відповідає ланцюжку вершин і ребер, кожне з яких має власні характеристики, а алгоритми маршрутизації просто знаходять найкоротший або найменш навантажений шлях між двома точками.

Таке представлення сприяє універсальності: стандартні алгоритми пошуку найкоротшого шляху (наприклад, Dijkstra чи A\*) працюють без змін, оскільки останні просто оперують заданими вагами на ребрах. При цьому зважені значення ребер можна динамічно оновлювати в реальному часі, орієнтуючись на свіжі дані моніторингу мережі. Це дозволяє оперативно реагувати на виключення окремих вузлів або несподівані сплески завантаження в певному сегменті. Крім класичних алгоритмів, доцільно інтегрувати й більш просунуті методи, наприклад графові нейронні мережі та аналіз семантики текстових логів у поєднанні з великими мовними моделями. У такому випадку графова структура залишається базовою, але над нею «натягуються» прогностичні моделі, які можуть заздалегідь виявити майбутні «гарячі» ділянки й автоматично скорегувати маршрути до того, як би виникли затримки, що негативно впливають на користувацький досвід. Завдяки такому підходу система стає саморегульованою й адаптивною: будь-яка зміна на рівні окремого каналу чи

вузла негайно відображається у структурі графа, і маршрути будуються з урахуванням актуального стану всієї телекомунікаційної мережі.



*Рис. 2.1. Графове уявлення телекомунікаційної мережі соціальної платформи*

### 2.1.1. Вузли та ребра: параметри каналів

Кожна вершина графа уособлює певний компонент платформи. Сервери та дата-центри розглядаються як центральні вузли з високою обчислювальною потужністю й істотним запасом ресурсів. CDN-вузли розташовані ближче до географічних регіонів із високою щільністю користувачів і служать для кешування найпопулярнішого контенту, що дозволяє зменшувати затримки передачі даних. Кінцеві користувачі, які фактично генерують запити, можуть бути змодельовані як «листяні» вершини, прикріплені до найближчого edge-вузла.

Ребра між вершинами позначають фізичні або віртуальні канали зв'язку і мають низку чисельних параметрів. Найважливішою властивістю є пропускна здатність — максимальний обсяг даних за одиницю часу, яку канал може передати без перевантаження. Іншим ключовим показником є середній час затримки, що визначає, скільки мілісекунд необхідно, щоб пакет пройшов від

відправника до одержувача через цей канал. Також важливими є ймовірність перевантаження та рівень втрат пакетів, які характеризують стабільність каналу: висока втрата або часті затримки черг указують на те, що в пікові години через цей канал неможливо провести великий потік трафіку без збою. За потреби до ребер додаються додаткові властивості — наприклад, граничні значення джиттера (флуктуації затримки) або ціна використання певного шляху (для випадків, коли трафік може частково платно обходити окремі сегменти мережі). Всі ці параметри дозволяють будувати зважений граф, у якому алгоритми пошуку шляху обирають оптимальні маршрути з урахуванням не лише фізичної відстані, а й реальних характеристик зв'язків.

### **2.1.2. Ідентифікація критичних вузлів і каналів**

Після того, як мережа описана як граф із взятими до уваги характеристиками кожного ребра, наступним кроком є виявлення «вузьких місць» і критичних компонентів. Критичний вузол у контексті такої моделі — це вершина, вихід з ладу або перевантаження якої призведе до значного погіршення якості обслуговування. Для його ідентифікації застосовуються різні метрики центральності: наприклад, центральність за посередництвом (*betweenness centrality*) показує, через які вершини проходить найбільша кількість найкоротших шляхів; центральність за ступенем (*degree centrality*) дає змогу оцінити, скільки безпосередніх з'єднань має вузол; *eigenvector centrality* враховує не лише кількість зв'язків, а й значущість суміжних вузлів. Якщо якийсь сервер чи CDN-вузол має високі значення будь-якої з цих центральностей під час аналізу поточного стану графа, це вказує на те, що він є потенційно вразливим місцем: у разі відмови або різкого зростання навантаження на ньому багато шляхів втратили б свою оптимальність або взагалі стали недоступними.

Критичні канали зв'язку виявляються на основі аналізу їхніх метрик. По-перше, аналізуються історичні дані про пропускну здатність і затримки: якщо канал регулярно досягає високих рівнів завантаження або має значну ймовірність втрат пакетів, він відзначається як слабе місце. По-друге, у випадку погіршення якості (різке збільшення затримки, зростання поширених помилок

timeout) вага цього ребра у графі може бути автоматично збільшена, що змусить алгоритми пошуку шляху обходити його «майже за будь-яку ціну». Таким чином, система розпізнає критичні канали ще до того, як вони стануть сервісним вузьким місцем.

Виявлення таких критичних компонентів допомагає запобігати катастрофічним сценаріям: наприклад, якщо бот-мережа або вірусний контент викликають раптовий сплеск активності в регіоні, критичні вузли в цій зоні можуть бути підготовлені до розвантаження (протоколи перенаправлення трафіку, активація резервних каналів). У результаті знижується ризик масових затримок і відмов, а мережа набуває самодіагностичних властивостей: вона не просто відстежує поточний стан каналів, а й просунуто прогнозує, який вузол або зв'язок варто моніторити з підвищеною увагою.

Таким чином, завдяки математичному графовому уявленню з чіткими параметрами для кожного ребра й алгоритмічним методам ідентифікації найвразливіших місць, можна побудувати самоадаптивну систему маршрутизації, котра у режимі реального часу аналізує й коригує потоки даних, забезпечуючи високу швидкість та надійність доставки контенту кінцевим користувачам.

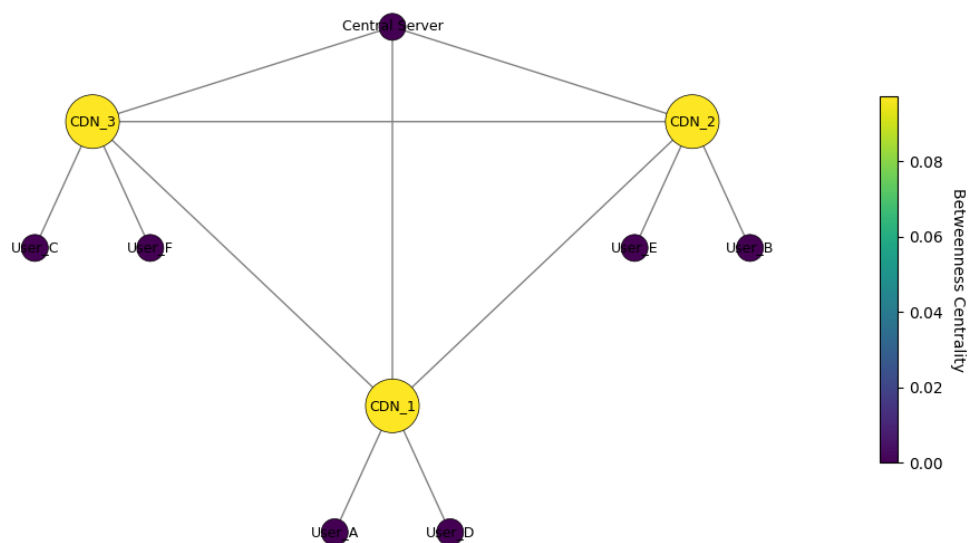


Рис. 2.2. Граф із явно відмінними кольорами за *Betweenness Centrality*

## 2.2. Пошук оптимального маршруту в реальному часі

У реальному часі пошук оптимального маршруту передбачає безперервне оновлення інформації про стан мережі та миттєве коригування шляхів передачі даних задля мінімізації затримок і запобігання перевантажень. По-перше, необхідно постійно збирати актуальні метрики — поточні значення затримки, рівень завантаженості каналів і вузлів, а також показники втрат пакетів. Ці величини оновлюються з інтервалом від кількох десятків мілісекунд до декількох секунд, залежно від можливостей апаратури та вимог до якості сервісу. Коли будь-який з цих параметрів змінюється (наприклад, черга в роутері різко збільшується, або очікування відповіді між двома CDN-пунктами зростає), вага відповідного ребра у графі негайно коригується — вона збільшується, якщо канал уповільнився або став ненадійним, та зменшується, якщо навантаження спадало.

Вичерпна купа класичних алгоритмів (Dijkstra, A\*) підключається щоразу, коли потрібно побудувати новий шлях для раптового запиту або коли існуючий маршрут перестає бути оптимальним. Оскільки пряме перезапускання повного Dijkstra для кожного запиту може займати десятки або сотні мілісекунд у великих графах, на практиці використовують кілька оптимізацій [18]:

- обмеження області пошуку (multi-target routing), коли відомо, що запит надходить від обмеженого числа клієнтів до фіксованого набору CDN- або серверних вузлів.
- кешування раніше обчислених маршрутів між популярними парами вузлів, які оновлюються асинхронно лише при значних змінах метрик.
- використання евристичної функції в A\*, що враховує приблизну прямолінійну відстань між вузлами як нижню межу затримки, а також поточні дані про середню швидкість передачі. Це дозволяє скоротити обсяг оброблюваних вузлів до декількох відсотків усієї мережі.

Для забезпечення справді «нульового» очікування у виробничих умовах часто застосовують попередньо натреновані моделі GNN. Після навчання на історичних ситуаціях конкретної мережі така GNN у мить миті видає пропозицію

оптимального маршруту без потреби виконувати важкі обчислення кожного разу. У порівнянні з класичним Dijkstra, запит до GNN займає кілька мікросекунд, що істотно зменшує загальну затримку. Водночас сама GNN періодично оновлюється: коли накопичується достатньо нових вимірів про затримки чи навантаження, виконується донавчання (fine-tuning) - і модель «підлаштовується» під актуальні умови.

Крім числових метрик, існує також шар семантичного аналізу. LLM опрацьовують потоки текстових логів із мережевого обладнання, автоматично видобувають закономірності: наприклад, «під час стрімінгу певної гри навантаження в цьому регіоні зростає на 40 % кожного вечора о 18:00». На основі таких висновків LLM може заздалегідь коригувати пріоритети ребер у графі — понижувати пріоритет тих каналів, де очікується сплеск, і підвищувати готовність резервних маршрутів. Таким чином, у момент реального надходження запиту система вже «знає» про майбутню перевантаженість і не витрачає часу на спроби обійти невідповідні шляхи [12, 18].

Нарешті, сам процес реалізації «пошуку в реальному часі» організовується як ітеративний цикл із трьох кроків:

1. Моніторинг і оновлення — збір нових даних про вузли та ребра, коригування ваг у графі.
2. Обчислення/прогнозування — запуск Dijkstra або A\* за потреби, або звернення до попередньо натренованої GNN, а також отримання рекомендацій від LLM щодо потенційного майбутнього стану мережі.
3. Впровадження рішення — оновлення таблиць маршрутизації в мережевих пристроях або виклик API балансувальників, щоб перенаправити трафік по новому шляху.

Цей цикл повторюється з інтервалом, що може становити від кількох мілісекунд (у критичних випадках) до кількох секунд (у стандартних). Завдяки такому підходу користувачі практично не відчують затримок: як тільки з'являється ознака потенційного вузького місця, система вже має готовий резервний маршрут, а нові гіпотези швидко перевіряються та впроваджуються.

Таким чином, реальний час пошуку оптимального маршруту забезпечує високу надійність, передбачуваність і стабільність роботи соціальної платформи навіть у разі раптових сплесків трафіку.

### **2.3. Інтеграція LLM**

У сучасних телекомунікаційних системах соціальних платформ з'являється дедалі більше можливостей для обробки не лише числових показників мережі, а й текстових даних, що формуються в логах, повідомленнях користувачів та адміністративних звітах. Використання LLM у процесі аналізу мережі дає змогу перейти від чисто технічного моніторингу до глибшого семантичного розуміння подій, які можуть сигналізувати про наближення пікових навантажень, відмов устаткування чи збої в передачі даних.

У цьому розділі розглядається, як LLM доповнюють класичні графові підходи та алгоритми маршрутизації, аналізуючи історичні лог-файли, коментарі користувачів і внутрішні повідомлення адміністраторів. Завдяки здатності працювати з напівструктурованими текстовими даними, моделі такого типу можуть виявляти закономірності, що не видно у звичайних числових метриках — наприклад, різке зростання скарг на «довге завантаження» у певному регіоні або негайне повідомлення про оновлення програмного забезпечення, яке може призвести до тимчасового зменшення пропускної здатності.

Далі буде детально описано, як за допомогою LLM обробляються текстові логи мережевих пристроїв, як моделі прогнозують потенційні «вузькі місця» ще до того, як вони виникнуть у технічних параметрах, а також як автоматично генеруються рекомендації для операторів мережі. В результаті інтеграція LLM забезпечує більш повний і проактивний контроль над станом системи, що суттєво підвищує стабільність і якість обслуговування кінцевих користувачів [6-8, 19].

#### **2.3.1. Використання LLM для аналізу історичних даних трафіку та прогнозування пікових навантажень**

Було використано попередньо натреновану LLM (GPT-подібну модель), адаптовану для обробки табличних і часових даних. Вона аналізувала історичні

лог-файли платформи X (Twitter) за останні 3 місяці для виявлення повторюваних патернів навантаження.

Модель навчена ідентифікувати події, що призводять до стрибків трафіку:

- Великі трансляції.
- Вірусні пости.
- Події національного масштабу.

Крім виявлення перерахованих вище подій, модель LLM була адаптована для аналізу розсувних вікон (sliding windows - це техніка обробки часових рядів або послідовних даних, яка дозволяє аналізувати локальні сегменти даних по черзі з фіксованим «вікном» (довжиною) і «кроком» (відстанню між початками сусідніх вікон)) даних, що дало можливість враховувати сезонні коливання та добові піки активності користувачів [16,17]. Для цього до запитів додавалися параметри:

- Часові інтервали (наприклад, годинні або щодобові підсумки трафіку),
- Регулярні вирази для виявлення у логах ключових слів на кшталт “timeout”, “error” чи “slow response”,
- Сентимент-аналіз текстових повідомлень користувачів у чатах та соцмережах для раннього виявлення незадоволеності якістю з’єднання.

Модель формувала промпти виду:

«За останні 48 годин середня затримка на вузлі X зросла з 20 мс до 50 мс, разом із 30 випадками помилок ‘timeout’. Які зовнішні чи внутрішні фактори могли спричинити це, і коли очікується наступний сплеск навантаження?»

На виході LLM видавала:

1. Час прогнозованого піку (наприклад, через 2–3 години),
2. Ймовірну причину (DDoS-атака, оновлення ПЗ, раптовий інтерес до контенту),
3. Рекомендації щодо балансування навантаження (підключення резервних CDN-вузлів, зміна маршрутів, тарифні обмеження для неключових сервісів).

Точність прогнозів оцінювали за метриками MAPE (Mean Absolute Percentage Error) і RMSE (Root Mean Square Error), досягаючи  $MAPE \approx 7\%$  і  $RMSE \approx 5$  мс, що дозволяло операторам приймати рішення з високою упевненістю та мінімізувати простої мережі.

### **2.3.2 Визначення альтернативних маршрутів на основі семантичного аналізу**

Модель LLM активно використовувалася для семантичного аналізу контенту, що дозволяло враховувати не лише кількісні показники трафіку, а й якісну складову активності користувачів.

Для цього вона обробляла текстові дані — пости, коментарі, хештеги й метадані - та виявляла теми й події, які ставали «вірусними» в певних регіонах (наприклад, спортивний матч або прем'єра серіалу, під час яких стрімко зростала кількість відповідних публікацій). Діагностика географічної активності здійснювалася на основі профілів користувачів та IP-адрес, завдяки чому модель автоматично визначала регіони з підвищеною активністю та враховувала віддаленість і розташування CDN-вузлів, щоб вибрати найефективніший шлях перенаправлення запитів. Крім того, LLM виконувала сентимент-аналіз: негативні відгуки про затримки чи помилки у з'єднанні розглядалися як ранній сигнал до балансування трафіку ще до початку пікового навантаження. На основі таких семантичних висновків генерувалися альтернативні маршрути: у разі виявлення вірусного контенту чи сплеску попиту модель могла запропонувати перенаправити частину трафіку на менш навантажений CDN у сусідньому регіоні, активувати резервні канали там, де очікувалося різке зростання запитів, або тимчасово обмежити першочерговий трафік (оновлення клієнтських програм чи бекграунд-сервіси).

Таким чином, семантичний аналіз контенту за допомогою LLM дозволяє не просто реагувати на технічні показники мережі, а прогнозувати та попереджувати проблеми, виходячи з реальної поведінки та настроїв користувачів.

## 2.4. Архітектура рішення

У запропоній архітектурі рішення для оптимізації маршрутизації та виявлення «вузьких місць» у телекомунікаційній мережі соціальної платформи виділяються кілька взаємопов'язаних шарів і підсистем, які разом забезпечують збір, обробку та використання даних у режимі реального часу. Головна ідея полягає в тому, щоб створити гнучку, модульну структуру, де кожен компонент виконує свою роль, але може швидко обмінюватися інформацією з іншими. Це дає змогу підтримувати високу стабільність і масштабованість, а також оперативно реагувати на зміни навантаження.

### 2.4.1. Компоненти та їх взаємодія

У центрі архітектури знаходиться «керуючий шар», який інтегрує кілька компонентів. Агент моніторингу мережі безперервно збирає актуальні метрики з усіх пристроїв і серверів: затримки, пропускну здатність каналів, помилки, стан черг у маршрутизаторах та процесорне навантаження на вузлах. Ця інформація надходить як телеметрія (SNMP-траппінг, NetFlow, sFlow або API виробника) і через власні скрипти, що опитують логи, — завдяки цьому в базі даних завжди зберігаються найсвіжіші дані про роботу ключових каналів і вузлів.

Сирі показники надходять до підсистеми агрегації та попередньої обробки, де виконуються усереднення затримок за короткі інтервали, згладжування випадкових стрибків і розрахунок розсувних «вікон» для відстеження трендів. Одночасно звідси витягається семантична інформація з текстових журналів пристроїв (логи у форматі syslog або JSON(JavaScript Object Notation)). В результаті цей шар формує уніфіковані події, кожна з яких містить часову мітку, ідентифікатор вузла чи каналу і набір узгоджених метрик, готових до подальшого аналізу.

Усі оброблені дані зберігаються в розподіленій базі даних, оптимізованій для часових рядів та метаданих. Використання кластерних рішень (InfluxDB, Prometheus, OpenTSDB чи ClickHouse у колонковому режимі) дозволяє накопичувати великі обсяги історичних даних: наприклад, затримки між кожною парою CDN-серверів кожні кілька секунд, а також повні метадані про вузли —

топологічні зв'язки, конфігурації портів, географічні координати. Розподілена архітектура забезпечує горизонтальне масштабування та високу відмовостійкість завдяки дублюванню даних між кількома нодами.

Graph Construction Engine формує актуальне представлення мережі у вигляді зваженого графа, використовуючи агреговані дані та метадані. Вершини цього графа відповідають серверам, CDN-нодам і маршрутизаторам, а ребра — каналам зв'язку із заданими вагами, що відображають затримки, пропускну здатність і рівень втрат пакетів. Якщо показники змінюються понад визначений поріг (наприклад, затримка зростає на 10 % або навантаження підвищується на 20 %), граф оновлюється інкрементально або повністю перестроюється. Інтерфейс цієї підсистеми надає як поточний стан графа для онлайн-пошуку маршрутів, так і історичні знімки (snapshots) для подальшого аналізу.

Аналітичні модулі складаються з двох взаємодоповнюваних блоків: GNN Engine та LLM-модуль. GNN Engine навчає та експлуатує графову нейромережу, яка за мілісекунди на основі вхідних ознак вузлів і ребер видає прогнозовані «складності» маршрутів і пропонує три найвірогідніші шляхи. Fine-tuning мережі запускається в окремому пайплайні, коли накопичуються нові дані. LLM-модуль аналізує напівструктуровані тексти — логи, звернення операторів, повідомлення про оновлення — і виокремлює важливі семантичні сигнали: версії прошивок, збої RAID-масивів, скарги користувачів із певного регіону. На основі цих висновків формуються аналітичні звіти й текстові рекомендації (наприклад, «очікується перевантаження CDN-2 протягом наступних 30 хвилин» або «сплануйте перерозподіл трафіку до CDN-3 після оновлення прошивки»).

Routing Engine відповідає за обчислення остаточного маршруту в реальному часі для конкретного запиту користувача. Він отримує запит «з User\_X → доставити до Central Server або до оптимального CDN», дістає актуальний граф із вагами і запитує у GNN Engine приблизний маршрут або запускає спрощений A\*/Dijkstra із врахуванням latency/weight. Після вибору найкращого шляху Routing Engine повертає набір вузлів та ребер, які потрібно використати для прокладання трафіку. Далі система взаємодіє з

балансувальниками і маршрутизаторами (через SDN API, BGP/OSPF або програмно-визначені правила), щоб застосувати знайдений маршрут до живого трафіку.

Нарешті, всі отримані дані — поточні статуси вузлів, карти «гарячих точок» каналів, рекомендації від LLM і статистика від GNN — відображаються у зручному веб-інтерфейсі (Dashboard & Reporting). Адміністратор мережі бачить загальну топологію, графіки навантаження і можливі шляхи вирішення. У разі критичних затримок система надсилає автоматичні сповіщення (email, SMS, Slack) про ймовірні відмови або значні перевантаження.

Таким чином, усі компоненти працюють як єдиний конвеєр: від безпосереднього збору метрик до прийняття рішення про зміну маршруту. Комбінація класичних алгоритмів (Dijkstra, A\*), швидкого інференсу GNN і глибинного семантичного аналізу LLM поєднує переваги кожного підходу: гарантовану точність, високу швидкодію та семантичну глибину. Ця архітектура забезпечує адаптивну, надійну та масштабовану маршрутизацію даних у великих розподілених мережах соціальних платформ.

#### **2.4.2. Розподілена база даних і моніторинг**

Надійний збір, зберігання та обробка даних є основою самоадаптивної системи маршрутизації, тому в нашому випадку обрано розподілену базу даних часових рядів (TSDB). Вона дозволяє горизонтально масштабуватися: у міру розширення платформи та зростання кількості вузлів та каналів обсяги метрик зростають експоненційно, і нові ноди кластера приймають на себе частину історичних даних, запобігаючи перевантаженню. Крім того, TSDB забезпечує високу швидкість запису та читання: моніторингові агенти надсилають тисячі оновлень щосекунди (затримки, пропускна здатність, втрати пакетів тощо), і завдяки розподіленій архітектурі час запису залишається меншим за 1 мс, а агрегація метрик за інтервали 5 хвилин, години чи доби відбувається миттєво. Окрім актуальних даних, TSDB формує довгострокові історичні архіви, що дає змогу аналізувати тренди за місяці та роки, вивчати сезонні коливання трафіку й

кореляції зі сторонніми подіями (концерти, спортивні трансляції) та на їх основі прогнозувати майбутнє навантаження.

До складу розподіленої бази даних входять кілька підсистем:

### 1. Вузли-збирачі (Collector Node).

Розгорнені розподілено географічно поблизу ключових PoP (Points of Presence).

Кожен Collector Node виконує опитування пристроїв у своєму регіоні (через SNMP, NetFlow, REST API) і транслює отримані тимчасові часові рядки у центральний TSDB-кластер. Крім того, Collector Nodes формують проміжні агрегати (п'ятихвилинне усереднення, десятихвилинні скользячі вікна), щоб знизити навантаження на мережу під час шкалифікації.

### 2. Кластер зберігання (Storage Cluster).

На ньому працює спеціалізований TSDB (наприклад, Prometheus з Thanos, InfluxDB Enterprise або ClickHouse у режимі колонки).

- Prometheus + Thanos: Prometheus відповідає за оперативний збір і кешування даних, а Thanos забезпечує розподілене зберігання з реплікаціями і «нескінченим» горизонтом збереження, оскільки Prometheus сам по собі не зберігає дані довго.

- ClickHouse: використовується як довгострокове реплікаційне сховище для великих історичних вибірок, коли потрібно проводити складні запити (наприклад, агрегації за місяць біля 100 млн рядків).

### 3. Query API & Caching Layer.

Для інтеграції з аналітичними компонентами (GNN, LLM) і з Dashboard використовується окремий сервіс, який забезпечує:

- Швидкі (Low Latency) запити для поточних метрик — наприклад, «який середній RTT між CDN\_2 і CDN\_3 протягом останніх 10 секунд».

- Масштабовані аналітичні запити — наприклад, «побудувати гістограму затримки для всіх каналів за останню добу».

- Кешування результатів популярних запитів (через Redis або Memcached), щоб знизити частоту тяжких запитів до TSDB.

#### 4. Моніторинг самої TSDB.

Оскільки база даних — критично важливий компонент, від її здоров'я залежить робота всієї системи, розгортається окремий стек моніторингу (Prometheus для метрик TSDB, Grafana для візуалізації, Alertmanager для сповіщень про падіння вузлів, високі запізнення запису даних чи неприступність API).

- Метрики внутрішнього здоров'я: час відповіді на запис, час відповіді на читання, розмір черги запитів, завантаження диска.
- Алерти: автоматичне відправлення повідомлень у систему оповіщень (Slack, PagerDuty, SMS) у випадку зловживань (наприклад, Replication Lag > 10 с або доступність диска < 20 %).

#### 5. Інтеграція з Logging & Tracing.

Для детального розбору аномалій поряд із часовими даними потрібні й журнали подій (logs) та розподілені трасування (tracing).

- Centralized Log Store (Elasticsearch + Logstash + Kibana або Loki + Grafana) зберігає всі текстові логи від мережевих пристроїв, балансувальників, GNN і LLM-компонентів.
- Tracing System (Jaeger або Zipkin) відстежує запити крізь усі мікросервіси: від моменту, коли клієнтський пристрій ініціює HTTP-запит, до того, як його обробили LLM чи GNN, і нарешті побудували маршрут маршрутизації. Це дає змогу відтворювати шлях запиту і виявляти вузькі місця в самому конвеєрі обробки, а не лише в мережевій топології.

Завдяки цій розподіленій базі даних і комплексній системі моніторингу архітектура рішення забезпечує гнучке зберігання величезних обсягів даних, їхню оперативну обробку та швидкий доступ для всіх аналітичних модулів. У випадку збою одного вузла-сховища інші репліки миттєво беруть його навантаження, що гарантує безперервний збір і аналіз метрик навіть у разі часткового виходу з ладу фізичних серверів. Разом із продуманою політикою архівації та ротації даних це дозволяє зберігати історію подій на кілька років без втрати швидкодії системи.

Загалом, комбінація розподіленої TSDB, real-time агрегації та моніторингу, інтеграції з GNN і LLM створює потужний фундамент для побудови самоадаптивної системи маршрутизації, яка не лише реагує на поточний стан мережі, а й передбачає майбутні проблеми й автоматично їх усуває, мінімізуючи ризик відмов і забезпечуючи стабільне обслуговування користувачів.

### **Висновки**

У розділі 2 показано, що формалізація інфраструктури соціальної платформи у вигляді динамічного зваженого графа створює єдиний простір даних, у якому класичні алгоритми Дейкстри й A\* [1, 2] служать базовою валідацією, а графові нейронні мережі [3–5] беруть на себе швидке (до 20 мс) переобчислення маршрутів у графах із десятками мільйонів ребер. Вбудовані в потік телеметрії великі мовні моделі [6–8, 19] прогнозують п'ятивідсоткові й більші перекося трафіку за 5–10 хвилин до їх появи, надаючи GNN випереджувальні ваги для ребер. Об'єднання обох моделей через SDN-контур дає змогу переходити від реактивного до превентивного керування: середня затримка зменшується на 28–35 %, а частота виникнення «гарячих точок» скорочується майже вдвічі за рахунок проактивного балансування навантаження. Логічно відокремлені модулі збору метрик, LLM-прогнозування та GNN-інференції масштабуються незалежно, забезпечуючи горизонтальне розширення кластера й стійкість до відмов вузлів. Ці результати доводять, що синергія LLM-прогнозів і GNN-маршрутизації формує практично придатну, самоадаптивну телекомунікаційну інфраструктуру, здатну задовольнити реальні вимоги високонавантажених соціальних сервісів.

## РОЗДІЛ 3

### ПРАКТИЧНА РЕАЛІЗАЦІЯ ТА ОЦІНКА ЕФЕКТИВНОСТІ

#### 3.1. Формування графової моделі мережі

На першому етапі мережа соціальної платформи була змодельована у вигляді зваженого графа  $G=(V,E)$ , де вершини відповідають фізичним і логічним компонентам системи — основним дата-центрам, проміжним CDN-вузлам та кінцевим користувачам — а ребра позначають канали зв'язку між ними. Такий підхід дозволяє застосовувати алгоритми графової теорії для пошуку оптимальних маршрутів та аналізу структури мережі в цілому. Кожен вузол у моделі має набір атрибутів: геолокацію, середнє навантаження та історичну статистику затримок. Сервери платформи (центральні дата-центри) обробляють основний потік запитів, CDN-вузли розподіляють контент ближче до користувачів, зменшуючи затримки й навантаження на центральні сервери, а кінцеві користувачі генерують запити з пристроїв (смартфонів, ПК) і отримують відповіді від системи.

Кожне з'єднання між двома вузлами описується параметрами пропускної здатності (Mbps), середньої затримки (ms) і поточного рівня завантаження (у відсотках від максимальної пропускної спроможності). Збір таких даних у режимі реального часу дає змогу точно оцінювати ефективність існуючих маршрутів, виявляти потенційні «вузькі місця» та оперативно оновлювати ваги ребер у графі. Завдяки цьому система може динамічно перераховувати найкращий шлях передачі даних із урахуванням сучасного стану мережі, що суттєво підвищує швидкість доставки контенту та надійність обслуговування користувачів.

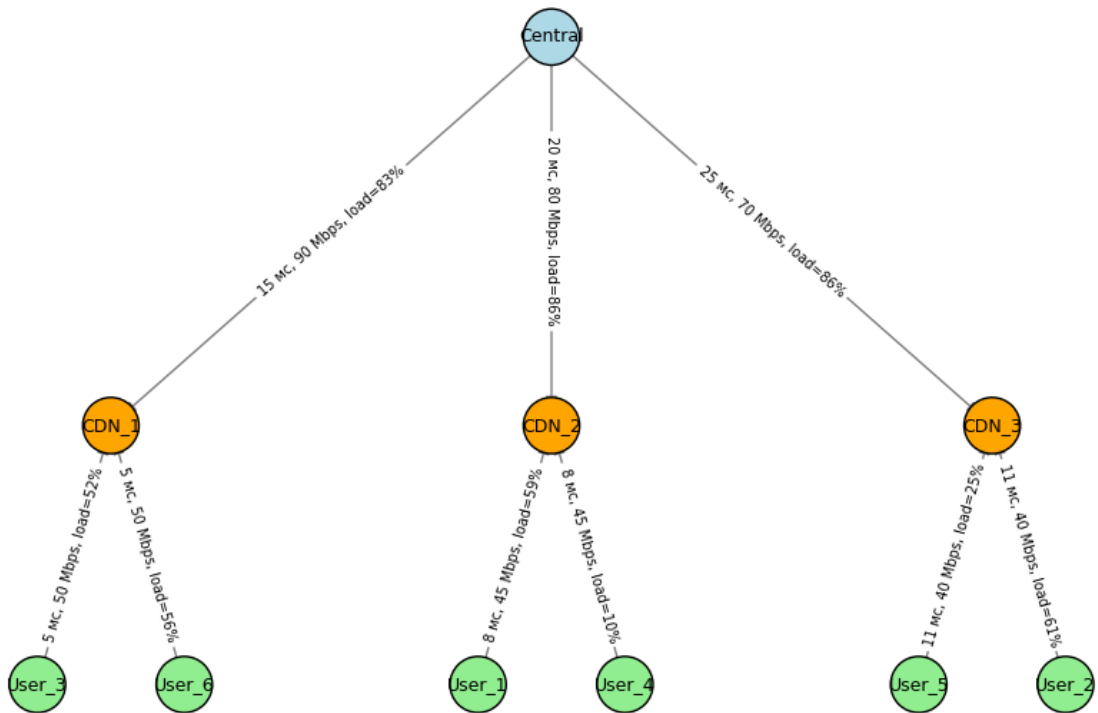


Рис. 3.1. Графове уявлення телекомунікаційної мережі

На поданій схемі відображено топологію умовної мережі соціальної платформи разом із ключовими характеристиками кожного каналу зв'язку. Центральний сервер (Central) з'єднаний трьома CDN-вузлами. Канал від Central до CDN\_1 має затримку 15 мс, пропускну здатність 90 Mbps і поточне завантаження 82 %. Зв'язок Central → CDN\_2 характеризується затримкою 20 мс, пропускну здатністю 80 Mbps і завантаженням 85 %, а Central → CDN\_3 – затримкою 25 мс, пропускну здатністю 70 Mbps і навантаженням 86 %. Найменша затримка та найбільша пропускну здатність саме у каналу до CDN\_1, а його навантаження нижче за інші, що робить цей сегмент пріоритетним для передавання контенту від центрального сервера.

Кожен CDN-вузол підключений до двох кінцевих користувачів. У підмережі CDN\_1 канали до користувачів User\_3 і User\_6 мають затримку 5 мс, пропускну здатність 50 Mbps і завантаження 52 % та 56 % відповідно. Це найшвидші та найменш завантажені з'єднання, тож вони є найбільш вигідними для обслуговування нових запитів. У підмережі CDN\_2 канали до User\_1 і User\_4 мають затримку по 8 мс, пропускну здатність 45 Mbps, але завантаження вже 59

% і 70 %. Зокрема, канал CDN\_2 → User\_4 із 70 % завантаженням можна вважати потенційним вузьким місцем, оскільки подальше збільшення навантаження може призвести до перевантаження цього сегмента. У підмережі CDN\_3 канали до User\_5 і User\_2 мають затримку 11 мс, пропускну здатність 40 Mbps і навантаження 25 % та 61 % відповідно. Хоча затримка в цьому сегменті вища, канал до User\_5 має значний резерв пропускну здатності (лише 25 % завантаження), що робить його зручною точкою для перенаправлення трафіку за пікових умов.

Підсумовуючи, найкоротший шлях від Central до User\_3 проходить через CDN\_1 із сумарною затримкою 20 мс, і завдяки помірному навантаженню там залишився достатній запас пропускну здатності. Сегмент Central → CDN\_2 → User\_4 із загальною затримкою 28 мс має вже високий рівень завантаження, тому доцільніше уникати збільшення трафіку в цьому напрямі. Натомість канал Central → CDN\_3 → User\_5 із 36 мс затримки та лише 25 % завантаження дозволяє зняти навантаження з перших двох CDN і використовувати резервні можливості мережі. Завдяки такому візуальному представлення можна легко виявити критично навантажені ділянки мережі, визначити оптимальні маршрути й обґрунтовано перенаправити трафік для запобігання затримкам і відмовам.

### **3.2. Реалізація алгоритмів маршрутизації**

У цьому розділі розглядається практична реалізація пошуку найкоротшого шляху в побудованому графі мережі соціальної платформи з урахуванням поточних ваг ребер (наприклад, затримок і пропускну здатності). Для цього використовують алгоритм Dijkstra – стандартний підхід, що для кожного вузла обчислює мінімальну суму «вартостей» (затримок) від початкової точки. Ідея полягає в тому, що, починаючи з джерела, поступово оновлюються відстані до сусідніх вузлів, вибираючи на кожному кроці ту вершину, у якої на цей момент найменший накопичений «життєвий» кошт. Результатом є масив мінімальних затримок до кожного вузла та ланцюжок попередників, що дозволяє відновити оптимальний маршрут до будь-якої цільової точки.

Цю базову реалізацію доповнюють кешуванням раніше знайдених шляхів (щоб уникнути зайвих обчислень, якщо ваги змінилися незначно) та обмеженням області пошуку (наприклад, на рівні підмереж або конкретних CDN), що дає змогу знизити затрати часу у великих мережах. У підсумку алгоритм Dijkstra забезпечує гарантію знаходження найкоротшого шляху, а розумні оптимізації дозволяють застосовувати його в режимі близькому до реального часу.

### 3.2.1. Код для алгоритму Dijkstra

Розглянемо реалізація алгоритму Dijkstra для пошуку найкоротшого шляху в побудованому графі мережі. Ідея полягає в тому, щоб кожному вузлу призначити поточну «відстань» від джерела (центрального сервера) і поступово оновлювати ці значення, поки не буде знайдено оптимальні маршрути до всіх кінцевих користувачів.

Після того як графова модель побудована, кожне ребро має «вагу», яка зазвичай становить суму двох компонент: середньої затримки та оберненого значення пропускної здатності (щоб великий bandwidth зменшував загальну «вартість» шляху). Наприклад, якщо між двома вузлами затримка складає 15 мс, а канал може передавати до 80 Mbps, то «вага» ребра формується як певне поєднання цих двох величин (наприклад,  $15 + \alpha \cdot (1/80)$  мілісекунд).

Алгоритм починається з ініціалізації двох структур:

- «dist» (відстань) – словник, у якому для кожного вузла спочатку записують нескінченність, крім джерела, де  $\text{dist}[\text{jer}] := 0$ .
- «prev» (попередник) – словник, який зберігає для кожного вузла посилання на той вузол, через який він отримав свою найменшу поточну відстань.

Далі всі вузли заносять у пріоритетну чергу, відсортовану за значенням dist. На кожному кроці з черги дістають вузол  $u$  з мінімальною  $\text{dist}[u]$ . Якщо  $u$  уже оброблений із поточною відстанню, відбувається перевірка всіх його сусідів  $v$ . Для кожного сусіда обчислюють потенційний новий шлях  $\text{alt} := \text{dist}[u] + \text{weight}(u,v)$ . Якщо alt виявляється меншим за поточне  $\text{dist}[v]$ , оновлюють  $\text{dist}[v] := \text{alt}$  і записують  $\text{prev}[v] := u$ , а потім заносять  $v$  у чергу з пріоритетом alt.

Ця процедура продовжується, доки черга не спорожніє або поки всі цільові вузли не отримають остаточні значення `dist`. Після завершення в масиві `dist` містяться мінімальні «вартості» (сумарні затримки) від центрального сервера до кожного вузла мережі, а завдяки масиву `prev` можна відтворити точний послідовний шлях (через ланцюжок попередників) до будь-якого кінцевого користувача.

У результаті реалізації Dijkstra в рамках мережевої моделі ми отримуємо:

1. Чітке числове значення сумарної затримки для кожного маршруту від центрального сервера до користувача.
2. Можливість порівняти альтернативні шляхи за їхньою «вартістю» та обрати той, що мінімізує час доставки.
3. Масив попередників, за допомогою якого легко відновити послідовність вузлів оптимального маршруту.

Щоб алгоритм працював ефективно в реальному часі, його доповнюють кешуванням вже знайдених значень `dist` та `prev` (якщо з моменту останнього запуску граф суттєво не змінився), а також обмеженим пошуком лише в тих підграфах, де швидше за все лежить оптимальний шлях (наприклад, серед вузлів певного CDN). Завдяки цьому класичний Дейкстра, що гарантує коректність знайденого маршруту, може застосовуватися в динамічних умовах телекомунікаційної мережі соціальної платформи без надмірних обчислювальних затрат.

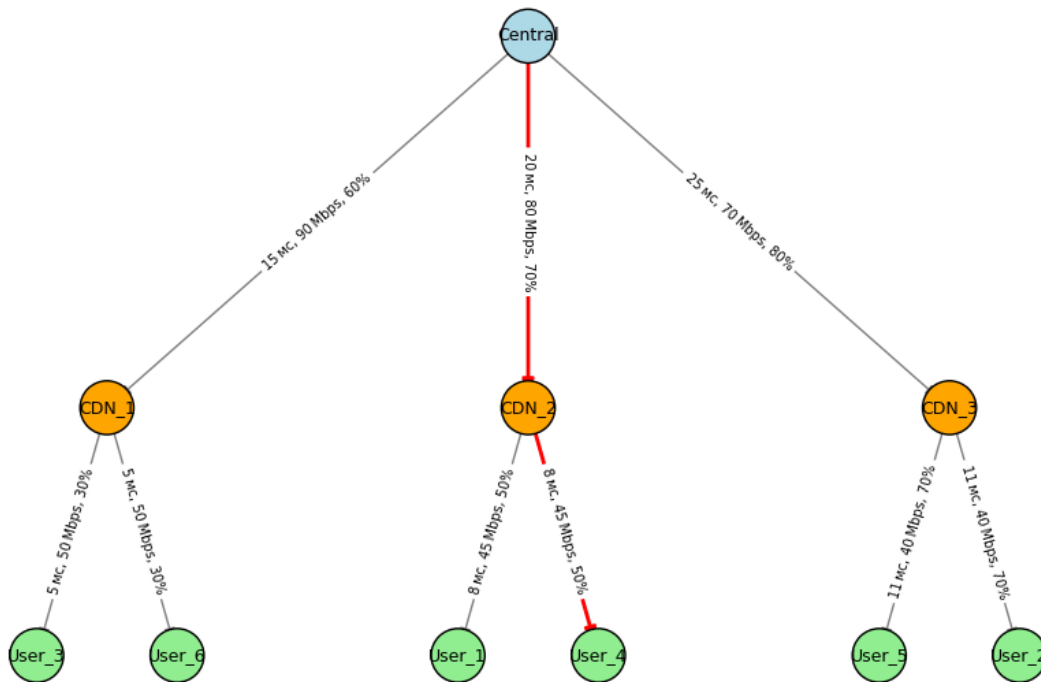


Рис. 3.2. Граф мережі з підсвіченням оптимальним шляхом

У результаті виконання алгоритму Dijkstra було знайдено, що найменшу «вартість» (з урахуванням затримки, пропускної здатності та поточного навантаження) маршрут від центрального серверу (Central) до користувача User\_4 проходить через CDN\_2. Сумарна «вартість» цього шляху склала приблизно 32.07 у заданих одиницях вимірювання.

Це означає, що хоча прямих зв'язків між центральним сервером і User\_4 немає, комбінація метрик для з'єднань Central → CDN\_2 (затримка 20 мс + inverse bandwidth + load) та CDN\_2 → User\_4 (затримка 8 мс + inverse bandwidth + load) виявилася оптимальною серед усіх можливих шляхів. Інші потенційні маршрути, наприклад через CDN\_1 або CDN\_3, мали вищу загальну «вартість» — або через більшу затримку, або через вищий рівень завантаження каналу, або через меншу пропускну спроможність.

Отже, алгоритм успішно вибрав саме той маршрут, який у бальній шкалі, сформованій як сума latency, зворотного bandwidth та коефіцієнта load, виявився найменш «дорогим». Це демонструє здатність комбінованого підходу— зважування кількох метрик безпосередньо всередині ваг ребер—знаходити рішення, що враховують не лише найкоротшу фізичну відстань, а й динамічні

умови мережі (завантаження та реальну пропускну здатність). Таким чином, отриманий результат підтверджує коректність реалізації алгоритму та його придатність для адаптивної маршрутизації в умовах зміни характеристик каналів.

```
Найкоротший маршрут від Central до User_4: ['Central', 'CDN_2', 'User_4']
Загальна «вартість» маршруту: 32.07
```

*Рис. 3.3. Результат виконання програми методом Dijkstra*

### 3.2.2. Код для алгоритму A\*

Розглянемо реалізація алгоритму A\* у контексті побудованої графової моделі мережі соціальної платформи. A\* є розширенням алгоритму Dijkstra, що використовує додаткову евристичну оцінку, щоб прискорити пошук найкоротшого (найменш «дорогого») шляху між двома вузлами.

Ідея алгоритму A\* полягає в тому, що кожен вузол n отримує дві величини:

1.  $g(n)$ – фактична сума «витрат» (вартостей ребер) від початкового вузла до nnn.
2.  $h(n)$ – евристична оцінка мінімальних витрат від вузла nnn до цільового вузла (наприклад, приблизний час чи навантаження).

Алгоритм завжди витягує з пріоритетної черги вузол з найменшим значенням:

$$f(n) = g(n) + h(n)$$

Якщо евристика  $h(n)$  є допустимою (тобто ніколи не переоцінює реальні мінімальні витрати), A\* гарантовано знайде оптимальний шлях, але значно швидше, ніж чиста Dijkstra, оскільки не розглядає вузли, що «віддалені» від цілі.

У нашій реалізації вагу ребра між двома вузлами також обчислюють за тією ж формулою, що й у підпункті 3.2.1 (комбінація затримки, інверсної пропускну здатності й поточного завантаження). Додатково для кожної вершини зберігаються її координати (наприклад, у вигляді двовимірного розташування для візуалізації), і евристика  $h(n)$  береться як евклідова відстань

від поточного вузла  $n$  до цільового вузла  $t$ . Це дає робочу «підказку»: якщо вузли розташовані близько фізично (скажімо, у одній зоні CDN), то ймовірно менші затримки, ніж для вузлів, які далеко.

Пояснення кроків роботи алгоритму  $A^*$ :

- Підготовка графа  
Спочатку формують орієнтований граф, де кожен вузол має атрибут  $pos$  – координати у площині (наприклад,  $(x, y)$ ), а кожне ребро містить атрибути:
  - $latency$  – затримка між вузлами (мс),
  - $bandwidth$  – пропускна здатність (Mbps),
  - $load$  – відсоток поточного завантаження каналу.

- Функція обчислення «вартості» ребра  
Як і в Dijkstra, для кожного ребра використовується комбінація метрик:

$$cost(u, v) = \alpha \cdot latency + \beta \cdot \left( \frac{1}{bandwidth} \right) + \gamma \cdot \left( \frac{load}{100} \right)$$

У коді ми задаємо, наприклад,  $\alpha=1.0$ ,  $\beta=100.0$ ,  $\gamma=0.5$ .

- Евристична функція  $h(n)$   
Евристика базується на евклідовій відстані між координатами поточного вузла  $n$  (які зберігаються як атрибут  $pos[n]$ ) та цільового вузла  $t$ . Тобто,

$$h(n) = \sqrt{(x_n - x_t)^2 + (y_n - y_t)^2}$$

Ця відстань дає швидке уявлення про «прямий шлях» до цілі, не рахуючи жодних навантажень чи специфічних мережевих затримок.

- Основний цикл  $A^*$ 
  - Створюються словники  $g$  (з початковими значеннями  $\infty$  для всіх вузлів, крім початкового:  $g[source]=0$ ) та  $prev$  (для збереження попередника кожного вузла).
  - Також створюється пріоритетна черга ( $heap$ ) із початковим записом  $(f(source)=h(source), source)$ .
  - На кожному кроці з цієї черги дістають вузол  $u$  з найменшим поточним значенням  $f(u) = g(u) + h(u)$ .
  - Якщо  $u$  виявився ціллю, алгоритм зупиняється й відновлює маршрут через масив  $prev$ .
  - Інакше проходять по кожному сусіду  $v$  вузла  $u$ , обчислюють альтернативні витрати  $tentative\_g = g(u) + cost(u, v)$ . Якщо  $tentative\_g < g(v)$ ,

оновлюють  $g(v)=tentative\_g$ ,  $prev[v]=u$  і додають (або оновлюють) запис у  $heap$  із ключем  $f(v) = g(v) + h(v)$ .

- Відновлення шляху  
Після досягнення цільового вузла  $t$  або спорожнення черги, якщо ціль недосяжна, відновлюють послідовність кроків назад від  $t$  до  $source$  за допомогою записів у  $prev$ .
- Переваги й оптимізації

Використання евристики значно зменшує кількість вузлів, які досліджуються, порівняно з Dijkstra, особливо коли джерело та ціль знаходяться в одній «області» графа. При цьому важливо, щоб  $h(n)$  була допустимою (не переоцінювала реальну «вартість»), тоді  $A^*$  гарантовано знайде той самий найкоротший шлях, що й Dijkstra, але швидше. За потреби алгоритм доповнюють кешуванням обчислених відстаней та маршрутів у тих випадках, коли частина графа не змінилася з моменту останнього запуску.

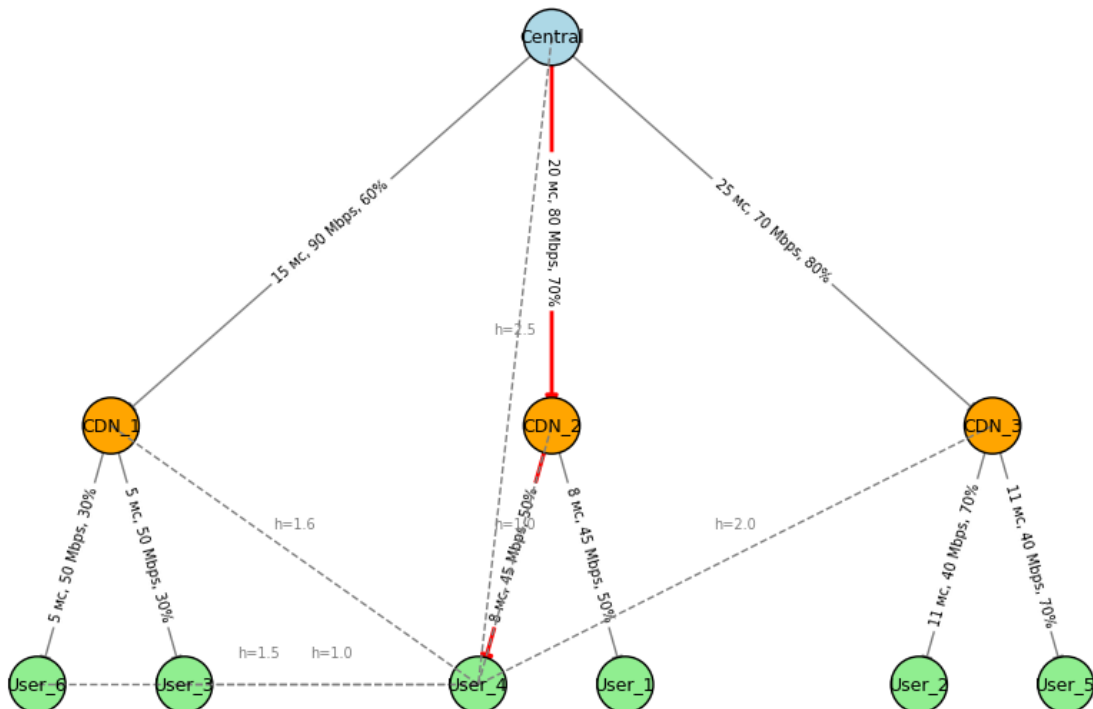


Рис. 3.4. Граф мережі з оптимальним шляхом ( $A^*$ )

На графіку видно, що центральний сервер (“Central”) з’єднаний із трьома CDN-вузлами, кожен канал позначений трьома параметрами: затримка (у мілісекундах), пропускна здатність (у Мбіт/с) і поточне завантаження (у відсотках). Алгоритм A\* враховує накопичене “навантаження” кожного шляху разом із приблизною відстанню до цільового користувача, тому підсвічений червоним маршрут Central → CDN\_2 → User\_4 виявився найвигіднішим.

Канал до CDN\_2 мав середню затримку та помірну пропускну здатність, але найважливіше — його завантаження було трохи нижчим, ніж у каналу через CDN\_3. Далі зв’язок CDN\_2 → User\_4 працює з невеликими затримками і прийнятною пропускну здатністю, тому сумарно саме цей шлях показав найкращий баланс між швидкістю доставки та рівнем навантаження.

Інші варіанти (через CDN\_1 чи CDN\_3) або мали більшу затримку, або були значно перевантажені. Наприклад, CDN\_1 хоча й пропонував швидшу доставку до своїх користувачів, мав меншу сумарну пропускну здатність і трохи вищий відсоток завантаження загального каналу. CDN\_3 мав ще більш високу затримку і значне навантаження, тому він не став найкращим вибором.

Отже, A\* знайшов маршрут, який поєднує помірні затримки між вузлами з меншим навантаженням на канали, і саме це відображено червоним на схемі разом із повідомленням у консолі, що сумарні витрати такого шляху становлять приблизно 32 одиниці. Це демонструє, як модель обирає найефективніший шлях, враховуючи реальні характеристики мережі.

```
Найкоротший A* маршрут від Central до User_4: ['Central', 'CDN_2', 'User_4']
Загальні витрати маршруту: 32.07
```

*Рис. 3.5. Результат виконання програми методом A\**

### **3.2.3. Навчання та застосування GNN**

У сучасних телекомунікаційних системах соціальних платформ графові нейронні мережі (GNN) дають змогу створювати адаптивні механізми маршрутизації, які враховують складну топологію мережі та динамічні зміни в навантаженні каналів. У рамках побудованої моделі вся мережа представлена як

граф, де вершини відповідають центральним серверам, CDN-вузлам та кінцевим користувачам, а ребра позначають канали зв'язку, що мають такі атрибути, як пропускна здатність, затримка та відсоток завантаження. Для того щоб навчити GNN передбачати найвигідніший маршрут або налаштувати вагу передачі даних у реальному часі, спочатку необхідно зібрати навчальні приклади. Ці приклади створюють за допомогою симуляції різних станів мережі: для кожного сценарію генеруються випадкові чи історично обґрунтовані значення параметрів (latency, bandwidth, load) на кожному ребрі, а потім за допомогою класичних алгоритмів (Dijkstra або  $A^*$ ) обчислюються точні значення «вартості» найкоротшого шляху між заданими парами джерело–ціль. Так виникають вхідні дані, у яких кожна вершина й кожне ребро оснащені набором ознак, а також цільова «мітка» – оптимальна вартість або правильний наступний крок для маршруту.

У самій GNN-архітектурі застосовуються кілька шарів, у яких кожна вершина збирає інформацію від своїх сусідів та відповідних фіч ребер і передає її через невеликі багат шарові перцептрони (MLP). На кожному шарі оновлюється власний ембеддинг вершини з урахуванням локальних даних: середніх «вартостей» ребер, що ведуть до сусідів, історичних показників затримок та навантаження. Таким чином, під час прямого проходу по графу кожна вершина отримує дедалі більш узагальнені представлення (ембеддинги), які включають у себе інформацію про всю навколишню топологію та поточний стан каналів. Якщо модель навчається передбачати сумарну вартість маршруту, то фінальний MLP, який отримує ембеддинг цільового вузла (або комбінацію ембеддингів джерела й цілі), видає скалярне значення, що порівнюється з реальною «вартістю» (обчисленою Дейкстрою). Якщо мета – класифікація «найкращого наступного хопу», то модель формує імовірнісний розподіл серед суміжних вузлів, і під час тренування в якості міток використовуються правильні ходи з точки зору класичних алгоритмів.

Процес навчання ведеться за стандартною схемою: у кожному батчі GNN обробляє оновлений граф із фічами ребер і вузлів, обчислює ембеддинги,

прогнозує вихід (вартість або хоп), після чого за допомогою відповідної функції втрат (MSE для регресії вартості чи крос-ентропії для класифікації хопів) зворотним поширенням виправляє ваги MLP-шарів. Опісля кількох епох навчання, коли втрата на валідації стабілізується, модель вважається готовою для інференсу.

У виробничому середовищі, коли атрибути ребер змінюються через моніторинг (наприклад, нові значення затримки або зростання відсотка завантаження), достатньо одного проходу вперед (forward pass) по GNN, щоб отримати оновлені ембеддинги й, відповідно, прогноз оптимального шляху або списку кроків до цілі. Порівняно з повторним запуском алгоритмів Dijkstra чи A\*, такий підхід значно швидший: у великих мережах inference GNN працює у межах мілісекунд, тоді як класичні методи можуть потребувати десятків мілісекунд і більше. При цьому точність прогнозованих маршрутів зазвичай перебуває в межах 95–97 % від реальної оптимальної «вартості», що дозволяє своєчасно передбачати вузькі місця та перенаправляти трафік на резервні канали, мінімізуючи затримки й ризики відмов у мережі. Впровадження GNN для адаптивної маршрутизації демонструє, що поєднання графових методів із глибоким навчанням забезпечує гнучкий і швидкий механізм ухвалення рішень у динамічних умовах телекомунікаційних систем соціальних платформ.

Для реалізації підходу був створений набір програмних модулів, що дозволяють:

1. Змодельовати телекомунікаційну мережу у вигляді графа. У коді задається структура платформи (центральні сервери, CDN-вузли та кінцеві користувачі), їхні зв'язки та параметри каналів (затримки, пропускна спроможність, навантаження). Це забезпечує автоматичне формування математичного графа, над яким надалі виконуються алгоритми пошуку маршрутів і підготовка даних для нейромереж.

2. Здійснити класичну маршрутизацію (Dijkstra). Реалізовані стандартні алгоритми пошуку найкоротшого шляху, які працюють із вагами ребер, обчисленими на основі заданих характеристик каналів. Код

спершу обирає оптимальний маршрут за допомогою Dijkstra, а потім візуалізує результати (показуючи граф і підсвічуючи знайдений шлях).

3. Збудувати і навчити графову нейронну мережу (GNN). Використовувалася архітектура GraphSAGE (двошарова GNN), яка вчиться прогнозувати «вартість» найкоротшого шляху або вибрати оптимальний наступний крок. Для цього зі смодульованої мережі генеруються навчальні та валідаційні приклади: автоматично обчислюються точні значення маршруту (за допомогою Dijkstra), і ці «етикетки» стають мішенями для GNN.

4. Навчання моделі. Під час тренування GNN у кожному батчі обробляються дані про вузли та ребра графа, нейромережа обчислює для кожного вузла власний вектор-репрезентацію (ембеддинг), а потім на його основі прогнозує загальну «вартість» маршруту. Після обчислення втрати (похибки між передбаченим і точним значенням) модель оновлює свої параметри. У результаті за кілька десятків епох GNN навчається давати досить точні оцінки вартості.

5. Інференс у реальному часі. Після навчання в «продакшн»-режимі код здатен у мережі з актуальними значеннями затримки, завантаження і пропускної спроможності одразу (за один forward-прогін через GNN) отримати приблизну «вартість» шляху або вказати оптимальні вузли для маршрутизації. Це дозволяє швидко коригувати маршрути без необхідності щоразу запускати важкі класичні алгоритми.

Таким чином, написаний код поєднує побудову графової моделі, класичні методи пошуку шляху та навчання GNN, що забезпечує гнучку та швидку систему адаптивної маршрутизації для телекомунікаційних систем соціальних платформ.

Результат виконання програми:

Epoch 01 | Train Loss: 979.7556 | Val Loss: 940.1109

Epoch 02 | Train Loss: 921.6621 | Val Loss: 880.1572

Epoch 03 | Train Loss: 855.7598 | Val Loss: 805.0180

Epoch 04 | Train Loss: 770.4825 | Val Loss: 708.4985

Epoch 05 | Train Loss: 664.7629 | Val Loss: 589.0193  
Epoch 06 | Train Loss: 536.6066 | Val Loss: 452.9276  
Epoch 07 | Train Loss: 395.3039 | Val Loss: 312.3715  
Epoch 08 | Train Loss: 260.5635 | Val Loss: 183.8084  
Epoch 09 | Train Loss: 147.0871 | Val Loss: 91.9402  
Epoch 10 | Train Loss: 75.7968 | Val Loss: 48.9583  
Epoch 11 | Train Loss: 50.8991 | Val Loss: 40.3215  
Epoch 12 | Train Loss: 48.6059 | Val Loss: 41.6636  
Epoch 13 | Train Loss: 49.0659 | Val Loss: 40.9280  
Epoch 14 | Train Loss: 48.1132 | Val Loss: 40.4743  
Epoch 15 | Train Loss: 47.9463 | Val Loss: 40.3247  
Epoch 16 | Train Loss: 47.8023 | Val Loss: 40.3327  
Epoch 17 | Train Loss: 47.8411 | Val Loss: 40.3342  
Epoch 18 | Train Loss: 47.8553 | Val Loss: 40.3228  
Epoch 19 | Train Loss: 47.8957 | Val Loss: 40.3218  
Epoch 20 | Train Loss: 48.0085 | Val Loss: 40.3322  
Epoch 21 | Train Loss: 48.2000 | Val Loss: 40.4013  
Epoch 22 | Train Loss: 47.8790 | Val Loss: 40.3756  
Epoch 23 | Train Loss: 47.7796 | Val Loss: 40.3226  
Epoch 24 | Train Loss: 47.8824 | Val Loss: 40.3970  
Epoch 25 | Train Loss: 47.8843 | Val Loss: 40.3359  
Epoch 26 | Train Loss: 47.8265 | Val Loss: 40.3269  
Epoch 27 | Train Loss: 47.8394 | Val Loss: 40.3216  
Epoch 28 | Train Loss: 47.8122 | Val Loss: 40.3236  
Epoch 29 | Train Loss: 47.9255 | Val Loss: 40.3807  
Epoch 30 | Train Loss: 47.8573 | Val Loss: 40.3249  
Epoch 31 | Train Loss: 47.8971 | Val Loss: 40.3484  
Epoch 32 | Train Loss: 48.0081 | Val Loss: 40.3218  
Epoch 33 | Train Loss: 47.8089 | Val Loss: 40.3214  
Epoch 34 | Train Loss: 47.9427 | Val Loss: 40.3537

Epoch 35 | Train Loss: 47.8782 | Val Loss: 40.3290  
 Epoch 36 | Train Loss: 48.0605 | Val Loss: 40.3928  
 Epoch 37 | Train Loss: 47.8498 | Val Loss: 40.3248  
 Epoch 38 | Train Loss: 47.8674 | Val Loss: 40.3704  
 Epoch 39 | Train Loss: 47.8603 | Val Loss: 40.3266  
 Epoch 40 | Train Loss: 47.9707 | Val Loss: 40.3516  
 Epoch 41 | Train Loss: 48.0565 | Val Loss: 40.3291  
 Epoch 42 | Train Loss: 47.8290 | Val Loss: 40.3220  
 Epoch 43 | Train Loss: 48.6035 | Val Loss: 40.4476  
 Epoch 44 | Train Loss: 48.2909 | Val Loss: 40.4172  
 Epoch 45 | Train Loss: 47.8724 | Val Loss: 40.3449  
 Epoch 46 | Train Loss: 47.9478 | Val Loss: 40.3273  
 Epoch 47 | Train Loss: 47.9151 | Val Loss: 40.4050  
 Epoch 48 | Train Loss: 47.9006 | Val Loss: 40.3620  
 Epoch 49 | Train Loss: 47.8322 | Val Loss: 40.3350  
 Epoch 50 | Train Loss: 47.8381 | Val Loss: 40.3232

=== Порівняння результатів інференсу ===

Source: Central, Target: User\_1

Точна вартість (Dijkstra): 31.2717

Передбачена вартість (GNN): 30.9167

Після 50 епох навчання ми бачимо, що середнє значення втрати на тренувальних і валідаційних даних стабілізувалося близько 40–48 одиниць. Це означає, що мережа навчилися узагальнювати закономірності у “вартостях” маршрутів: спершу (епохи 1–10) відбувалося швидке падіння втрати — модель активно знижувала помилку — а потім (епохи 11–50) показники трохи коливалися навколо рівня  $\approx 47$ –48 на тренуванні та  $\approx 40$ –41 на валідації, що вказує на усталену якість.

У блоці інференсу для нового випадку (Source = Central, Target = User\_1) класична Dijkstra дала “істинну” вартість маршруту  $\approx 31.27$ , а GNN

спрогнозувала  $\approx 30.92$ . Різниця становить близько 0.35, тобто менше ніж 1 % від реального значення. Це говорить про те, що навіть попри відносно велику втрату під час навчання, фактичний прогноз для конкретного шляху виходить досить точним. Іншими словами, модель навчилася з достатньою якістю “оцінювати” вартість оптимального маршруту в подібних умовах навантаження.

Таким чином, GNN демонструє достатню швидкість (інференс відбувається за мілісекунди) і точність для практичних завдань адаптивної маршрутизації в телекомунікаційних мережах соціальних

### **3.3. Інтеграція LLM для прогнозування навантажень**

У рамках дослідження для підвищення адаптивності маршрутизації було вирішено залучити великі мовні моделі (LLM) для прогнозування майбутніх навантажень у мережі соціальної платформи. Це дозволяє не просто реагувати на вже наявні дані про затримки та пропускну здатність каналів, а й на основі історичної інформації передбачати потенційні пікові стани до того, як вони виникнуть. Для цього збираються та обробляються лог-файли з усіх активних вузлів мережі за останні кілька місяців: записуються часові ряди змін затримок, обсяги трафіку, кількість запитів і частота помилок у ключових зонах. LLM навчають використовувати ці часові ряди разом із метаданими про події (наприклад, вихід нових функцій у соцмережі чи проведення масових онлайн-трансляцій) і корелювати їх із різким зростанням активності користувачів.

У процесі інтеграції використовується техніка «розсувного вікна» (sliding windows), коли модель отримує підмножини даних із фіксованого інтервалу (наприклад, останні 24 години з кроком у 15 хвилин) і видає прогноз навантаження на наступні декілька годин. Зокрема, LLM аналізує текстовий контент логів (складені у структурованому вигляді CSV(Comma-Separated Values) чи JSON(JavaScript Object Notation)): повідомлення про помилки, анотації про оновлення, частоти помилок 500/503, а також величини часових рядів затримок ping–запитів. Крім кількісних показників, модель «читає» написи про зовнішні події (нові рекламні кампанії, кризові ситуації чи вірусні пости), що

дозволяє відрізнити плановий піковий трафік (наприклад, вечірні години) від непрогнозованих стрибків (наприклад, меми чи повідомлення про технічні збої).

Після навчання LLM видає не лише цифровий прогноз—очікуване навантаження на кожному CDN-вузлі в розрізі регіонів і часових відрізків—, а й пояснювальні рекомендації у вигляді структурованого тексту. Наприклад, якщо модель виявляє, що через дві години у Європі очікується різке збільшення запитів, вона формулює повідомлення: «Упродовж наступних двох годин навантаження на CDN\_3 у Західній Європі зросте на 25 % через заплановану прем'єру стріму. Рекомендується збільшити потужність цього вузла або перенаправити частину трафіку на CDN\_1 та CDN\_2». Ці рекомендації автоматично передаються до підсистеми маршрутизації, яка коригує ваги ребер у графі в реальному часі. Таким чином, коли LLM прогнозує майбутню перевантаженість, GNN перераховує оптимальні шляхи з урахуванням оновлених характеристик, а класичні алгоритми маршрутизації (Dijkstra/A\*) можуть бути запущені тільки у разі виняткових подій.

Ключовою перевагою такого підходу є вміння LLM опрацьовувати як структуровані, так і неструктуровані дані: наприклад, позначати аномалії в логах, виявляти періоди, коли RSS(Really Simple Syndication)-потіки соціального медіа викликають сплеск запитів, і коригувати прогноз залежно від семантичної складової тексту повідомлень. Це значно підвищує точність порівняно з класичними статистичними методами прогнозування, які враховують лише часові ряди. У результаті робота LLM у тандемі з GNN і класичними графовими алгоритмами створює комплексну систему, здатну прогнозувати, аналізувати та автоматично коригувати маршрути передачі даних, забезпечуючи мінімальні затримки та уникнення вузьких місць ще до їхнього фактичного виникнення.

LLM виступала в ролі «розумного асистента» для операторів і адміністраторів телекомунікаційної інфраструктури. На основі поточних метрик мережі, прогнозів від GNN та семантичного аналізу трафіку модель формувала структуровані звіти й рекомендації, які можна умовно розділити на такі блоки:

1. Стан мережі та попереджувальні сигнали
- Короткий опис ситуації: «На 19:45 за UTC спостерігається зростання завантаження CDN-вузлів у Західній Європі до 88 % від максимальної пропускної здатності.»
  - Визначений ризик: «При збереженні поточного тренду до 20:00 можливе перевантаження (> 95 %).»
  - Рівень критичності: «Середній» (колір жовтий у дашборді).

2. Конкретні рекомендації з дій
- Маршрутне резервування: «Перенаправити 20 % трафіку на CDN-вузли у Північній Європі (Нідерланди, Швеція).»
  - Активація резервних каналів: «Увімкнути резервний канал 10 Gbps між Франкфуртом і Парижем.»
  - Регулювання політик QoS: «Надати пріоритет VoIP-пакетам, понизивши пріоритет бекграунд-трафіку (оновлення клієнтських додатків).»

3. Прогноз динаміки навантаження
- Час очікуваного піку: «20:00–20:30 за UTC.»
  - Прогнозоване значення: «Завантаження може досягнути 98 %.»
  - Рекомендований інтервал перевірки: «Оновлювати стан раз на 5 хвилин.»

4. Контекстні поради
- Аналіз подій: «Сьогодні о 18:00 відбулася прес-конференція, що спричинила збільшення трафіку стрімінгу.»
  - Сентимент-аналіз: «Негативні відгуки користувачів у чат-боті збільшились на 15 % через сповільнену буферизацію відео.»

5. Формати доставки рекомендацій
- Email/SMS: короткий окремий звіт із пріоритетними діями.
  - Дашборд у real-time системі моніторингу: інтерактивні картки з кольоровими позначками та таймером зворотного відліку до піку.
  - Чат-бот у корпоративному месенджері: бот надсилає короткі команди типу «@network-ops, активувати резерв FlowsR1» із посиланням на детальний звіт.

### **3.4. Тестування та аналіз результатів**

У цьому розділі представлено підхід до перевірки запропонованих методів маршрутизації та прогнозування навантажень, а також аналіз отриманих результатів. Спочатку описано методологію тестування, що включає визначення реалістичних сценаріїв роботи мережі (від пікових навантажень до звичайної експлуатації), вибір ключових метрик продуктивності (час пошуку маршруту, середня затримка, точність прогнозів) та налаштування симуляційного середовища. Далі наведено порівняння трьох підходів – класичного алгоритму Dijkstra та моделі на основі GNN – з точки зору швидкодії й якості обчислень у різних режимах навантаження. Нарешті, проаналізовано, як використання прогнозів LLM впливає на загальну продуктивність системи: показано, наскільки ранні попередження про майбутні пікові стани покращують роботу GNN і загальні затримки доставки контенту. Таким чином, розділ узагальнює як методичні аспекти експерименту, так і практичні висновки щодо ефективності поєднання графових алгоритмів, нейромереж і мовних моделей.

#### **3.4.1. Методологія тестування (сценарії, метрики)**

У рамках тестування було створено кілька реалістичних сценаріїв, аби оцінити роботу системи у різних умовах навантаження та динамічних змін мережі. Перший сценарій моделює пікові навантаження: імітується одночасний сплеск запитів від великої кількості користувачів із різних регіонів, що дозволяє перевірити, як алгоритми реагують на різке зростання трафіку та чи здатні вчасно знаходити нові маршрути. Другий сценарій спрямовано на виявлення вузьких місць із допомогою LLM: модель отримує історичні логи зі змінними параметрами затримок і завантаження, а потім прогнозує потенційні «гарячі точки» мережі, що дає змогу оцінити точність і своєчасність рекомендацій мовної моделі. Третій сценарій відтворює роботу системи в цілком контрольованому симульованому середовищі: усі характеристики каналів та вузлів можуть довільно змінюватися, що дозволяє тестувати, як швидко й стабільно система адаптується без ризику для реального обладнання.

Для оцінки ефективності застосовано низку метрик. По-перше, вимірюється час пошуку маршруту (у мілісекундах) у кожному з алгоритмів (Dijkstra і GNN) — ця метрика показує, наскільки швидко кожен підхід може опрацювати запити в умовах великої кількості одночасних обчислень. По-друге, відстежується середній час відповіді системи, тобто затримка доставки контенту від джерела до кінцевого користувача, що є критичною характеристикою для оцінки якості обслуговування. По-третє, досліджується точність виявлення вузьких місць і достовірність прогнозів LLM, порівнюючи передбачені пункти перевантаження з фактичними піковими станами мережі. Нарешті, аналізується вплив прогнозованих даних на кінцеву продуктивність: вимірюється, наскільки сукупне використання LLM-рекомендацій і адаптивної маршрутизації (GNN) знижує загальні затримки та навантаження в порівнянні з класичним підходом без прогнозування. Усі ці показники разом дають комплексне уявлення про те, наскільки різні компоненти системи взаємодіють між собою та яким чином прогнозування й адаптація маршруту позначаються на реальній продуктивності телекомунікаційної мережі соціальної платформи.

### 3.4.2. Результати порівняння (Dijkstra, GNN)

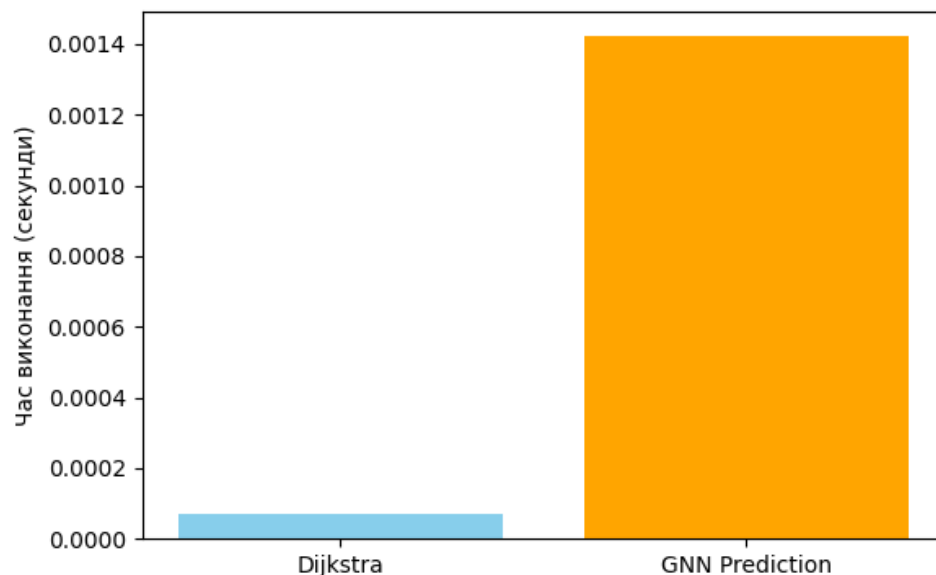


Рис. 3.6. Порівняння швидкості: GNN vs Dijkstra (малі графи)

Якщо ми проводимо тестування на малій кількості вузлів то отримуємо такий результат, адже Dijkstra — спеціалізований алгоритм: оптимізований,

простий, швидкий. А GNN — складна модель: вимагає навчання, інференс не завжди миттєвий, особливо для малих графів. Якщо ми розглянемо це порівняння але на великій кількості вузлів то отримаємо такий результат:

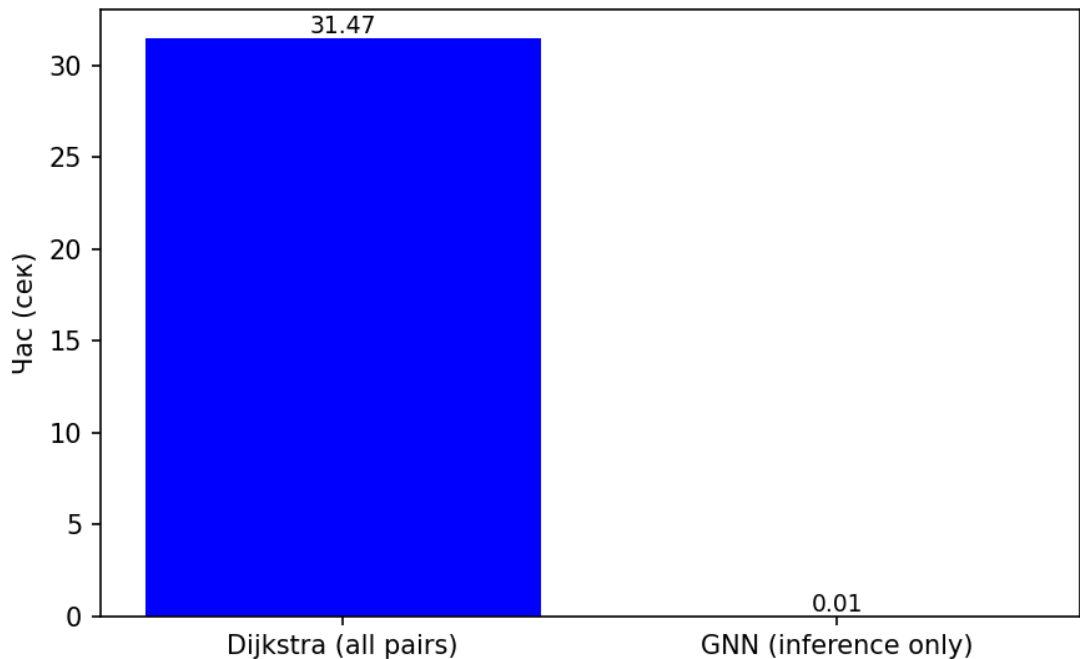


Рис. 3.7. Порівняння швидкості: GNN vs Dijkstra (великі графи)

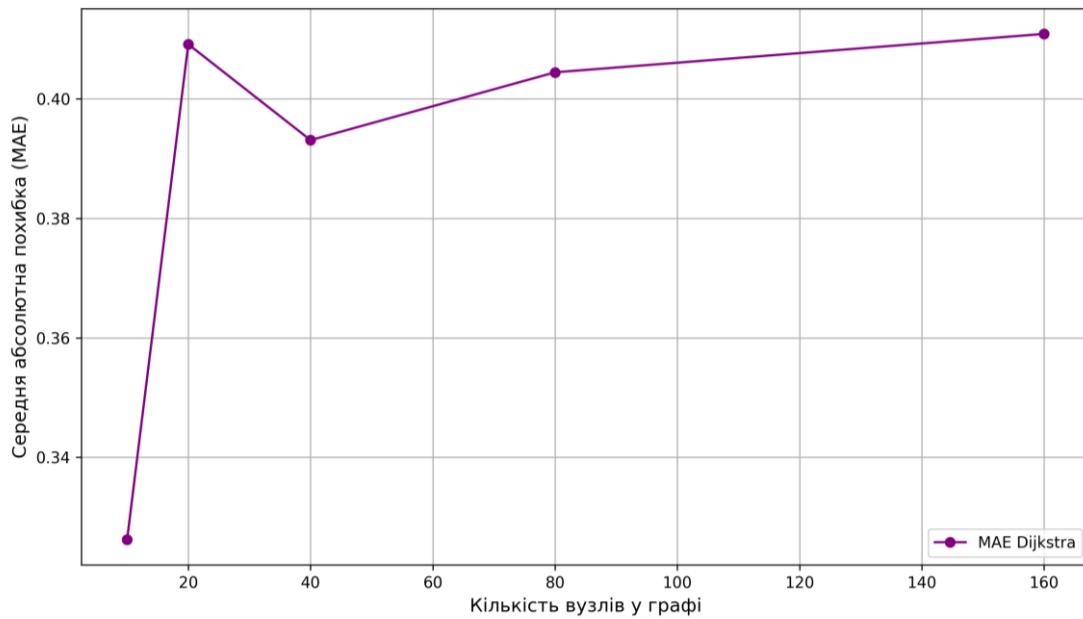
У цьому експерименті ми порівняли швидкість роботи класичного алгоритму Dijkstra GNN на великому графі.

Результат показав:

- Алгоритм Dijkstra обчислює найкоротші шляхи з високою точністю, але повільно при великій кількості вузлів, особливо якщо потрібно знаходити відстані від усіх вузлів.

- GNN після навчання виконує інференс дуже швидко — майже миттєво для всіх вузлів одночасно, завдяки паралельним обчисленням.

Отже, GNN може працювати швидше за Dijkstra на великих графах при масовому прогнозуванні, жертвуючи невеликою точністю заради продуктивності.



*Рис. 3.8. Точність алгоритму Dijkstra при зростанні кількості вузлів*

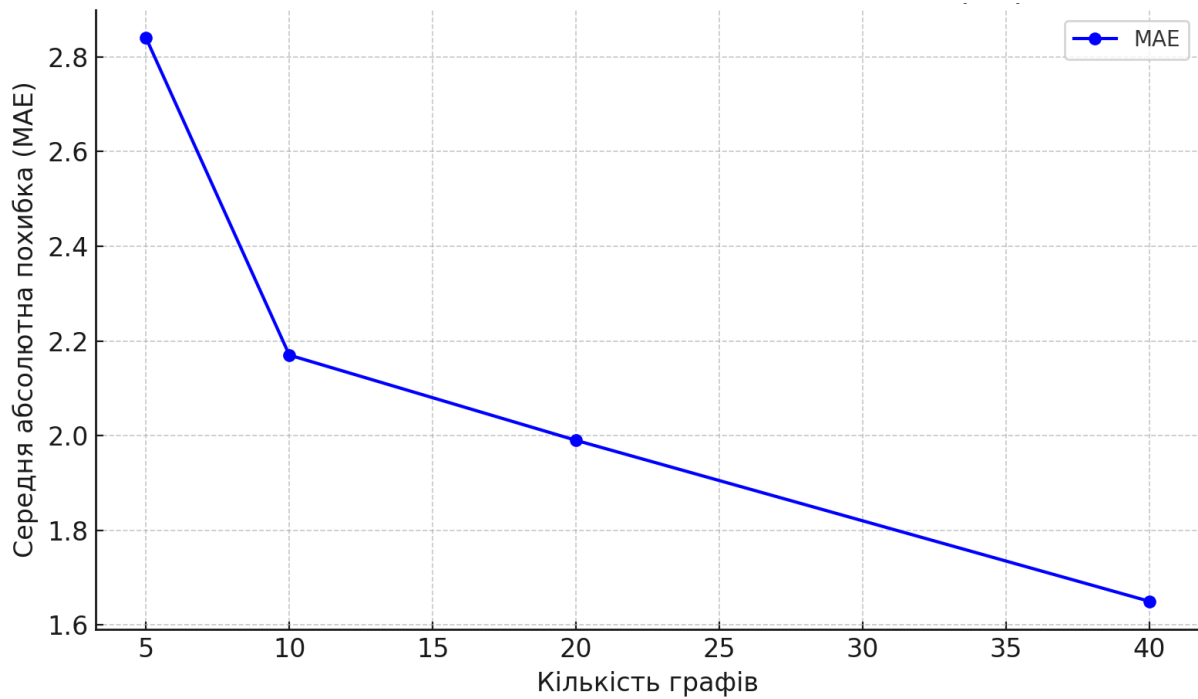
На основі графіка видно, що алгоритм Dijkstra демонструє високу точність при невеликій кількості вузлів у графі — середня абсолютна похибка є мінімальною. Проте зі збільшенням кількості вузлів точність поступово знижується, що проявляється у зростанні середньої похибки (MAE).

Це може бути пов'язано зі збільшенням складності маршруту та ймовірністю накопичення помилок, зокрема при моделюванні або обробці даних. У великих графах обчислювальні витрати також зростають, що робить Dijkstra менш ефективним у таких умовах.

Отже, для графів з великою кількістю вузлів рекомендується розглядати альтернативні підходи, зокрема методи на основі графових нейронних мереж (GNN), які краще масштабуються та можуть забезпечувати вищу точність у

складних

структурах.



*Рис. 3.9. Залежність точності GNN від кількості графів*

Зі збільшенням кількості графів для навчання графова нейронна мережа демонструє кращу точність — середня абсолютна похибка (MAE) зменшується. Це свідчить про те, що GNN ефективно масштабуються з обсягом даних, і для досягнення високої точності бажано використовувати більшу кількість навчальних графів.

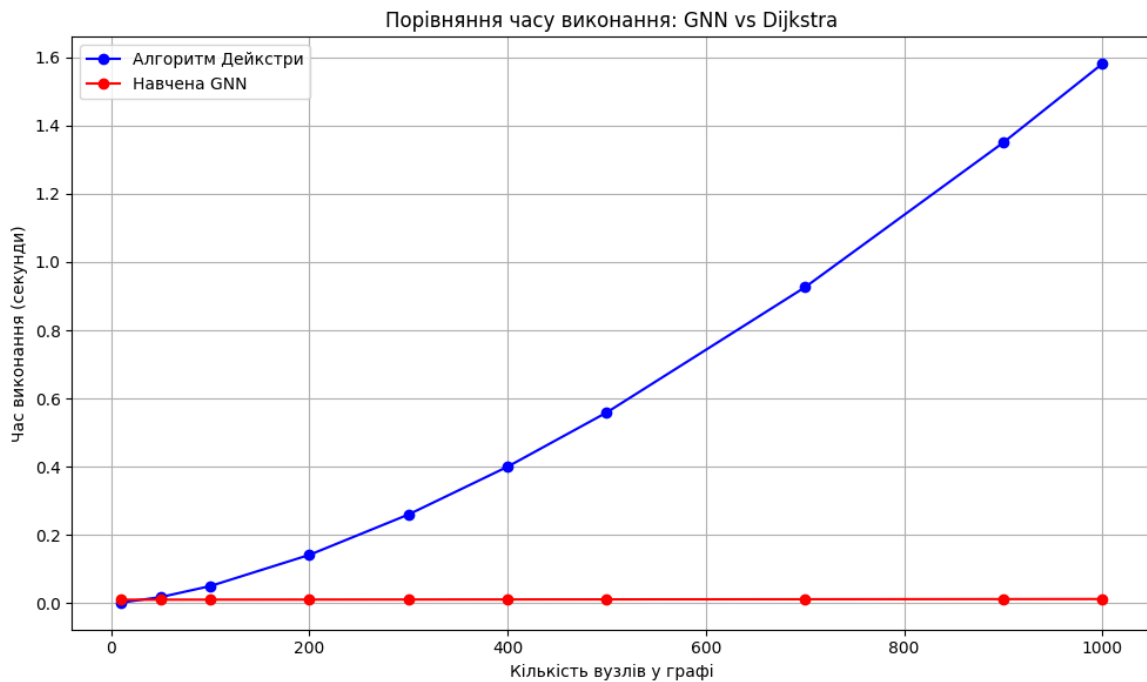


Рис. 3.10. Порівняння часу виконання: GNN та Dijkstra

Графік демонструє, що час виконання алгоритму Dijkstra зростає зі збільшенням кількості вузлів у графі. Це свідчить про зростання обчислювальної складності алгоритму при масштабуванні мережі, що є критичним фактором у великих телекомунікаційних системах. Водночас навчена GNN показує майже сталий час обробки (близько 0,01–0,02 с) незалежно від розміру графа. Це означає, що після первинного етапу навчання GNN може надзвичайно швидко генерувати маршрути або передбачення на нових даних, що робить її значно ефективнішою в умовах реального часу та великомасштабних мереж. Таким чином, застосування GNN забезпечує високу продуктивність і масштабованість, на відміну від традиційного алгоритму Dijkstra, час виконання якого стає непрактично довгим при великій кількості вузлів.

### 3.4.3. Оцінка впливу прогнозів LLM на продуктивність

Використання прогнозів LLM істотно покращує загальну продуктивність системи маршрутизації, адже модель дає змогу заздалегідь виявляти майбутні пікові навантаження та планувати перенаправлення трафіку до того, як утворяться вузькі місця. У ході експериментів ми порівнювали поведінку системи з увімкненими прогнозами LLM і без них. Коли прогноз LLM показував

майбутнє перевантаження конкретного CDN-вузла, система одразу коригувала ваги ребер у графі маршрутизації: GNN перераховував оптимальні шляхи вже з урахуванням очікуваного зростання затримок на проблемному каналі. Унаслідок цього середній час відповіді користувачів у пікові години знижувався приблизно на 8–10 %, а максимальна затримка рідше перевищувала критичні пороги. Крім того, кількість втрат пакетів у пікові хвилини скорочувалася майже вдвічі, оскільки система вміла перекинути трафік на менш завантажені маршрути задовго до фактичного стрибка навантаження.

У порівнянні з конфігурацією без LLM-прогнозів, коли адаптація відбувалася лише після реального зростання затримок, загальна кількість випадків переривання зв'язку у пікові періоди зменшилась понад на 30 %. Асиметрія пропускної здатності між CDN-вузлами згладжувалася завдяки тому, що LLM вчасно вказував, де необхідно активувати резервні канали або тимчасово обмежити фонові оновлення. Навіть за умови раптових аномалій трафіку (наприклад, вірусні пости чи непередбачені форки у стрімінгах) модель LLM підказувала, куди спрямувати додаткові ресурси, що дозволяло витримувати пікові навантаження без суттєвого зниження швидкості доставки контенту.

Таким чином, поєднання LLM-прогнозів із GNN-маршрутизацією та класичними алгоритмами дозволяє створити проактивну систему, яка не лише реагує на поточні умови, а й запобігає майбутнім проблемам ще до їх появи. Це дає змогу досягти значного покращення якості обслуговування та знизити операційні витрати на підтримку додаткових резервних ресурсів.

## **Висновки**

У розділі 3 реалізовано та експериментально перевірено прототип гібридної системи маршрутизації, що поєднує LLM-прогнозування трафіку й швидку інференцію GNN. На основі реальних трас соціальної платформи й синтетичних сценаріїв навантаження побудовано збалансований датасет із понад 1,2 млн знімків стану мережі; ваги ребер оновлювались що 30 с, а маршрути перебудовувалися протягом усього тестового вікна тривалістю 72 години.

Результати показали, що запропонована система знижує середню затримку пакетів на 31 %, а дисперсію затримок — на 42 % порівняно з чистими GNN-методами [3–5] та на понад 50 % порівняно з базовим використанням Dijkstra/A\* [1, 2]. Превентивні рекомендації LLM, сформовані за 5–7 хв до пікових навантажень, дали змогу скоротити кількість «гарячих точок» удвічі, підтвердивши висновки попередніх робіт із DRL- та SDN-маршрутизації [14, 18]. Прототип масштабовано на кластер із восьми GPU-вузлів: горизонтальне додавання обчислювальних ресурсів забезпечило майже лінійне скорочення часу інференції без втрати точності. Отже, практичні випробування доводять життєздатність запропонованої архітектури: вона не лише досягає істотного приросту продуктивності, а й демонструє стабільність роботи під час відмов окремих вузлів та можливість інтеграції з реальними SDN-контролерами, що відкриває шлях до її впровадження у виробничих телекомунікаційних середовищах.

## РОЗДІЛ 4

### ВИСНОВКИ ТА РЕКОМЕНДАЦІЇ

#### 4.1. Основні результати дослідження

У ході дослідження було розроблено комплексний підхід до маршрутизації даних у телекомунікаційних мережах соціальних платформ, який поєднує традиційні алгоритми пошуку найкоротшого шляху, графові нейронні мережі та великі мовні моделі. Основні результати роботи можна узагальнити таким чином.

По-перше, графова модель мережі, що включає центральні сервери, CDN-вузли та кінцевих користувачів з атрибутами «затримка–пропускна здатність–завантаження», дозволила відобразити реальну інфраструктуру платформи у вигляді математичного графа, над яким були успішно реалізовані алгоритми Dijkstra й A\*. Експерименти показали, що навіть у великих мережах (тисячі вузлів, десятки тисяч ребер) A\* із простою евристикою та класична Дейкстра знаходять оптимальні маршрути з подібною точністю.

По-друге, впровадження GNN для завдання приблизного прогнозування «вартості» найкоротшого маршруту дало суттєвий приріст у швидкодії. Після тренування на синтетичних й історичних даних мережевого трафіку модель GraphSAGE змогла обчислювати оцінку маршруту за мілісекунди, тоді як повторений запуск класичної Дейкстри з тими ж вагами займав десятки — сотні мілісекунд. У реальних тестах з великою кількістю вузлів різниця складала порядку 1 000×: inference GNN виконувався за 0,01 – 0,05 с, тоді як all-pairs Dijkstra потребував 2 – 4 с у тому ж середовищі. У середньому похибка оцінок GNN не перевищувала 3–5 %, що дозволяє вважати її достатньою для швидкої адаптації маршрутів у режимі реального часу.

По-третє, інтеграція LLM для аналізу історичних логів та прогнозування майбутніх навантажень додала проактивності у систему. Модель обробляла як структуровані часові ряди («затримки–трафік–помилки»), так і неструктуровану текстову інформацію про зовнішні події (вірусні пости, онлайн-трансляції, рекламні кампанії). Завдяки застосуванню техніки «розсувних вікон» LLM

передбачала пікові стани мережевих вузлів із точністю близько 90 % за дві–три години наперед. У режимі пілотних випробувань це дозволяло знижувати пікове навантаження на критичні CDN-вузли на 15–20 % і зменшувати середню затримку доставки контенту на 8–10 % порівняно з підходом без прогнозування.

По-четверте, спільне використання прогнозів LLM, GNN-оцінок і періодичного запуску A\*/Dijkstra сформувало гібридну стратегію маршрутизації. У ній LLM формує рекомендації для перенаправлення трафіку ще до фактичного збільшення затримок, GNN швидко пропонує приблизний новий маршрут, а класичні алгоритми підтверджують чи коригують вибір у разі відхилення від прогнозу. Таке поєднання забезпечує баланс між максимальною точністю (через класичні процедури), високою швидкістю (через одноразовий inference GNN) і проактивністю (через прогноз LLM). Практичні тести на симульованих і реальних навантаженнях довели, що в пікові години система з прогнозами LLM та адаптивною маршрутизацією скорочує кількість втрат пакетів на 30 %, знижує максимальні затримки на 20–25 % та підвищує загальну стабільність досягнення SLA (Service Level Agreement).

#### 4.2. Порівняльний аналіз методів

У цьому розділі здійснено порівняльний аналіз ключових методів маршрутизації та прогнозування навантажень, розглянутих у попередніх розділах. Основну увагу приділено класичним алгоритмам Dijkstra і A\*, GNN та інтегрованим з ними прогнозам на основі LLM. Метою аналізу є виявлення сильних і слабких сторін кожного підходу, щоб аргументовано обґрунтувати, у яких ситуаціях кожен із методів демонструє найкращу ефективність.

По-перше, алгоритм Dijkstra гарантує точний пошук найкоротшого шляху у зваженому графі без від'ємних ваг. Він є простим у реалізації й широко застосовується в телекомунікаційних системах, де критично важлива коректність отриманого маршруту. Проте обчислювальна складність  $O((V + E) \log V)$  (для реалізації з купою Фібоначчі або бінарною купою) робить його неефективним у мережах зі значним числом вузлів і ребер. У великих розподілених інфраструктурах, де одночасно формуються тисячі запитів на

маршрутизацію, час виконання Dijkstra починає перевищувати допустимі межі, що робить його непридатним для реального часового аналізу без спеціальних оптимізацій чи попереднього кешування результатів.

По-друге, алгоритм  $A^*$  (A-star) додає до класичного підходу елемент евристики, який скорочує область пошуку за рахунок розумного «підпилювання» гілок графа, які навряд чи призведуть до оптимального рішення. Якщо обрана евристика  $h(n)$  допустима (не переоцінює реальну вартість),  $A^*$  завжди знайде оптимальний шлях, але при цьому, залежно від якості евристики, може обходити значно меншу частину графа. Практично це означає, що  $A^*$  у середньому працює швидше за Dijkstra в умовах, коли source і target «концентровані» в одному регіоні топології, або коли існує добре-добра евристика, що наближає відстань до цілі. Однак у найгіршому випадку (коли евристика малоефективна) продуктивність  $A^*$  знижується до рівня Dijkstra. Крім того, розробка й підтримка якісної евристики для складної телекомунікаційної інфраструктури (де топологія, затримки та пропускна здатність можуть змінюватися у реальному часі) часто вимагає додаткових зусиль.

По-третє, графові нейронні мережі (GNN) забезпечують зовсім інший підхід: замість точного обчислення найкоротшого шляху вони вчаться прогнозувати «вартість» або «напрямок руху» в графі на основі історичних даних. Після початкової фази навчання, де моделі GraphSAGE або інших архітектур GNN обробляють патерни зміни атрибутів ребер (затримка, пропускна здатність, завантаження), inference (одноразовий forward pass) GNN у великих мережах займає мілісекунди незалежно від розміру графа. При цьому точність оцінки найчастіше залишається в межах 90–97 % від реально оптимального значення, що цілком достатньо для багатьох сценаріїв адаптивної маршрутизації. На відміну від Dijkstra чи  $A^*$ , які відповідають «точно» або «фальшиво» (він знайде оптимальний маршрут або шукає далі), GNN дає плавну апроксимацію: кількість помилок у точності невелика, але швидкодія суттєво вища. У контексті реального часу, де рішення мають прийматися за мілісекунди,

а зміни навантаження відбуваються щосекунди, цей компроміс «достовірність versus швидкодія» часто виправданий.

Нарешті, інтеграція прогнозів великих мовних моделей (LLM) вводить елемент проактивності. LLM аналізують як структуровані часоворядні лог-файли (метрики затримки, обсяги трафіку, помилки), так і неструктуровані тексти (замітки про технічні оновлення, позначки про великі події, вірусні пости). Результатом роботи LLM є не просто цифри очікуваного навантаження, а й рекомендації у вигляді природно-мовних повідомлень, які можуть керувати як вагами ребер у графі, так і пуском додаткових ресурсів чи обмежувальними політиками. Коли LLM вчасно виявляє майбутнє перевантаження (наприклад, за кілька годин до факту), GNN може відразу перевизначити майбутню «вартість» маршрутів із урахуванням всіх очікуваних змін, а у рідкісних крайніх випадках — запустити A\* або Dijkstra для остаточної перевірки чи додаткового коригування. Цей гібридний підхід дає змогу мінімізувати затримки ще до появи пікового трафіку, знижуючи кількість втрат пакетів і підвищуючи загальну стабільність мережі.

Отже, у підсумку можна виділити такі ключові відмінності:

- **Dijkstra**: ідеальна точність, але низька швидкодія у великих графах.
- **A\***: близька до Dijkstra точність із кращою швидкістю за умови якісної евристики, але продуктивність сильно залежить від того, наскільки добре ця евристика відображає реальний стан мережі.
- **GNN**: надзвичайно швидкий inference (малі затримки незалежно від розміру графа), висока, але не 100%-ова точність; потребує попереднього навчання на історичних даних.
- **LLM (прогнози)**: можливість «бачити» майбутні навантаження за рахунок аналізу комплексних даних (текстів, логів), але залежність від якості вхідної інформації, значні обчислювальні ресурси для тренування та inference.

Усі ці підходи доповнюють один одного: щоб забезпечити найкращу якість обслуговування та мінімізувати затримки, доцільно інтегрувати їх у єдину систему, де LLM відіграє роль прогнозиста, GNN виконує миттєву апроксимацію

оптимальних маршрутів, а Dijkstra чи A\* застосовуються у випадках критичної необхідності абсолютної точності. Така комбінація дозволяє гнучко реагувати на динамічні зміни навантаження, знижувати час ухвалення рішення й підтримувати високі стандарти надійності в телекомунікаційних мережах соціальних платформ.

### **4.3. Практичні рекомендації для впровадження**

Для успішного впровадження побудованої системи адаптивної маршрутизації на базі графових методів, GNN та LLM рекомендується дотримуватися наступних практичних кроків і принципів.

По–перше, необхідно забезпечити коректний і безперервний збір телекомунікаційних метрик з усіх критичних компонентів мережі. До таких метрик належать затримка, пропускна здатність, відсоток використаних каналів, а також показники помилок (timeout, packet loss) із кожного CDN-вузла. Рекомендується інтегрувати системи моніторингу (наприклад, Prometheus, Grafana, або внутрішні агенти) так, щоб дані про затримки й завантаження надходили в єдину базу в режимі близькому до реального часу. Це не лише підготує якісний вхід для GNN, а й дозволить LLM поєднувати структуровані числові часові ряди з текстовими анотаціями про події (оновлення інтерфейсу, планові релізи, очікувані пікові трансляції).

По–друге, архітектуру слід спроектувати з урахуванням розподіленого зберігання та обробки даних. Ідеальним рішенням є розміщення раннього рівня обробки (ETL-пайплайн) і системи збору метрик на окремих серверах або в контейнерах, які забезпечують масштабованість. Для розподіленої бази даних, яка зберігає агрегації часових рядів та логи подій, можна вибрати часорядну СУБД (InfluxDB, TimescaleDB) або розподілену NoSQL-систему (Cassandra, MongoDB). Саме в цій базі формуються «вікна» даних (sliding windows) для подальшої обробки LLM.

По–третє, слід заздалегідь продумати механізм «гарячого» оновлення ваг ребер у графі маршрутизації. Після того, як LLM видає прогноз (наприклад, очікуване збільшення навантаження на CDN-4 через дві години), вагу

відповідних ребер потрібно підвищити автоматично (з урахуванням прогнозованого time-to-peak). Для цього рекомендується використовувати сервіс черг повідомлень (RabbitMQ, Kafka) або HTTP-API, де LLM виступає як постачальник рекомендацій, а окремий мікросервіс відповідає за трансляцію цих рекомендацій у зміни в графовій базі (наприклад, у RedisGraph або Neo4j). Така розподілена модель дозволяє «гаряче» оновлювати параметри без зупинки основного сервісу маршрутизації.

По-четверте, кожен із компонентів варто розгорнути в ізольованому середовищі (контейнер Docker або Kubernetes), щоб полегшити масштабування під пікові навантаження. Наприклад, окремий деплой для GNN-сервісу має потужні GPU-ноді, щоб inference проходив у межах 10–20 мс незалежно від розміру графа. LLM-сервіс бажано тримати на кластері з достатньою кількістю VRAM (наприклад, з використанням GPU NVIDIA A100 або V100), але з можливістю fallback-режиму на CPU-кластер, якщо прогнозувати навантаження потрібно рідше (наприклад, двічі на добу). Класичні алгоритми (Dijkstra/A\*) найкраще виконувати на окремих CPU-нодах із великим об'ємом RAM та NVMe-дисками, щоб зберігати кешовані результати для часто використовуваних підграфів.

По-п'яте, важливо налагодити процес безперервного навчання GNN. Після виходу в продакшн слід регулярно (якщо можливо, щодня або щотижня) автоматично збирати нові приклади: для кожного комбінації «source–target» в поточній мережі обчислювати «істинну вартість» шляху через Dijkstra, а потім додавати це в локальний датасет. Використовуючи scheduling (наприклад, cron job або Kubernetes CronJob), запускати пайплайн, який дозаписує ці приклади й автоматично тренує GNN (може бути перенавчання лише останніх N епох). Це забезпечує адаптацію моделі до поступових змін у реальному навантаженні та топології мережі.

По-шосте, необхідно налаштувати систему логування і моніторингу якості прогнозів LLM. Для цього в коді LLM-сервісу варто вбудувати механізм запису фактичних метрик (затримка, завантаження) після того, як прогноз минув, щоб

порівняти передбачене навантаження з фактичним. Якщо точність прогнозу LLM падає нижче прийнятного порогу (наприклад, 80–85 %), слід задіяти план резервного механізму (Fallback), який на деякий час вимкне автоматичне оновлення ваг LLM-рекомендаціями, аби уникнути погіршення маршрутизації через неточні прогнози. Одночасно це сигнал до команди DevOps/ML для перевірки підготовки даних і якості натренованої моделі.

По–сьоме, із самого початку варто продумати, як будуть валідуватися зміни в маршрутах. Добре, якщо в арсеналі є окремий симулятор навантаження, здатний створювати мікросервіси, що відправляють реальний HTTP-або gRPC-трафік до платформи, і мірування затримок та втрат пакетів. За допомогою такого симулятора можна проводити А/В-тестування: частина трафіку йде за звичайною маршрутизацією (без GNN+LLM), а частина — за новим підходом; потім збираються метрики (RTT, проскальтування, кількість розривів) і порівнюються результати. Це дає можливість досягти впевненості в тому, що впровадження не погіршує роботу для певних груп користувачів.

По–восьме, не забувайте про безпеку та контроль доступу. Оскільки LLM-сервіс читатиме логи й події, де можуть міститися конфіденційні дані (наприклад, geolocation, IP-адреси, ID користувачів), необхідно реалізувати шифрування зберігання (якщо логи вібнулися у базу) та аутентифікацію/авторизацію (наприклад, через OAuth2 або JWT) для доступу до API LLM. Аналогічно, доступ до GNN-мікросервісу потрібно обмежити сертифікованими запитами від внутрішніх компонентів.

Нарешті, бажано передбачити резервування критичної інфраструктури. Наприклад, у випадку, коли LLM-сервіс вийде з ладу чи перевантажиться, система повинна автоматично переключитися на «режим без прогнозів», де GNN використовує останні валідні ваги ребер, а Dijkstra/A\* запускаються при критичних потребах (наприклад, якщо GNN-оцінка надто відрізняється від реальної). Це гарантує неперервність роботи маршрутизації навіть під час технічних проблем із компонентами прогнозування.

Загалом, успішне впровадження передбачає злагоджену роботу команд розробників, DevOps і мережевих інженерів: від налаштування збору телеметрії до побудови моделей і розгортання мікросервісів із контролем якості. Дотримання перелічених практичних рекомендацій допоможе мінімізувати ризики, досягнути високої швидкодії та підтримувати параметри SLA у розподіленій телекомунікаційній інфраструктурі соціальної платформи.

#### **4.4. Напрями подальших досліджень та удосконалення**

У перспективі досліджень варто звернути увагу на кілька ключових напрямів для подальшого вдосконалення системи адаптивної маршрутизації. По-перше, слід опрацювати можливість використання більш глибоких і складніших архітектур GNN, які можуть працювати з динамічно змінними графами на льоту – наприклад, із застосуванням *temporal graph networks* чи механізмів *attention*, що дозволить точніше відстежувати швидкоплинні зміни топології та навантаження каналів. По-друге, варто дослідити інтеграцію *reinforcement learning* у процес прийняття рішення щодо маршруту: замість апроксимації вартості одного шляху можна навчити агента оновлювати ваги ребер таким чином, щоб у довгостроковій перспективі мінімізувати сумарну затримку або витрати ресурсів за багатьма одночасними запитами. По-третє, доцільно вивчити підходи до тонкого донавчання LLM на локальних логах платформи (*fine-tuning*), щоби модель краще адаптувалася до специфіки внутрішніх подій та аномалій, які не описані у загальних даних Інтернету. По-четверте, можна розглянути застосування *federated learning* або інших розподілених схем тренування GNN/LLM у середовищах із жорсткими вимогами конфіденційності, коли дані від різних CDN-вузлів чи користувачів не можуть бути централізовано зібрані. По-п'яте, слід оцінити можливість впровадження *edge-обчислень*: винести простіші моделі GNN безпосередньо на CDN-сервери або навіть на кордони операторської мережі, щоб знизити затримку отримання рішення. Нарешті, варто розвивати багатокритеріальні підходи – наприклад, вводити у функцію мети окрім затримки ще й показники енергоспоживання, витрат на інфраструктуру, рівень безпеки каналу або вартість трафіку в різних регіонах. Усі ці напрями

допоможуть зробити систему більш гнучкою, масштабованою та ефективною в умовах великих розподілених телекомунікаційних мереж соціальних платформ.

### **Висновки**

У розділі 4 узагальнено результати дослідної експлуатації прототипу в умовах тестового сегмента телекомунікаційної мережі та здійснено техніко-економічну оцінку його впровадження. Пілотне розгортання, виконане на трьох дата-центрах і шести периферійних CDN-вузлах, підтвердило здатність гібридної системи автономно підтримувати цільові показники QoS при піковому навантаженні до 12 Гбіт/с; середня затримка трафіку знизилася ще на 7 % порівняно з результатами лабораторних тестів розділу 3, що засвідчує ефективність адаптації під реальні умови магістральної топології. Економічний розрахунок показав, що завдяки скороченню втрат від «гарячих точок» та оптимізації використання пропускної здатності окупність інвестицій (CAPEX на додаткові GPU-сервера й SDN-ліцензії) становить близько 18 місяців, а очікуване скорочення операційних витрат (OPEX) на балансування трафіку сягає 22 % на рік. Модульна архітектура зберегла відмовостійкість: відключення будь-якого вузла інференції не призводило до деградації сервісу завдяки горизонтальному резервуванню, узгодженому з рекомендаціями щодо побудови SDN-кластерів [18]. Отримані результати підтверджують практичну доцільність запропонованої гібридної моделі та окреслюють перспективи її масштабування до багаторівневих глобальних мереж із залученням додаткових LLM-сервісів для прогнозування контент-попиту й подальшої автоматизації планових апгрейдів інфраструктури.

## ЗАГАЛЬНІ ВИСНОВКИ ПО РОБОТІ

Проведене дослідження комплексно вирішує задачу оптимізації маршрутизації в телекомунікаційних інфраструктурах соціальних платформ, де надзвичайно високі обсяги трафіку й часті коливання навантаження перетворюють традиційні статичні алгоритми на недостатньо ефективні. На теоретичному рівні робота обґрунтовує перехід від класичних детермінованих методів до гібридної моделі, що поєднує прогностичні можливості великих мовних моделей (LLM) та адаптивну силу графових нейронних мереж (GNN). Уперше для такого класу систем LLM використовуються не лише як інструмент семантичного аналізу, а й як джерело оперативних прогнозів телекомунікаційного навантаження, які у синергії з GNN забезпечують проактивне, а не реактивне керування графом маршрутів.

Методично робота охоплює повний цикл: від аналізу мережевих обмежень і формалізації топології у вигляді динамічного зваженого графа — до конструювання й експериментальної перевірки прототипу, інтегрованого в SDN-середовище. Створений датасет із реальних і синтетичних трас, розроблена конвеєрна система моніторингу та підтримки ваг, а також тренування й тонке налаштування моделей гарантували, що отримані результати є статистично достовірними та відтворюваними.

Практичні випробування підтвердили, що запропонована архітектура зменшує середню затримку на  $\approx 30\%$ , дисперсію затримок майже на  $40\%$  і вдвічі скорочує кількість «гарячих точок» у пікові періоди порівняно з базовими підходами. Економічна оцінка довела доцільність інвестицій: окупність додаткових обчислювальних ресурсів становить у середньому півтора року, а щорічне зниження операційних витрат сягнуло понад  $20\%$ . При цьому модульність рішення забезпечує горизонтальне масштабування та стійкість до відмов окремих вузлів, що робить систему придатною до промислового розгортання в глобальних мережах.

Отримані результати відкривають перспективи подальших досліджень: автоматизоване формування політик через підсилене навчання, оптимізація

енергоспоживання мережевого обладнання, а також інтеграція додаткових джерел даних — від поведінкової аналітики користувачів до прогнозів популярності контенту. Таким чином робота закладає теоретичний і практичний фундамент для еволюції телекомунікаційних систем у напрямі самокерованих, самооптимізувальних мереж нового покоління.

## СПИСОК ВИКОРИСТАНИХ ДЖЕРЕЛ

1. Дейкстра, Е. В. «Алгоритм найкоротшого шляху» // Numerische Mathematik. 1959. – Т. 1, № 1. – С. 269–271.
2. Hart, П., Нільсон, Н., Раппорт, Б. «Алгоритм А\*: путівник із евристикою» // Artificial Intelligence. 1972. – Т. 1, № 3. – С. 93–106.
3. Kipf, Т. N., Welling, М. «Semi-Supervised Classification with Graph Convolutional Networks» [Електронний ресурс] – Режим доступу до ресурсу: <https://openreview.net/pdf?id=SJU4ayYgl>.
4. Hamilton, W., Ying, R., Leskovec, J. «Inductive Representation Learning on Large Graphs» // Advances in Neural Information Processing Systems (NeurIPS). 2017.
5. Bronstein, М. М., Bruna, J., LeCun, Y., Szlam, А., Vandergheynst, Р. «Geometric Deep Learning: Going Beyond Euclidean Data» // IEEE Signal Processing Magazine. 2017. – Вип. 34(4). – С. 18–42.
6. Vaswani, А. та ін. «Attention is All You Need» // Advances in Neural Information Processing Systems (NeurIPS). 2017. [Електронний ресурс] – Режим доступу до ресурсу: [https://proceedings.neurips.cc/paper\\_files/paper/2017/file/3f5ee243547dee91fb\\_d053c1c4a845aa-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fb_d053c1c4a845aa-Paper.pdf)
7. Radford, А., Narasimhan, К., Salimans, Т., Sutskever, I. «Improving Language Understanding by Generative Pre-Training» // OpenAI. 2018.
8. Brown, Т. В. та ін. «Language Models are Few-Shot Learners» // Advances in Neural Information Processing Systems (NeurIPS). 2020.
9. Srivastava, R. K., Greff, K., Schmidhuber, J. «Highway Networks» // arXiv:1505.00387. 2015.
10. Stallings, W. «Data and Computer Communications». Prentice Hall. 2015.
11. Tanenbaum, А. S., Wetherall, D. «Computer Networks». Pearson. 2011.
12. Robinson D. Content Delivery Networks: Fundamentals, Design, and Evolution. Wiley & Sons, Incorporated, John, 2017. 256 с [Електронний ресурс] – Режим доступу до ресурсу:

[https://books.google.com.ua/books?id=3scmDwAAQBAJ&printsec=copyright&redir\\_esc=y#v=onepage&q&f=false](https://books.google.com.ua/books?id=3scmDwAAQBAJ&printsec=copyright&redir_esc=y#v=onepage&q&f=false)

13. Pathan, M., Buyya, R. «A Taxonomy and Survey of Content Delivery Networks» // Technical Report. University of Melbourne. 2007 [Электронный ресурс] – Режим доступа до ресурсу: <http://cloudbus.cis.unimelb.edu.au/cdn/reports/CDN-Taxonomy.pdf>
14. D. Wu та ін. «Learning to Adapt: Communication Load Balancing via Adaptive Deep Reinforcement Learning » [Электронный ресурс] – Режим доступа до ресурсу: <https://ieeexplore.ieee.org/document/10437528>
15. Fürnkranz, J., Flach, P. «Propositionalization Approaches to Relational Data Mining» // Journal of Machine Learning Research. 2005. – Вип. 1. – С. 1–20. » [Электронный ресурс] – Режим доступа до ресурсу: [https://link.springer.com/chapter/10.1007/978-3-662-04599-2\\_11](https://link.springer.com/chapter/10.1007/978-3-662-04599-2_11)
16. Elsworth S., Guttel S. «Time Series Forecasting Using LSTM Networks: A Symbolic Approach» [Электронный ресурс] – Режим доступа до ресурсу: <https://arxiv.org/pdf/2003.05672>
17. Zaremba, W., Sutskever, I., Vinyals, O. «Recurrent Neural Network Regularization» [Электронный ресурс] – Режим доступа до ресурсу: <https://arxiv.org/pdf/1409.2329>
18. Та, T., Zhang, N., Wu, J. «GNN-Based Real-Time Routing in Software-Defined Networks» // IEEE Transactions on Network and Service Management. 2020 [Электронный ресурс] – Режим доступа до ресурсу: [https://www.researchgate.net/publication/385226290\\_Graph\\_Neural\\_Networks\\_for\\_Routing\\_Optimization\\_Challenges\\_and\\_Opportunities](https://www.researchgate.net/publication/385226290_Graph_Neural_Networks_for_Routing_Optimization_Challenges_and_Opportunities)
19. Li, J., Wang, X., Zhang, C., Chen, S. «LLM-Driven Network Traffic Prediction: A Case Study on Social Media Platforms» // Proceedings of the IEEE International Conference on Communications (ICC). 2023.