

МАГІСТЕРСЬКА РОБОТА

МР. ШМ - 51.00.00.000 ПЗ

Група ШМ-23-2

Михайлов Степан

2024

Івано-Франківський національний технічний університет нафти і газу

Інститут інформаційних технологій

Кафедра інженерії програмного забезпечення

Михайлов Степан Сергійович

(прізвище, ім'я, по батькові)

УДК 004.942
(індекс)

МАГІСТЕРСЬКА РОБОТА

Моделі та методи динамічного навчання на графових

представленнях стрімінгу

(назва роботи)

Інженерія програмного забезпечення

(назва освітньої програми)

121 - Інженерія програмного забезпечення

(шифр і назва спеціальності)

Михайлов С.С.

(підпис, ініціали та прізвище здобувача освітнього ступеня)

Науковий керівник Храбатин Роман Ігорович, к.ф.-м.н., доцент

(прізвище, ім'я, по батькові, науковий ступінь, вчене звання)

Допущено до захисту

Завідувач кафедри

доц. Бандура В.В.

(посада) (підпис) (дата) (ініціали та прізвище)

Нормоконтроль

доц. Вовк Р.Б.

(посада) (підпис) (дата) (ініціали та прізвище)

Робота містить результати власних досліджень. Використання ідей, результатів і текстів інших авторів мають посилання на відповідне джерело

Івано-Франківськ – 2024

Івано-Франківський національний технічний університет нафти і газу

Інститут інформаційних технологій

Кафедра інженерії програмного забезпечення

Освітній рівень магістр

Спеціальність 121 – Інженерія програмного забезпечення

ЗАТВЕРДЖУЮ:

Зав. кафедрою

ІІЗ

доц.

В.В. Бандура

“ 04 ” вересня 2024 р.

ЗАВДАННЯ

НА МАГІСТЕРСЬКУ РОБОТУ СТУДЕНТУ

Михайлову Степану Сергійовичу

(прізвище, ім'я, по-батькові)

1. Тема магістерської роботи “ Моделі та методи динамічного навчання на графових представленнях стрімінгу”

керівник проекту (роботи) Храбатин Роман Ігорович, к.ф.-м.н., доцент

затверджені наказом закладу вищої освіти від “ 22 ” листопада 2024 р. № 781/7

2. Строк подання студентом проекту (роботи) 15 грудня 2024 р.

3. Вихідні дані до проекту (роботи) Теоретичні концепції та формальні моделі побудови та функціонування інформаційних технологій динамічного навчання

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

1. Аналіз предметної області застосування відео стрімінгу та потокового відео

2. Моделі та алгоритми динамічного та статичного навчання на графових представленнях

3. Методологія, моделі та методи динамічного навчання на графових представленнях стрімінгу

4. Прогнозування пропускнуої здатності за допомогою навчання динамічних графів

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

1. Приклад трансляції відео в реальному часі для компанії з двома офісами (рис. 1.1)

2. Архітектура потокового сервісу Nive (рис. 1.2)

3. Архітектура PLS агента (рис. 1.3)

4. Приклади оверлейних мереж (рис. 1.4)

5. Структура RTP пакету показана в 32-бітних (4-байтових) блоках (рис. 2.1)

6. Консультанти розділів проекту (роботи)

Розділ	Консультант	Підпис, дата
Перевірка на плагіат	доц., к.т.н. Вовк Р.Б.	

7. Дата видачі завдання 04 вересня 2024 р.

Керівник _____
(підпис)

Завдання прийняв до виконання _____
(підпис)

КАЛЕНДАРНИЙ ПЛАН

№ п/п	Назви етапів магістерської роботи	Строк виконання етапів роботи	Примітка
1	Підбір і вивчення літератури по темі магістерської роботи	15.09.2024	виконано
2	Аналіз концепцій та алгоритмів предметної області	29.09.2024	виконано
3	Дослідження предметної області застосування відео стрімінгу та потокового відео	15.10.2024	виконано
4	Моделі та алгоритми динамічного та статичного навчання на графових представленнях	08.11.2024	виконано
5	Методологія, моделі та методи динамічного навчання на графових представленнях стрімінгу	20.11.2024	виконано
6	Прогнозування пропускнуої здатності за допомогою навчання динамічних графів	01.12.2024	виконано
7	Затвердження пояснювальної записки роботи завідувачем кафедри	15.12.2024	виконано

Студент – магістр _____
(підпис)

Керівник роботи _____
(підпис)

АНОТАЦІЯ

Магістерська робота: 78 с., 24 рис., 2 табл., 51 джерел.

Тема: Моделі та методи динамічного навчання на графових представленнях стрімінгу

Об'єкт дослідження: процеси передачі даних у системах потокового відео та відеострімінгу.

Мета роботи: розробка методології, моделей і методів динамічного навчання на графових представленнях для прогнозування пропускної здатності мереж потокового відео.

Предмет дослідження: моделі та методи динамічного навчання на графових представленнях, що використовуються для оптимізації передачі даних у потоковому відео.

Результати дослідження

В роботі запропоновано комплексну методологію для прогнозування пропускної здатності мереж потокового відео на основі динамічного навчання графів та розроблено адаптивну модель (ABD) з використанням механізму уваги графа, що підвищує точність аналізу динамічних мережевих процесів.

Висновок

Визначено, що інтеграція динамічного навчання на графах із механізмами уваги є ефективним підходом до прогнозування пропускної здатності стрімінгових систем.

**ДИНАМІЧНЕ НАВЧАННЯ, ПРЕДСТАВЛЕННЯ ГРАФІВ,
ПОТОКОВЕ ВІДЕО, ВІДЕОСТРІМІНГ, ПРОГНОЗУВАННЯ
ПРОПУСКНОЇ ЗДАТНОСТІ, МЕХАНІЗМ УВАГИ, НЕЙРОННІ
МЕРЕЖІ, ОПТИМІЗАЦІЯ ПЕРЕДАЧІ ДАНИХ**

ABSTRACT

Master Thesis: 78 pp., 24 fig., 2 tab., 51 sources.

Thesis Subject: Models and methods of dynamic learning on graph representations of streaming

Object of research: data transfer processes in video streaming and video streaming systems.

Purpose of work: development of methodology, models and methods of dynamic learning on graph representations for predicting the bandwidth of streaming video networks.

Subject of research: models and methods of dynamic learning on graph representations used to optimize data transmission in streaming video.

Research results

The paper proposes a comprehensive methodology for predicting the bandwidth of streaming video networks based on dynamic graph learning and develops an adaptive model (ABD) using the graph attention mechanism, which increases the accuracy of the analysis of dynamic network processes.

Conclusion

It is determined that the integration of dynamic learning on graphs with attention mechanisms is an effective approach to predicting the bandwidth of streaming systems.

**DYNAMIC LEARNING, GRAPH REPRESENTATION,
STREAMING VIDEO, VIDEO STREAMING, BANDWIDTH
PREDICTION, ATTENTION MECHANISM, NEURAL NETWORKS,
DATA TRANSMISSION OPTIMIZATION**

ЗМІСТ

ПЕРЕЛІК УМОВНИХ СКОРОЧЕНЬ	9
ВСТУП.....	10
РОЗДІЛ 1. АНАЛІЗ ПРЕДМЕТНОЇ ОБЛАСТІ ЗАСТОСУВАННЯ ВІДЕО СТРІМІНГУ ТА ПОТОКОВОГО ВІДЕО	13
1.1. Особливості трансляції відео в реальному часі як засобу для спілкування	13
1.1.1. Використання відео стрімінгу в компаніях	14
1.2. Корпоративні послуги розповсюдження потокового відео	15
1.3. Проблема накладених (оверлейних) мереж.....	19
1.4. Опис Gossip протоколів	23
1.5. Оцінка пропускної здатності за допомогою Dynamic Graph Representation Learning	25
1.6. Цілі та методологія дослідження	26
Висновки до розділу	27
РОЗДІЛ 2. МОДЕЛІ ТА АЛГОРИТМИ ДИНАМІЧНОГО ТА СТАТИЧНОГО НАВЧАННЯ НА ГРАФОВИХ ПРЕДСТАВЛЕННЯХ	29
2.1. Алгоритми передачі даних, які використовуються в трансляції відео в реальному часі.....	29
2.2. Концепція графового представлення навчання.....	32
2.3. Навчання статичного представлення графів	35
2.3.1. Методи факторизації матриць.....	35
2.3.2. Методи випадкового блукання	36
2.3.3. Підходи нейронних мереж	38
2.4 Концепція навчання динамічного представлення графів	41
2.4.1. Методи факторизації матриць.....	41
2.4.2. Методи випадкового блукання	41

2.4.3. Підходи на основі нейронних мереж	41
Висновки до розділу	44
РОЗДІЛ 3. МЕТОДОЛОГІЯ, МОДЕЛІ ТА МЕТОДИ ДИНАМІЧНОГО НАВЧАННЯ НА ГРАФОВИХ ПРЕДСТАВЛЕННЯХ СТРІМІНГУ	46
3.1. Представлення методології	46
3.1.1. Збір та опис даних.....	47
3.1.2. Аналіз даних.....	48
3.2. Прогнозування пропускної здатності за допомогою навчання представлення динамічних графів	49
3.2.1. Визначення проблеми.....	50
3.2.2. Запропонований підхід (ABD)	51
3.2.3. Механізм уваги графа	52
3.2.3. Навчання моделі	55
3.3. Оцінка продуктивності запропонованої моделі	57
3.3.1. Набори даних	57
3.3.2. Проведення імітаційних експериментів	61
3.3.3. Аналіз параметрів	64
Висновки до розділу	69
ВИСНОВКИ	71
ПЕРЕЛІК ВИКОРИСТАНОЇ ЛІТЕРАТУРИ.....	73

ПЕРЕЛІК УМОВНИХ СКОРОЧЕНЬ

CCPA - California consumer privacy act.

CDN - Content delivery network.

DNGR - Deep neural networks for learning graph representations.

DVGAE - Dynamic variational graph autoencoder.

EGCN - Evolve graph convolutional network.

GAN - Generative adversarial network.

GAT - Graph attention network.

Gbps - Gigabits per second.

GCN - Graph convolutional network.

GDPR - General data protection regulation.

LINE - Large-scale information network embedding.

LSTM - Long short-term memory.

MAE - Mean absolute error.

Mbps - Megabits per second.

MSE - Mean square error.

P2P - Peer to peer.

PCA - Principal component analysis.

PPMI - Positive pointwise mutual information.

RMSE - Root mean square error.

SDNE - Structural deep network embedding.

VGAE - Variational graph autoencoder.

WebRTC - Web real-time communication.

ВСТУП

Актуальність теми.

Сучасний розвиток технологій передачі даних, включаючи потокове відео та відеострімінг, значною мірою визначає ефективність цифрової взаємодії в багатьох галузях, таких як корпоративна комунікація, онлайн-освіта, розваги та телемедицина. В умовах стрімкого зростання обсягів відеоконтенту та кількості користувачів, постачальники послуг зіштовхуються з низкою технічних викликів, серед яких найбільш критичними є:

- Проблема масштабованості. Із збільшенням кількості користувачів та пристроїв необхідно забезпечувати стабільну роботу мереж навіть за умов пікових навантажень.

- Динамічність умов мережі. Нестабільність пропускної здатності, затримки та втрата пакетів значно впливають на якість потокового відео, особливо в реальному часі.

- Ефективність управління трафіком. Оптимізація використання ресурсів мережі є важливим завданням для мінімізації витрат і забезпечення якості обслуговування (QoS).

Традиційні методи моделювання та управління мережевими процесами часто не враховують динамічну природу змін у топології та поведінці мережі. У цьому контексті, застосування підходів на основі графових представлень, особливо з використанням динамічного навчання, є перспективним напрямом.

Динамічне навчання на графах дозволяє:

- ефективно моделювати складні взаємозв'язки між елементами мережі;
- адаптуватися до змін мережеских умов у реальному часі;
- прогнозувати поведінку мережі для підвищення її продуктивності.

Особливу актуальність тема має для розробки систем відеострімінгу, які повинні забезпечувати не лише високу якість передачі даних, а й

адаптивність до умов нестабільної мережі. Зокрема, запропоновані в рамках дослідження моделі та методи можуть бути застосовані в корпоративних і комерційних середовищах для підвищення ефективності систем потокового відео.

Таким чином, дослідження спрямоване на вирішення однієї з ключових проблем сучасних мережевих технологій — забезпечення надійності та адаптивності систем потокового відео в умовах динамічних змін, що визначає його високу актуальність.

Мета дослідження – розробка методології, моделей і методів динамічного навчання на графових представленнях для прогнозування пропускної здатності мереж потокового відео.

Об'єкт дослідження – процеси передачі даних у системах потокового відео та відеострімінгу.

Предмет дослідження – моделі та методи динамічного навчання на графових представленнях, що використовуються для оптимізації передачі даних у потоковому відео.

Завдання дослідження:

1. Проаналізувати предметну область застосування потокового відео та відеострімінгу, зокрема, технології корпоративного використання та основні проблеми, пов'язані з мережевою інфраструктурою.

2. Дослідити існуючі моделі та алгоритми навчання на статичних і динамічних графових представленнях.

3. Розробити методологію та модель для прогнозування пропускної здатності мереж з використанням динамічного навчання на графах.

4. Провести оцінку продуктивності запропонованої моделі за допомогою імітаційних експериментів.

Методи дослідження:

- Теоретичний аналіз літератури та існуючих моделей графового представлення.

- Методи навчання статичних і динамічних графів, включаючи факторизацію матриць, випадкові блукання та нейронні мережі.
- Механізми уваги для підвищення точності аналізу динамічних графів.
- Імітаційні експерименти для оцінки продуктивності моделі.

Наукова новизна отриманих результатів

Запропоновано комплексну методологію для прогнозування пропускну здатності мереж потокового відео на основі динамічного навчання графів та розроблено адаптивну модель (ABD) з використанням механізму уваги графа, що підвищує точність аналізу динамічних мережевих процесів.

Практичне значення результатів

Запропоновані моделі та методи можуть бути застосовані в системах потокового відео для: зниження затримок і підвищення якості передачі даних; адаптації до динамічних умов мережі; створення інтелектуальних систем управління трафіком у корпоративних та комерційних середовищах.

Структура магістерської роботи. Робота складається зі вступу, трьох розділів та висновків. Загальний обсяг роботи становить 78 сторінок, і містить 24 рисунки, 2 таблиці, список використаних джерел із 51 найменування.

РОЗДІЛ 1. АНАЛІЗ ПРЕДМЕТНОЇ ОБЛАСТІ ЗАСТОСУВАННЯ ВІДЕО СТІМІНГУ ТА ПОТОКОВОГО ВІДЕО

1.1. Особливості трансляції відео в реальному часі як засобу для спілкування

Трансляція відео в реальному часі (лайв стрімінг) стрімко набирає популярність як засіб комунікації, пропонуючи унікальні можливості, яких немає в інших форматах. Ось деякі ключові особливості:

- Миттєвий зворотній зв'язок: Глядачі можуть взаємодіяти зі стримером у реальному часі через чат, коментарі, реакції. Це створює відчуття спільної присутності та активної участі.

- Непідробність та автентичність: На відміну від відредагованих відео, прямі трансляції демонструють події такими, якими вони є, що сприяє більшій довірі та емоційному зв'язку з аудиторією.

- Широке охоплення: Лайв стріми доступні будь-де та будь-коли, забезпечуючи глобальний охоплення аудиторії.

- Різноманітність контенту: Від ігрових стрімів та освітніх лекцій до віртуальних подорожей та кулінарних шоу – лайв стрімінг пропонує безмежні можливості для творчості та самореалізації.

- Доступність: Для початку трансляції не потрібно дорогого обладнання. Смартфон та стабільний інтернет – це все, що потрібно для багатьох видів стрімів.

Важливе людське спілкування. У повсякденному житті ми сприймаємо як належне регулярний обмін інформацією різними способами. Одним із найновіших способів спілкування є живе відео. Безсумнівно, це спілкування було доведено ефективним, оскільки воно імітує оригінальне спілкування віч-на-віч. Живі відеорозмови всюди, вони використовуються в щоденних дискусіях між сім'ями та друзями, в навчальних закладах для онлайн-навчання та багато іншого. Підприємства також визнали важливість живого

відео та переваги, які воно приносить у робоче середовище. Це практичний інструмент спілкування, корисний для онлайн-зустрічей між членами компанії, а також зустрічей у різних компаніях.

1.1.1. Використання відео стрімінгу в компаніях

Оскільки компанії складаються з сотень або навіть тисяч співробітників, технологія, яка забезпечує живий відеозв'язок, повинна мати можливість масштабування. На жаль, передача високоякісного відеовмісту в корпоративній мережі за своєю суттю є складною проблемою.

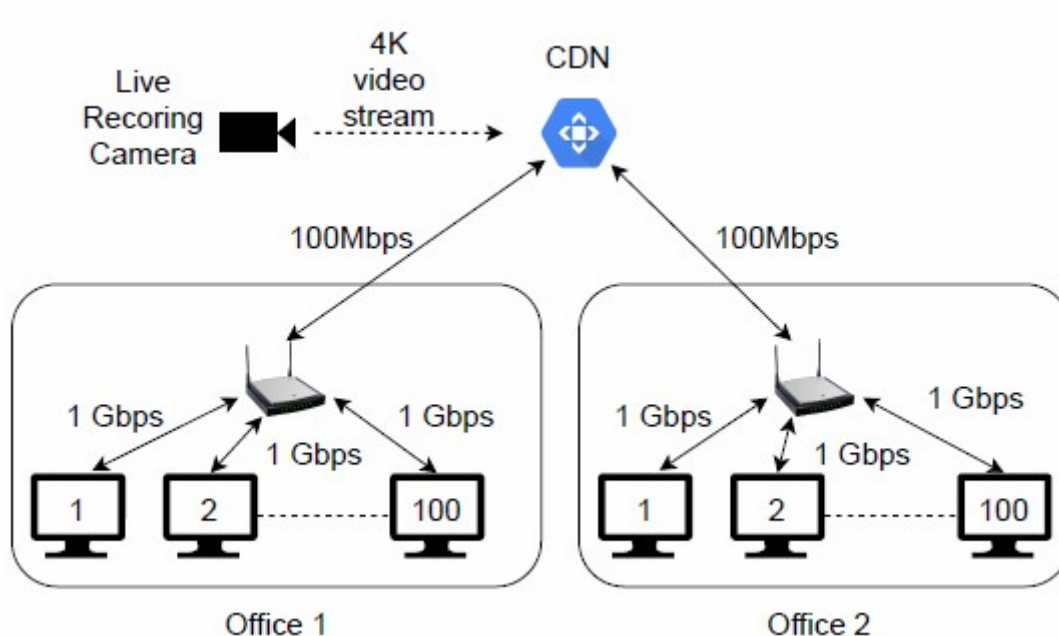


Рис. 1.1. Приклад трансляції відео в реальному часі для компанії з двома офісами. Кожен офіс має високошвидкісну локальну мережу, тоді як пропускна здатність за межами офісу зазвичай менша

Як показано на рисунку 1.1, процес розповсюдження відео можна розділити на дві частини:

- У першій частині відео доставляється з джерела (тобто камери для запису в реальному часі) до мережі доставки контенту (CDN). CDN — це

розподілена система, яка виконує всю необхідну обробку відео, доки вміст не буде готовий до доставки клієнтам.

- У другій частині оброблений контент доставляється з CDN на пристрої користувача. Наша основна увага — друга частина, яка передбачає розповсюдження контенту на велику кількість пристроїв.

Друга частина процесу вимагає підключення кожного окремого пристрою безпосередньо до CDN. Потім CDN передає той самий вміст на кожен пристрій окремо. У цьому сценарії дуже ймовірно, що корпоративна мережа буде перевантажена, оскільки її пропускна здатність нижча, ніж вимоги до мережі. Пояснимо це на прикладі. Припустимо, що ми хочемо передати потокове відео 4K в обидва офіси на рисунку 1.1. Якщо припустити, що відеопотік 4K потребує 20 Мбіт/с, то кожен офіс зі 100 пристроями загалом потребує 2 Гбіт/с. Однак за межами кожного офісу пропускна здатність обмежена 100 Мбіт/с. У цьому сценарії очевидно, що вимоги до мережі перевищують пропускну здатність мережі. Це призводить до перевантаження мережі, що вплине не лише на якість відеопотоку, але й на інші послуги (наприклад, VoIP, хмарне сховище).

1.2. Корпоративні послуги розповсюдження потокового відео

Попередню проблему не можна легко вирішити. У той же час трансляція відео в прямому ефірі в корпоративних мережах займає велику частку ринку, яка з часом зростає, оскільки все більше компаній використовують технології живого відео для повсякденного спілкування. Щодо цього, існує кілька компаній, які націлені на ринок корпоративного потокового відео в прямому ефірі. Наприклад, Nive Streaming є такою компанією, яка показала великий успіх за останні роки.

Nive Streaming – це компанія, яка надає масштабовані рішення для корпоративного розповсюдження відео в реальному часі. Їхніми основними цільовими клієнтами є 500 багатих компаній, у яких працюють тисячі чи

навіть сотні тисяч співробітників. Основним продуктом компанії є сервіс потокової передачі контенту, який може задовольнити потреби клієнтів, не вимагаючи додаткового обладнання. Складність проблеми пов'язана з тим фактом, що пропускна здатність зовнішнього з'єднання значно нижча, ніж фактична пропускна здатність, необхідна для розповсюдження відеовмісту на всі кінцеві точки всередині офісу (рис. 1.1). У зв'язку з цим розповсюдження відео має бути здійснено більш винахідливо.

Щоб уникнути вищезгаданої ситуації, Hive Streaming надає високоефективне масштабоване рішення. Їхній продукт дозволяє інтелектуально передавати вміст на численні пристрої, щоб уникнути перевантаження мережі. Загалом, ключова ідея Hive Streaming полягає в тому, щоб мінімізувати кількість разів, коли вміст доставляється з CDN на пристрої, що є причиною заторів у мережі. Щоб досягти цього, ідентичний вміст розподіляється між пристроями в одній мережі. Оскільки пропускна здатність усередині локальної мережі зазвичай вища, вміст можна поширювати без проблем. Все це відбувається з обмеженням, що якість відео залишається незмінною.

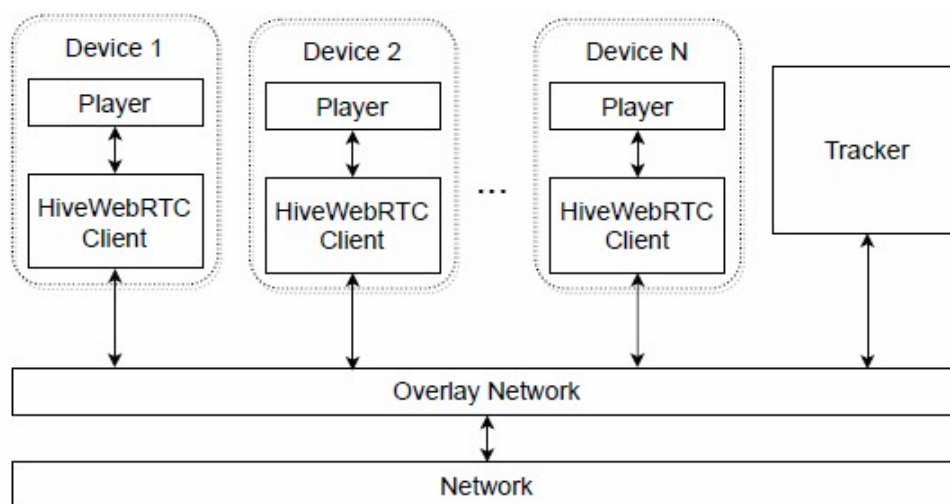


Рис. 1.2. Архітектура потокового сервісу Hive

Архітектуру Hive Streaming Service можна побачити на рисунку 1.2. Контент доставляється в одноранговому режимі (P2P), де пристрої

утворюють накладену мережу P2P. Все починається з пристроїв користувача, які містять відеоплеєр і клієнт Hive WebRTC. Відеопрогравач використовує вміст, наданий клієнтом Hive WebRTC. Клієнт Hive WebRTC є основним компонентом сервісу. Саме модуль формує оверлейну мережу та розподіляє відео між пристроями. Розповсюдження передбачає не лише отримання вмісту кожному одноранговому вузлу, але й обмін вмістом у накладній мережі. Окрім накладення P2P, архітектура містить Tracker, який діє як служба оркестровки. Його основні обов'язки включають завантаження та надання метаданих клієнтам Hive WebRTC. Деталі сервісу Tracker опущено, оскільки вони не стосуються цього проекту.

Як згадувалося, клієнти Hive WebRTC доставляють вміст, утворюючи накладену мережу. З точки зору програмного забезпечення, Hive WebRTC є розширенням технології WebRTC, адаптованим до потреб компанії. Великою новинкою Hive WebRTC є те, що він доставляє вміст за допомогою SmoothCache 2.0 [1].

SmoothCache 2.0 — це платформа розподіленого кешу для адаптивного HTTP трансляція контенту в прямому ефірі на основі однорангових мереж. Основна мета SmoothCache 2.0:

- 1) сформувати накладену мережу,
- 2) динамічно оновлювати з'єднання всередині накладення в міру зміни мережі,
- 3) забезпечити своєчасну доставку вмісту на кожен пристрій.

На відміну від попередніх систем, це дозволяє динамічно адаптувати бітрейт вмісту для однорангового вузла відповідно до його пропускної здатності та апаратного забезпечення. Ця функція реалізована шляхом покладання відповідальності за вибірку вмісту на сам одноранговий пристрій, а не на віддалений сервер.

SmoothCache - це P2P-агент для потокової трансляції в реальному часі (PLS), який виглядає як HTTP-проксі для відео, встановлений на кожному комп'ютері глядача. За наявності SmoothCache всі запити програвача

перенаправляються до локального агента PLS, замість того, щоб надходити безпосередньо до джерела потоку. Це єдина зміна, яка потрібна порівняно зі звичайною роботою без нашого агента.

SmoothCache використовує ряд допоміжних служб, які представляють нашу централізовану інфраструктуру. Важливо зазначити, що наша система розроблена таким чином, щоб бути прозорою для вже існуючої інфраструктури потокової передачі, щоб забезпечити як нейтралітет постачальника, так і спрощене розгортання.

Таким чином, ми не керуємо власним джерелом потоку, а натомість припускаємо, що існує вже опублікований потік, наприклад, з мережі доставки контенту (CDN), який можна відтворювати безпосередньо з його джерела, і наші клієнти прозоро використовуватимуть нашу систему для мінімізації навантаження на це джерело. Тобто, ми вирішуємо проблему зменшення навантаження на джерело потоку так само, як це робить CDN: шляхом побудови схеми кешування.

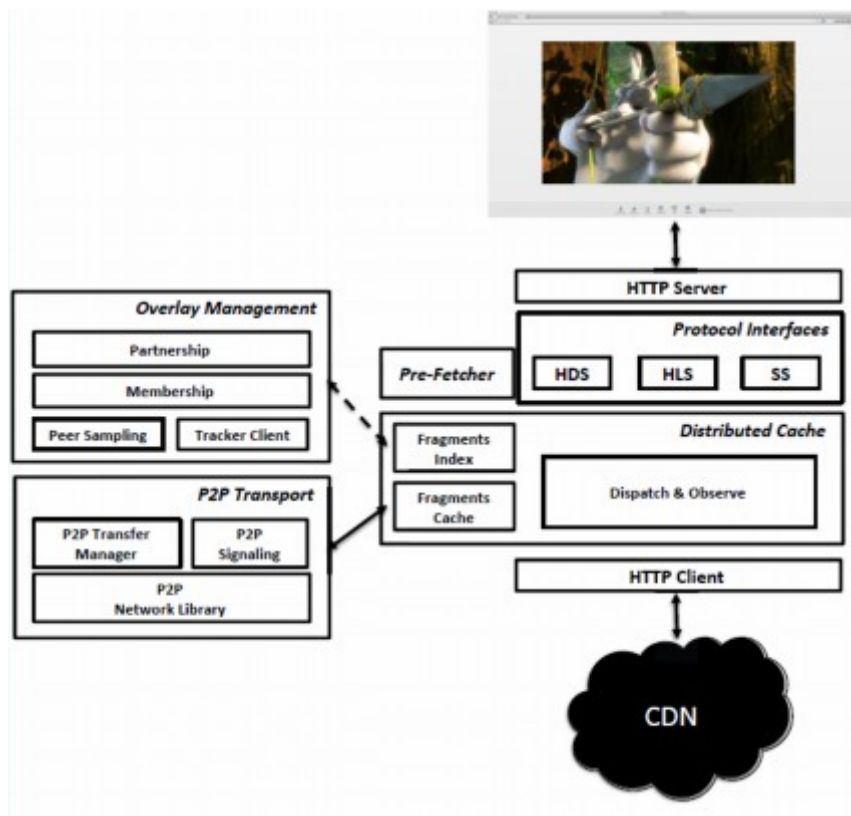


Рис. 1.3. Архітектура PLS агента

Внутрішньо, як показано на рисунку 1.3, агент PLS має HTTP-сервер, який отримує послідовні HTTP-запити від програвача. Він визначає, для якого протоколу зроблені запити, будь то Apple HTTP Live Streaming (HLS) [6], Adobe HDS [2] або Microsoft Smooth Streaming (SS) [13], і пересилає їх до компонента, специфічного для протоколу, який витягує відповідну інформацію та генерує нейтральні до протоколу запити до компонента розподіленого кешу.

Останній вирішує, чи можна обслуговувати запит із джерела або з P2P-мережі, використовуючи індекс фрагментів - компонент, який відстежує фрагменти, доступні на інших вузлах. Файли маніфесту завжди отримуються з джерела і ніколи не кешуються, оскільки вони часто оновлюються і містять актуальну інформацію про нові фрагменти. Описаний вище процес має реактивний характер, тобто всі запити до розподіленого кешу надходять від програвача після перетворення, специфічного для протоколу. На додаток до цього, SmoothCache реалізує процес, який проактивно завантажує фрагменти заздалегідь, перш ніж вони будуть запитані програвачем. Компонент, відповідальний за цю операцію, називається Pre-Fetcher.

Для реалізації абстракції розподіленого кешу, описаної вище, набір компонентів для управління оверлеєм працює разом, щоб підтримувати сусідство агента та оновлювати індекс фрагментів. Для ефективного отримання даних від інших учасників HTTP-клієнт підтримує постійне з'єднання з джерелом потоку, тоді як транспортний рівень P2P використовує найсучасніший NAT traversal та протокол контролю перевантаження прикладного рівня зі змінним пріоритетом [15] для передачі даних до та від інших вузлів.

1.3. Проблема накладених (оверлейних) мереж

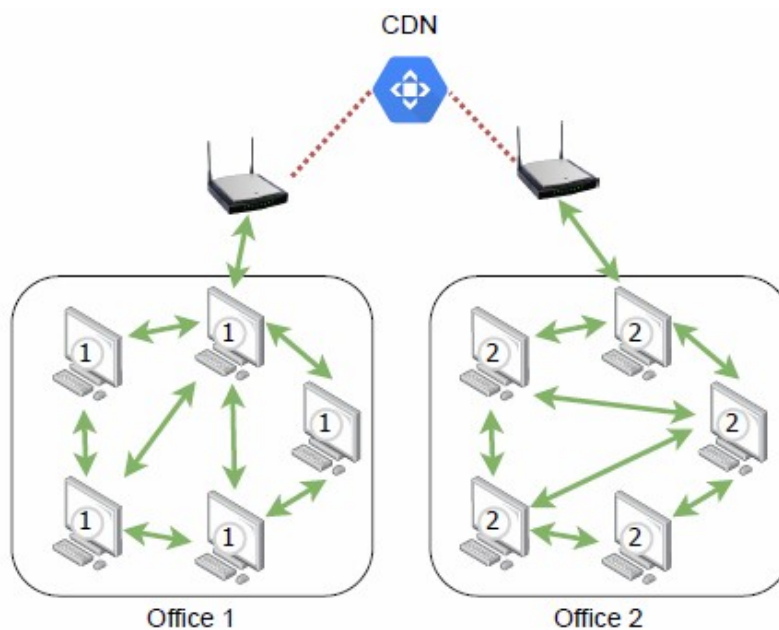
У цьому проекті ми зосередимо увагу на способі формування накладення. Кожен одноранговий вузол підтримує набір з'єднань з іншими

одноранговими вузлами, з якими можна обмінюватися вмістом. Грубо кажучи, ці з'єднання слід ретельно вибирати з огляду на пропускну здатність між одноранговими вузлами. Крім того, зв'язки між одноранговими вузлами мають бути надлишковими. Це необхідно для забезпечення того, що вміст завжди буде доставлено через всю мережу, навіть у разі часткових збоїв. Зв'язки кожного однорангового вузла можна класифікувати як у партнерів, які поведуться як постачальники контенту (розсівачі), і зовнішніх партнерів, які споживають вміст (п'явки). Наразі SmoothCache 2.0 використовує набір евристик для вибору партнерів кожного однорангового вузла. Ці евристики організують мережу, використовуючи інформацію про пропускну здатність між одноранговими вузлами.

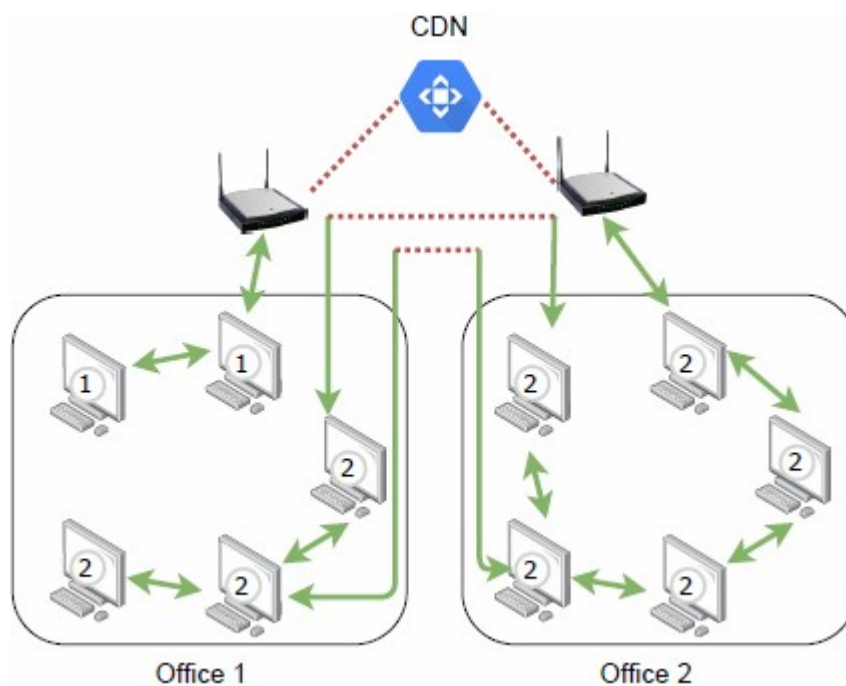
Загалом кажучи, обов'язки кожного однорангового вузла щодо доставки вмісту пропорційні смузі пропускання, яку він утримує. Наприклад, однорангові вузли з високою пропускну здатністю відповідають за доставку вмісту в різні області графіка, тоді як однорангові вузли з низькою пропускну здатністю гарантують, що вони можуть отримувати вміст без призупинення відтворення відео. Оскільки умови роботи мережі не гарантовано стабільні, накладена мережа не є статичною. Важливо, щоб оверлей розвивався з часом, щоб адаптуватися до змін підключення між пристроями.

Загалом оверлейну мережу можна розглядати як платформу розподіленого кешування, яка мінімізує витрати на розповсюдження вмісту без шкоди для якості досвіду. Безсумнівно, продуктивність такої системи значною мірою залежить від способу створення накладної мережі. Якщо партнери вибирають партнерів випадковим чином, загальна продуктивність постраждає через перевантаження мережі. Це тому, що мережа з низькою пропускну здатністю, яка існує за межами офісу, може бути перевантажена. З іншого боку, якщо партнери вибираються ретельно, вміст бездоганно транслюватиметься всередині офісів, а площа мережі залишатиметься мінімальною.

Приклад ефективного та неефективного накладання можна побачити на
рисунку 1.4.



а) Приклад гарної оверлейної мережі. З'єднання між офісом 1 та 2
відсутні



б) Приклад поганої оверлейної мережі. З'єднання між офісом 1 та 2
спричиняють перевантаження

Рис. 1.4. Приклади оверлейних мереж

На рисунку 1.4 зелені стрілки вказують на високошвидкісні з'єднання, тоді як червоні пунктирні лінії вказують на з'єднання, які можуть спричинити перевантаження мережі.

На цьому рисунку контент доставляється з CDN до двох офісів. У першому прикладі пристрої кожного офісу підключені один до одного. Зверніть увагу, що жодні пристрої не підключені між офісами. Оскільки пропускна здатність у кожному офісі висока, розподіл контенту між пристроями здійснюється ефективно. У цьому прикладі лише один пристрій на офіс вимагає потокової передачі даних безпосередньо з CDN, оскільки решта з них отримуватимуть той самий вміст через накладену мережу. У другому прикладі пристрої не підключено оптимальним чином. Ми бачимо, що пристрої з офісу 1 підключені до пристроїв з офісу 2. У результаті накладена мережа передає вміст між двома офісами. Така ситуація може спричинити перевантаження мережі, оскільки пропускна здатність між офісами 1 і 2 зазвичай обмежена.

До недавнього часу HIVE Streaming покладався на приватну IP-інформацію для створення накладеної мережі. У принципі, якщо два однорангових вузла знаходяться в одному діапазоні приватних IP-адрес, вони належать до одного офісу. Таким чином, вузли можуть бути згруповані в офіси за допомогою приватних діапазонів IP-адрес. Оскільки однорангові пристрої знають, які інші однорангові мережі знаходяться в одній мережі, вони можуть вибірково з'єднуватися з ними та уникати з'єднань між офісами. Однак ця ідея вже не може бути застосована.

З огляду на нещодавні закони про захист даних та конфіденційності, піри не можуть бути пов'язані з персональною інформацією. На відміну від попереднього прикладу, приватна IP-адреса була позначена як цифрова персональна інформація. Таким чином, вона юридично недоступна, і сучасні браузері блокують доступ до неї. Натомість, піри можуть бути пов'язані лише з їхньою публічною IP-адресою. На жаль, публічну адресу не можна використовувати так само, як приватну.

1.4. Опис Gossip протоколів

Оскільки піри не можуть бути ідентифіковані за їхньою приватною IP-адресою, Hive Streaming впроваджує протокол пліток, щоб допомогти у створенні оверлейної мережі. По суті, цей протокол виконує тести пропускної здатності між пірами. Мета цих тестів - визначити, чи знаходяться два піри в одному офісі, чи ні.

Наприклад, якщо два піри з'єднані через з'єднання з високою пропускною здатністю, це означає, що вони знаходяться в одному офісі. З іншого боку, якщо два піри мають низьку пропускну здатність між собою, вони, ймовірно, знаходяться в різних офісах.

Маючи цю інформацію, піри можна згрупувати в офіси за допомогою відомих методів кластеризації, таких як розбиття графів [2, 3] або спектральна кластеризація [4].

Теоретично, протокол пліток надає Hive Streaming необхідну інформацію для побудови оверлейної мережі. Однак цей протокол не можна наївно застосовувати на практиці. У реалістичному сценарії з $N = 5000$ пірів існує приблизно $N^2 = 25$ мільйонів з'єднань для тестування. Жодна компанія не захоче мати 5 мільйонів тестів пропускної здатності, що працюють в її мережі. Навіть якщо перевантаження мережі не було проблемою, протокол пліток зайняв би дні, щоб завершитись.

Щоб протокол пліток мав сенс у реалістичних сценаріях, кількість тестів має бути мінімальною. Кожен пір повинен вимірювати пропускну здатність з невеликою частиною корпоративної мережі (тобто 10 випадкових пірів). Ця альтернатива тримає мережеві вимоги під контролем.

Незважаючи на те, що вищезазначена ідея зменшує кількість тестів у мережі, вона створює іншу проблему. Зібрана інформація ($10N$) є лише крихітною частиною порівняно з кожним можливим з'єднанням мережі (N^2). Цієї інформації може бути недостатньо для виконання кластеризації, оскільки більшість з'єднань відсутня. В результаті алгоритм кластеризації

створить кластери, які не будуть дуже точними. Якщо клієнти HIVE WebRTC використовують ці кластери для формування оверлейної мережі, буде багато пірів, з'єднаних між офісами. Ці з'єднання використовуватимуть з'єднання з низькою пропускну здатністю мережі, що призведе до потенційного перевантаження мережі.

Загалом, весь процес страждає від ступеня розрідженості інформації. Можна стверджувати, що проблему розрідженості можна вирішити, вимірюючи пропускну здатність між більшою кількістю пірів мережі (тобто 20). Це дозволило б досягти кращих результатів, оскільки вхідні дані алгоритму кластеризації отримували б більше інформації. Використовуючи більше інформації, ми очікуємо, що отримані кластери будуть більш точними. На жаль, це рішення не є достатнім з двох причин. По-перше, воно не масштабується, оскільки ми обмінюємо мережеві ресурси на кращі кластери. По-друге, результати залежать від топології корпоративної мережі. У випадках, коли топологія проста, тобто є лише два офіси, зібраної інформації може бути достатньо для отримання надійних результатів. Однак, якщо топологія складна, що зазвичай буває у великих підприємствах, то збільшення кількості вибіркового з'єднань може бути неефективним.

Проблема стає помітною, коли ми беремо до уваги динамічність мережі. У типовому випадку HIVE Streaming нові пристрої можуть приєднуватися до мережі, а існуючі можуть відключатися в будь-який час. У таких випадках інформація про плитку зростає лінійно (тобто $10N$), тоді як кількість можливих з'єднань зростає квадратично (N^2). Отже, доки наша мережа зростає, зібрана інформація про пропускну здатність ($10N / N^2 = 10 / N$) з часом стає все більш розрідженою.

Врешті-решт, єдиним прийнятним рішенням є вирішення обох проблем одночасно. З одного боку, нам потрібно мінімізувати мережевий слід протоколу плиток. З іншого боку, нам потрібно достатньо інформації про пропускну здатність, щоб створити точні кластери. На щастя, є спосіб мати обидва. Враховуючи, що кожен пір виміряв пропускну здатність з невеликою

частиною інших пірів, гарною альтернативою є використання цієї інформації для прогнозування пропускну́ї здатності для решти пірів. Такі прогнози можна використовувати як приблизну оцінку фактичної пропускну́ї здатності, яку інакше ми ніколи не змогли б виміряти. Таким чином, інформація про плітки поєднується з прогнозованою інформацією, щоб створити кращі кластери. Цей метод зберігає мережеві витрати на тому ж рівні, створюючи більш ефективний оверлей.

1.5. Оцінка пропускну́ї здатності за допомогою Dynamic Graph Representation Learning

У наведеному вище контексті прогнозування пропускну́ї здатності між одноранговими вузлами є проблемою великої важливості. Цю проблему можна описати більш формально за допомогою теорії графів. Накладну мережу можна розглядати як графік. У цьому графіку однорангові вузли представлені як вузли, а ребра позначають зв'язки між ними. Пропускна здатність з'єднань відображається на вазі країв. Завдяки динамічності події HIVE Streaming графік також є динамічним. Нові вузли можуть входити в графік, тоді як існуючі вузли можуть бути несподівано видалені. Так само вага ребер може змінюватися з часом. У цьому динамічному графіку наша мета — передбачити вагу неіснуючих ребер.

По суті, задачу прогнозування ваги можна класифікувати як задачу прогнозування зв'язку. Для динамічного графа динамічне передбачення зв'язку є завданням передбачення значення ребра, про яке немає попередньої інформації. Прогнозування зв'язків – це добре досліджена область [5, 6]. Протягом багатьох років було запропоновано безліч алгоритмів і підходів, які відповідають різним постановкам задач, типам графів і вимогам [7]. Традиційні моделі прогнозування зв'язку погано працюють на динамічних графіках, створених за допомогою протоколу gossip. Це пов'язано з такими властивостями:

1. Висока розрідженість: створений динамічний графік за своєю суттю розріджений. Прості моделі, такі як логістична регресія, призводять до неточних прогнозів, коли вхідні екземпляри переповнені нулями.

2. Висока динамічність: динамічність вхідного графіка не обмежена, що означає, що динамічний графік може швидко змінюватися в значних межах. Наприклад, під час трансляції потокового відео в прямому ефірі більшість однолітків заходять на подію протягом перших кількох хвилин. За ці кілька хвилин структура динамічного графіка h зазнає великих змін, які нелегко вловити.

3. Ненормальний розподіл: загалом, підходи, які стосуються зважених графіків, неявно припускають, що ваги відповідають (перекошеному) нормальному розподілу. У нашому випадку це припущення не відповідає дійсності. Це ускладнює завдання прогнозування, оскільки може спричинити проблеми в процесі навчання моделі.

З цих причин ми плануємо використовувати навчання динамічного представлення графів для підтримки процесу прогнозування пропускної здатності.

1.6. Цілі та методологія дослідження

Мета цього проекту - знайти рішення проблеми оцінки пропускної здатності. Ми пропонуємо переглянути існуюче дослідження вивчення динамічного представлення графів і розробити нову модель, яка адаптована до випадку використання Hive Streaming. Враховуючи цю мету, ми можемо чітко визначити три підцілі для цього проекту:

1. Вивчити ефективність існуючих найсучасніших підходів до навчання представлення динамічних графів для сценарію використання Hive Streaming.
2. Запропонувати інноваційний динамічний підхід, який відповідає потребам Hive Streaming.

3. Виконати ретельну оцінку, яка покаже, що запропонований підхід досягає більшої продуктивності порівняно з існуючими найсучаснішими динамічними підходами.

У цьому проекті ми дотримуємося емпіричної методології. Новий підхід розроблено шляхом отримання знань із різних експериментів. Грунтуючись на наших спостереженнях, запропонована модель сформована так, щоб чітко відповідати випадку використання Hive Streaming. Найважливішим аспектом є емпіричне доведення того, що запропонований підхід може перевершити існуючі динамічні підходи. Це вимагає поглибленого аналізу результатів з точки зору чітко визначених показників ефективності.

Як згадувалося вище, кінцева мета Hive Streaming — визначити, які пристрої знаходяться в кожному офісі. Щоб досягти цього, першим кроком є збір даних за допомогою протоколу пліток. Потім ці дані будуть використані як вхідні дані для моделі цього проекту, яка спрямована на отримання оцінок пропускної здатності для відсутніх з'єднань. Нарешті, вихідні дані використовуватимуться для кластеризації пристроїв в офісах. Ми заявляємо, що обсяг цього проекту обмежений оцінками пропускної здатності. Продуктивність моделі буде оцінено у вмісті передбачення посилянь. Процедура кластеризації, яка використовує модель цього проекту, є ще одним поточним завданням Hive Streaming, яке виходить за межі нашої сфери.

Висновки до розділу

У розділі проведено всебічний аналіз предметної області застосування відеострімінгу та потокового відео, що дозволяє визначити основні аспекти, переваги та проблеми, пов'язані з цією технологією. Відеострімінг як засіб для комунікації в реальному часі надає значні переваги в корпоративному та соціальному контекстах. Зокрема, можливість оперативного обміну

інформацією сприяє підвищенню ефективності бізнес-процесів і покращенню взаємодії між учасниками процесу.

Розвиток спеціалізованих сервісів для потокового відео відповідає вимогам сучасних компаній до масштабованості, безпеки та ефективності. Водночас залишаються виклики, пов'язані з обмеженням пропускної здатності мереж і оптимізацією передачі даних. Виявлено, що оверлейні мережі, які часто використовуються для розповсюдження потокового контенту, стикаються з проблемами ефективного маршрутизації, затримок і перевантажень.

Використання методів динамічного графового навчання дозволяє підвищити точність оцінки пропускної здатності мережі, що є критично важливим для оптимізації потокового відео. Сформульовані основні цілі дослідження, які передбачають удосконалення методів передачі потокового відео, мінімізацію затримок і підвищення ефективності мережевих протоколів. Обрана методологія орієнтована на аналіз існуючих технологій та розробку нових підходів для їх покращення.

Узагальнення проведеного аналізу вказує на важливість розробки інноваційних методів для вирішення сучасних проблем потокового відео та вдосконалення наявних технологій.

РОЗДІЛ 2. МОДЕЛІ ТА АЛГОРИТМИ ДИНАМІЧНОГО ТА СТАТИЧНОГО НАВЧАННЯ НА ГРАФОВИХ ПРЕДСТАВЛЕННЯХ

2.1. Алгоритми передачі даних, які використовуються в трансляції відео в реальному часі

Розглянемо основні протоколи та алгоритми трансляції відео в реальному часі, зокрема про протоколи RTP, RTSP та WebRTC:

1. RTP (Real-time Transport Protocol)

RTP - це стандартний протокол для передачі аудіо- та відеоданих в реальному часі через мережу IP. Він забезпечує доставку пакетів даних з урахуванням часових міток, що дозволяє відтворювати медіа-контент синхронно.

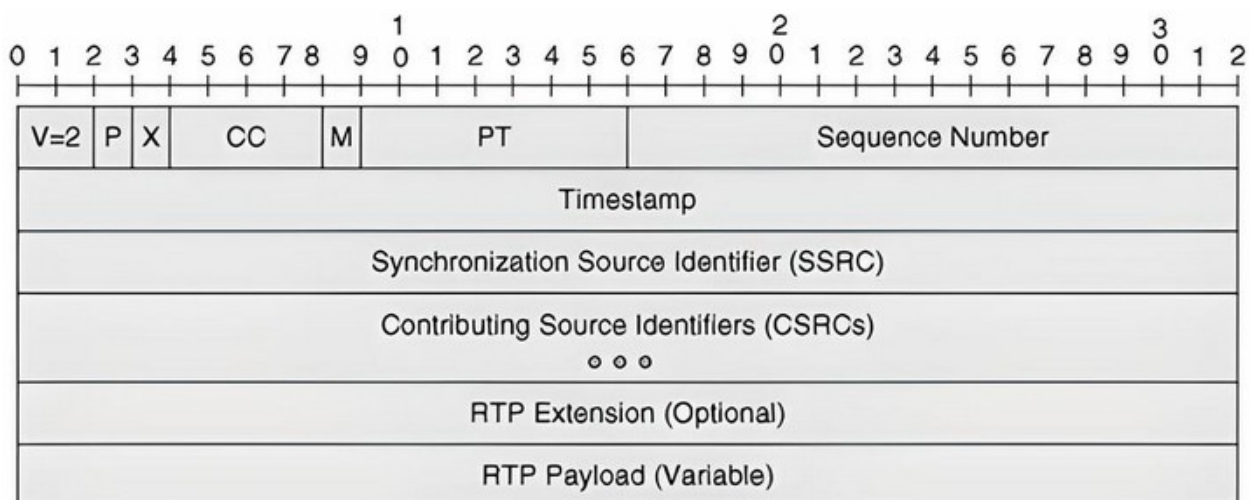


Рис. 2.1. Структура RTP пакету показана в 32-бітних (4-байтових) блоках

Ключові особливості RTP:

- Передача в реальному часі: RTP оптимізований для доставки медіаданих з мінімальною затримкою.
- Синхронізація: Часові мітки в RTP-пакетах допомагають синхронізувати аудіо та відео потоки.

- Ідентифікація джерела: RTP-пакети містять інформацію про джерело медіаданих.

- Контроль якості: RTP працює в парі з RTCP (Real-time Transport Control Protocol), який використовується для моніторингу якості передачі та надання зворотнього зв'язку.

2. RTSP (Real Time Streaming Protocol)

RTSP - це протокол прикладного рівня, який використовується для керування потоками медіаданих. Він дозволяє клієнтам (наприклад, відеоплеєрам) підключатися до медіасерверів, запускати та зупиняти відтворення, перемотувати відео тощо.

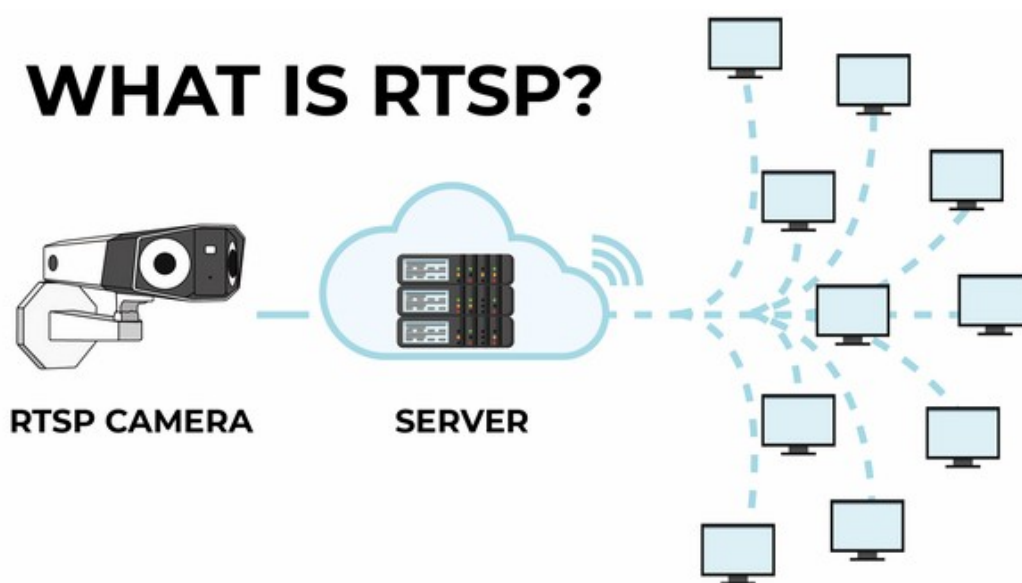


Рис. 2.2. Зв'язок RTSP між клієнтом і сервером

Ключові особливості RTSP:

- Керування сеансом: RTSP надає команди для встановлення, керування та завершення сеансів потокової передачі.

- Взаємодія з сервером: RTSP дозволяє клієнтам надсилати запити на сервер, наприклад, для отримання інформації про медіапотік або зміни параметрів відтворення.

- Незалежність від транспорту: RTSP може використовувати різні транспортні протоколи, такі як RTP, UDP або TCP, для передачі медіаданих.

3. WebRTC (Web Real-Time Communication)

WebRTC - це відкрита технологія, яка дозволяє здійснювати аудіо- та відеозв'язок безпосередньо в веб-браузері. Вона використовує peer-to-peer (P2P) з'єднання для передачі даних, що дозволяє зменшити затримку та навантаження на сервер.

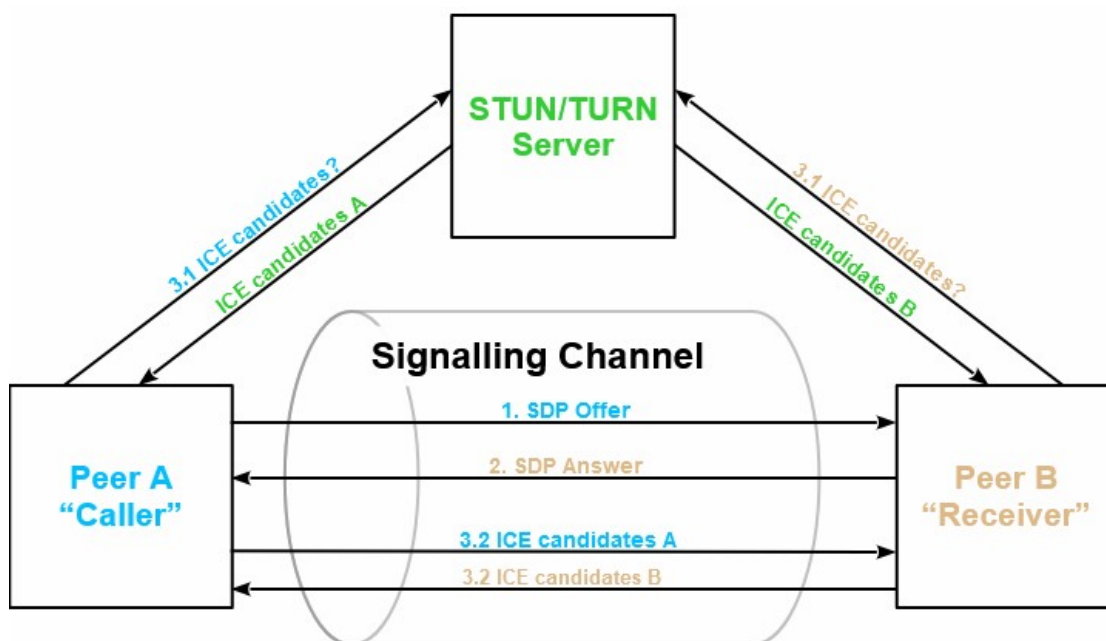


Рис. 2.3. Однорангове підключення WebRTC

Ключові особливості WebRTC:

- P2P з'єднання: WebRTC дозволяє браузерам обмінюватися даними напряму, без посередників.

- Низька затримка: P2P з'єднання мінімізує затримку, що важливо для відеоконференцій та інших інтерактивних застосунків.

- Вбудована підтримка в браузерах: WebRTC підтримується в більшості сучасних веб-браузерів, що робить його легким у використанні.

- Шифрування: WebRTC використовує шифрування для захисту даних під час передачі.

Отже, RTP, RTSP та WebRTC - це важливі протоколи для трансляції відео в реальному часі. RTP забезпечує доставку медіаданих, RTSP дозволяє керувати потоками, а WebRTC надає можливість здійснювати P2P зв'язок в браузері. Вибір протоколу залежить від конкретних потреб та вимог до застосування.

2.2. Концепція графового представлення навчання

Графове представлення навчання – це концепція, яка створює латентне (або просто альтернативне) подання оригінального графа. Потім це представлення використовується як вхідні дані для виконання інших завдань машинного навчання, таких як прогнозування зв'язків, класифікація, кластеризація тощо. Безсумнівно, було показано, що приховане представлення може досягти значного приросту продуктивності [8].

Ключова ідея вивчення графового представлення полягає в тому, що розріджена інформація кожного вузла перетворюється в інше представлення, яке називається вбудовуванням вузла. Щоб бути корисними для інших завдань машинного навчання, вбудовані вузли повинні мати певні властивості:

- Вбудовані вузли знаходяться в меншому розмірному просторі. На графі з N вузлами будь-який вузол можна виразити через його зв'язки з рештою мережі. Отже, будь-який вузол можна представити вектором $v \in \mathbb{R}^N$. У випадку вбудовування вузлів для прихованого представлення потрібно менше розмірів. Вузол представлений його вектором вкладення $z \in \mathbb{R}^d$, де $d \ll N$. Ці d розмірності можна інтерпретувати як приховані особливості вхідного графа в d -вимірному просторі.
- Вбудовані вузли щільні. На відміну від вихідного представлення вузла, яке містить безліч нулів, приховане представлення використовує кожен вимір прихованого простору.

- Вбудовані вузли повинні відображати структурні властивості вхідного графа. Іншими словами, вихідна інформація про кожен вузол v і відповідна вбудована інформація z повинні демонструвати однакоку поведінку. Наприклад, якщо існує зв'язок між вузлами i, j у вихідному графі, то їх вкладення z_i, z_j має відображати цей зв'язок. У цьому ж напрямку, якщо вузли i, j не з'єднані, то їх вкладення z_i, z_j має виражати цей розрив.

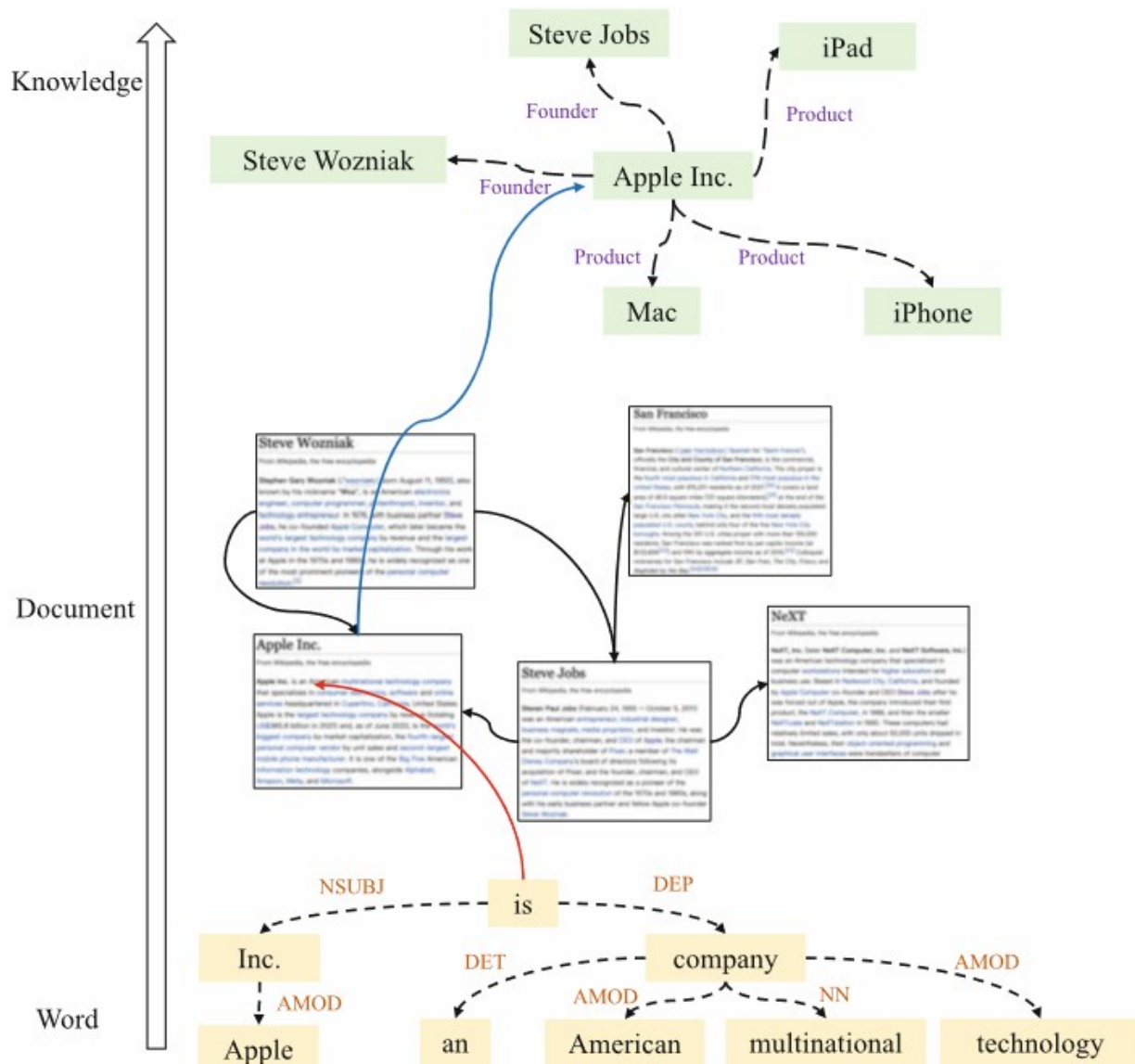


Рис. 2.4. Приклад текстових графів на різних рівнях. Документи, використані на рисунку, отримані з офіційного веб-сайту Вікіпедії

Загалом, навчання подання графів можна застосувати до графів зі статичною та динамічною поведінкою. Це ділить область дослідження на дві

менші. Навчання динамічним представленням графів – це область досліджень, яка стосується динамічних графіків, тоді як вивчення статичного представлення графів стосується графіків, які не змінюються (тобто залишаються статичними з часом). Для кожного типу існуючі підходи можна далі розділити на три основні категорії:

- Матричні підходи факторизації.
- Підходи випадкового блукання.
- Підходи глибокої нейронної мережі.

Підходи матричної факторизації створюють вкладення шляхом декомпозиції вхідного графа. Ці підходи працюють подібно до методів зменшення розмірності (тобто PCA). Матриця суміжності вхідного графа розкладається на добуток матриць низької розмірності. Ці маловимірні матриці містять багату інформацію про вузли графа, які можна використовувати як вкладення. Підходи випадкових блукань створюють вбудовування, покладаючись на інформацію, отриману випадковими блуканнями. По суті, вузли, які зустрічаються разом у випадкових блуканнях, вважаються подібними. Тому їх вбудовування мають форму, щоб відображати структурну подібність. Нарешті, підходи до нейронних мереж використовують складні нелінійні моделі для створення вставок. Ці моделі вчать створювати вбудовування, намагаючись реконструювати вихідну мережу. Моделі навчаються шляхом вимірювання похибки між оригінальною мережею та реконструйованою. Зводячи до мінімуму цю помилку, моделі можуть вивчати суттєві вбудовування.

Історично склалося так, що більшість підходів до вивчення представлення графів застосовуються до статичних графів. Однак графіки, створені на основі сценаріїв реального світу, за своєю суттю динамічні. Динамічні графіки зазвичай виражаються у вигляді послідовності знімків, кожен з яких є статичним виглядом графіка. Теоретично всі статичні підходи можна застосувати до динамічних графіків, застосовуючи відповідний підхід до кожного знімка окремо. Однак ця методологія вкрай неефективна,

оскільки кожен знімок обробляється незалежно. Тобто процес генерації вбудовувань починається з нуля для кожного знімка, без повторного використання будь-якої інформації з попередньої ітерації.

Протягом останніх кількох років було запропоновано різноманітні динамічні підходи. У динамічних підходах вхідними даними є серія знімків графіка, які визначають еволюцію графіка в часі. Динамічні підходи побудовані таким чином, щоб вони могли охопити цю еволюцію. Як правило, фіксація динамічності графіка ускладнює підходи. Основна відмінність від статичних підходів полягає в тому, що динамічні підходи здатні слідувати еволюції динамічного графіка. Замість того, щоб повторно створювати вбудовані вузли для кожного знімка, існуючі вбудовані адаптуються до поточного стану графіка. Це призводить до більш ефективних підходів, які досягають кращих результатів.

2.3. Навчання статичного представлення графів

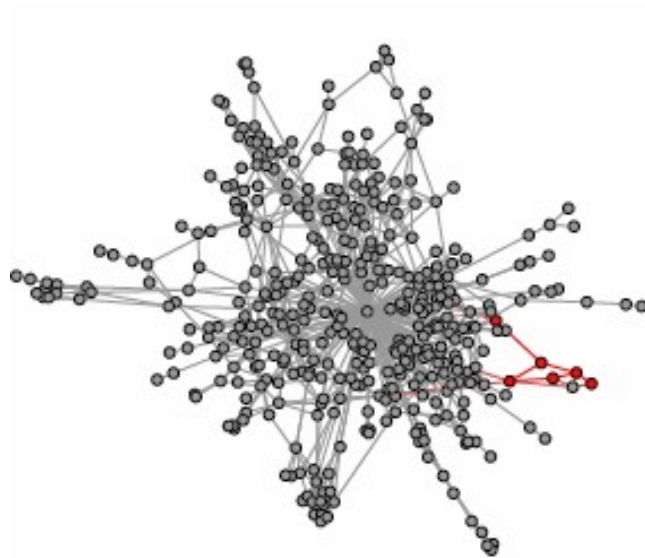
2.3.1. Методи факторизації матриць

Історично першими підходами до навчання представлення графів були методи факторизації матриць. В [10] виконали факторизацію матриці на матриці модульності. Матриця модульності - це матриця, яка вимірює різницю між вхідним графом та випадковим рівномірним графом з тією ж послідовністю ступенів. Потім вони витягли вбудовування вузлів з верхніх власних векторів матриці модульності та використали їх як ознаки для виконання логістичної регресії. В дослідженні [11] виконали регуляризовану факторизацію гаусової матриці. Їхньою метою було використати просту модель скалярного добутку для навчання вбудовування вузлів. Модель фокусується на подібності між вузлами, яка визначається матрицею суміжності графа. Подібно до [11], в [12] запропонували, що вбудовування вузлів повинно спиратися на різні степені матриці суміжності, оскільки сама матриця суміжності містить лише локальну інформацію. Використовуючи k -

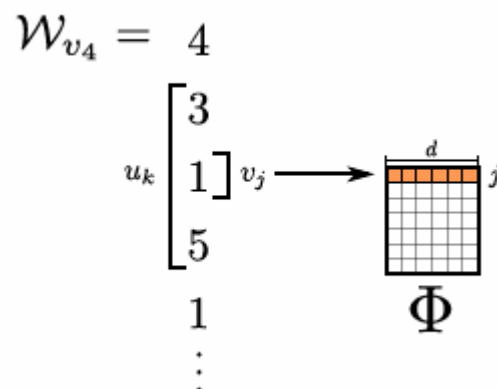
степені матриці суміжності, їхній метод може вивчати глобальну структурну інформацію графа. Нарешті, в [13] використали різні вимірювання близькості (тобто індекс Каца, Adamic-Adar та інші), які фіксують локальну або глобальну структуру графа. Використовуючи множинні близькості, їхній підхід факторизації може вивчати різні структурні зв'язки.

2.3.2. Методи випадкового блукання

Друга категорія навчання представлення графів відноситься до методів випадкового блукання. В [14] запропонували DeepWalk, перший алгоритм випадкового блукання, який здатний вивчати вбудовування вузлів.



а) Генерація випадкового блукання



б) Відображення представлення

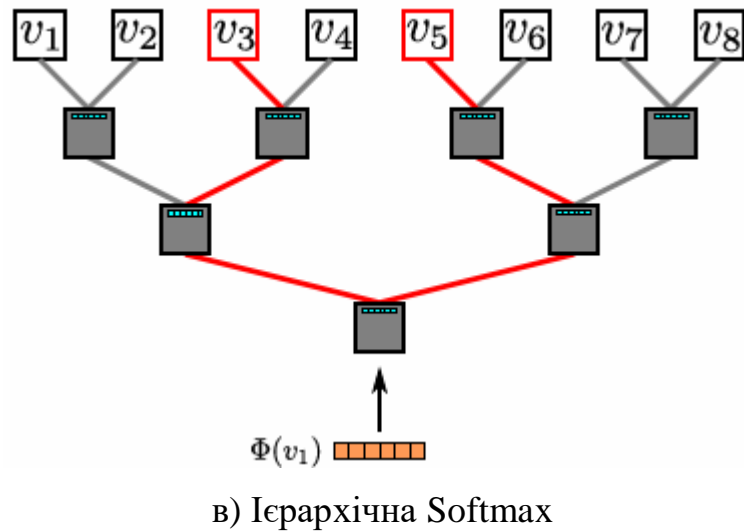


Рис. 2.5. Огляд DeepWalk.

На рис. 2.5. показано як відбувається ковзання вікном довжиною $2w + 1$ по випадковому блуканню Wv_4 , відображаючи центральну вершину v_1 на її представлення $\Phi(v_1)$. Ієрархічна Softmax розкладає $\Pr(v_3 | \Phi(v_1))$ та $\Pr(v_5 | \Phi(v_1))$ на послідовності розподілів ймовірностей, що відповідають шляхам, які починаються в корені та закінчуються в v_3 та v_5 . Представлення Φ оновлюється, щоб максимізувати ймовірність спільної появи v_1 з його контекстом $\{v_3, v_5\}$.

Натхненний обробкою природної мови, DeepWalk виконує невеликі випадкові блукання на вхідному графі. Потім ці блукання розглядаються як речення в моделі Skip-Gram [15]. Модель Skip-Gram вивчає вбудовування вузлів, які відображають ймовірність спільної появи вузлів. Пізніше в [16] запропонували Node2Vec, який є розширенням DeepWalk. Автори підкреслили, що DeepWalk виконує жорстку стратегію пошуку, яку не можна налаштувати. Вони запропонували параметризовану стратегію випадкового блукання, яка адаптується до унікальної топології графа. Node2Vec може фіксувати як відносини сусідства, так і структурні відносини.

Окрім цих основних підходів, було запропоновано різноманітні розширення випадкового блукання. В [17] заявили, що, хоча більшість моделей неявно припускають, що векторний простір є евклідовим, багато складних моделей можна краще описати за допомогою гіперболічної геометрії. Таким чином, вони запропонували підхід випадкового блукання в гіперболічному просторі. В дослідженні [18] запропонували *Struc2Vec*, який фокусується на фіксації структурної ідентичності вузлів. *Struc2Vec* обчислює контекстний граф, який є k -шаровим графом, що містить структурні подібності в k -розмірних шляхах. Вбудовування вузлів створюється шляхом запуску *Node2Vec* на контекстному графі. Нарешті, в [19] запропонували *LINE*, спробу вловити близькість першого та другого порядку. На відміну від *DeepWalk*, який поєднує подібності у випадкові блукання, близькості *LINE* явно визначені окремо. Хоча *LINE* не є підходом випадкового блукання, він вважається подібним до них через запропоновану техніку вибірки ребер.

2.3.3. Підходи нейронних мереж

Остання категорія цього розділу - це підходи глибоких нейронних мереж. Враховуючи останні досягнення в нейронних мережах, було розроблено безліч підходів для навчання вбудовування вузлів. В [20] запропонували *SDNE*, напівнавчальну модель, яка створює вбудовування вузлів. По суті, *SDNE* - це глибокий автокодер, який навчається реконструювати околиці кожного вузла. Під час навчання середній шар автокодера витягує найважливіші ознаки кожного вузла, які використовуються як вбудовування. Подібним чином в [21] запропонували *DNGR*, який є автокодером з шумозаглушенням, який отримує матрицю *PPMI* [22] як вхідні дані. Матриця *PPMI* будується шляхом застосування моделі випадкового серфінгу з перезапусками на вхідній матриці.

На рисунку 2.6 показана модель *DNGR*, що складається з трьох основних кроків. По-перше, тут введена модель випадкового блукання для захоплення структурної інформації графа та генерації матриці ймовірнісної

спільної появи. Далі тут обчислюється матриця PPMI на основі матриці ймовірнісної спільної появи. Після цього для навчання низьковимірних представлень вершин використовується багатoshаровий автокодер з шумозаглушенням.

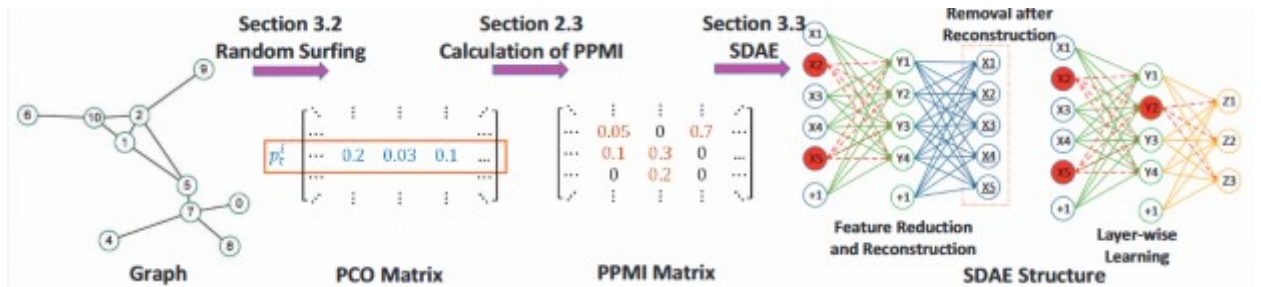


Рис. 2.6. Основні компоненти: випадкове блукання, обчислення матриці PPMI та зменшення розмірності ознак за допомогою SDAE

В [23] представили GraphSAGE, індуктивний алгоритм, який виконує агрегацію ознак сусідства. Ці агрегації є за своєю суттю складними, оскільки вони застосовуються до k -крокових шляхів, використовуючи різні функції агрегації. GraphSAGE також можна розглядати як динамічний підхід, оскільки ту саму ідею можна застосувати до динамічних графів.

Нова концепція навчання представлення графів народилася, коли в [24] застосували згорткові мережі на графах, щоб вивчати вбудовування вузлів. Вони представили Graph Convolutional Networks (GCN), нейронну мережу на основі графів, яка може вивчати складні зв'язки між локальною структурою графа та ознаками його вузлів. Незабаром після появи GCN було запропоновано різноманітні розширення GCN. Варіаційні графові автокодери (VGAE) [25], які також були запропоновані і є, мабуть, найвидатнішим розширенням GCN. VGAE використовує варіаційний кодер, який навчається шляхом вибірки вбудовування вузлів. В [26] заявили, що багато графів реального світу мають негативні ребра. З цієї причини вони запропонували модель GCN зі знаком, яка може обробляти інформацію про ребра зі знаком.

В [27] стверджували, що традиційна модель GCN обмежена у вивченні дельта-операторів. Через це вони пропонують новий згортковий шар графа, який використовує множинні степені входної матриці. І останнє, але не менш важливе, в [28] представили N-GCN, модель з кількома GCN. N-GCN здатний вивчати різні близькості за допомогою кількох моделей GCN.

Ще однією архітектурою нейронної мережі для навчання представлення графів є генеративні змагальні мережі [29]. Коротко кажучи, GAN складаються з двох моделей, генератора та дискримінатора, які навчаються взаємно. В [30] представили GraphGAN, першу архітектуру GAN для навчання ознак. Генератор GraphGAN виконує випадкові блукання для створення фальшивих ребер, тоді як дискримінатор навчається класифікувати ребра як справжні або фальшиві. NetGAN [31] подає випадкові блукання в LSTM [32] для вивчення топології та структури входу. KBGAN [33] вивчає вбудовування вузлів, використовуючи трійки графа як навчальні екземпляри. На жаль, більшість (якщо не всі) моделей GAN для навчання представлення графів вимагають попередньо навченого вбудовування вузлів як входних даних через їхню повільну збіжність.

Натхненні успіхом механізмів уваги, автори в [34] запропонували Graph Attention Networks (GAT), перший підхід, який використовує механізм уваги для створення вбудовування вузлів. Основна ідея GAT полягає в тому, що латентне представлення вузла можна вивчити, звертаючи увагу на його околиці. Модель уваги може зрозуміти, які сусіди є важливими, і створює вбудовування, які це відображають. В [35] запропонували Goat, який є контекстно-залежним алгоритмом, що поєднує комунікацію пліток з механізмом взаємної уваги. Комунікація пліток використовується для створення кількох представлень вузлів, які залежать від контекстів, частиною яких є кожен вузол. Ці різні представлення повинні збігатися для вузлів, які співіснують в одному контексті. Для цього механізм взаємної уваги гарантує, що вузли в межах одного контексту мають подібні представлення.

2.4 Концепція навчання динамічного представлення графів

2.4.1. Методи факторизації матриць

В останні кілька років багато зусиль було докладено до навчання представлення динамічних графів. Методи факторизації матриць стикаються зі складними проблемами, коли йдеться про динамічні графи. Оскільки виконання всього алгоритму факторизації для кожного знімка графа є обчислювально неефективним, динамічні підходи факторизації повинні бути інкрементними. Багато інкрементних підходів SVD були запропоновані в [36], але вони призводять або до високої складності, або до низької точності. Основна проблема цієї категорії полягає в тому, коли перезапускати SVD. А [37] використовували втрати реконструкції як спосіб вимірювання того, коли SVD слід перезапускати. В тому ж контексті в [38] відстежували різницю між втратами реконструкції останнього оновлення та теоретичними мінімальними втратами SVD.

2.4.2. Методи випадкового блукання

Враховуючи успіх статичних підходів випадкового блукання, в [39] представили інноваційну ідею для динамічних випадкових блукань. У своїй роботі вони припустили, що часова деталізація є важливим аспектом динамічних графів. З цієї причини вони запропонували StreamWalk та "Онлайн навчання подібності другого порядку", два динамічні підходи, які отримують безперервний потік ребер як вхідні дані. StreamWalk зосереджений на близькості першого порядку, тоді як другий підхід зосереджений на близькості другого порядку.

2.4.3. Підходи на основі нейронних мереж

Коли йдеться про архітектури нейронних мереж для навчання представлення динамічних графів, було розроблено багато моделей. В [40]

запропонували DynamicTriad, алгоритм, який моделює еволюцію мережі за допомогою триад. Зокрема, DynamicTriad вивчає динаміку мережі, спостерігаючи, як відкриті триади перетворюються на закриті. В [41] реалізували DynGEM, розширюваний глибокий автокодер, в якому розмір та кількість шарів можуть адаптуватися до динамічного входу. Загалом кажучи, у міру того, як вхідний граф збільшується в розмірі, автокодер стає ширшим і глибшим. В [42] також запропонували серію нейронних мереж, dyngraph2vecAE, dyngraph2vecRNN та dyngraph2vecAERNN. Перша модель - це простий автокодер, який використовує функцію втрат, яка може вивчати еволюцію графа. Для більш ефективного часового навчання друга модель використовує розріджено з'єднані комірочки RNN. Зокрема, друга модель складається з комірок LSTM, які здатні фіксувати довгострокові залежності. Третя модель - це комбінація автокодерів та LSTM, яка має на меті досягти більшої продуктивності.

В рамках GCN, в [43] запропонували EvolveGCN (рисунок 2.7). Це нова модель, яка поєднує GCN та RNN на динамічних графах. Замість безпосереднього навчання GCN, EvolveGCN навчає RNN створювати ваги для GCN. Отже, RNN вивчають еволюцію графа з часом. Нещодавно в [44] представили DynVGAE (рисунок 2.8), варіаційний графовий автокодер, який застосовується до динамічних графів. DynVGAE складається з серії VGAE, кожен з яких навчається на одному знімку графа. Новизна DynVGAE полягає в тому, що ваги спільних вузлів між двома послідовними знінками переносяться на наступний VGAE, що прискорює процес навчання.

Крім того, для динамічних графів були розроблені динамічні архітектури GAN. В [45] запропонували DynGAN та DynGAN-LSTM. DynGAN - це просто статичний GraphGAN, який навчається на кількох знімках графа. Подібно до DynVGAE, спільні вузли між послідовними знімками графа зберігають свої вбудовування на наступному етапі навчання. Автори стверджують, що вищезгадана модель не може фіксувати часові залежності між знімками графа. Згодом DynGAN-LSTM поєднує GraphGAN з

комірками LSTM, щоб подолати цю проблему. Далі DynGraphGAN [46] використовує підмножини оригінального графа для вивчення властивостей в околах вузлів. Ці підмножини генеруються з випадкових блукань за допомогою Node2Vec. Щоб фіксувати як структуру, так і еволюцію мережі, дискримінатор DynGraphGAN містить дві моделі, одну для залежностей близькості та одну для часових залежностей відповідно.

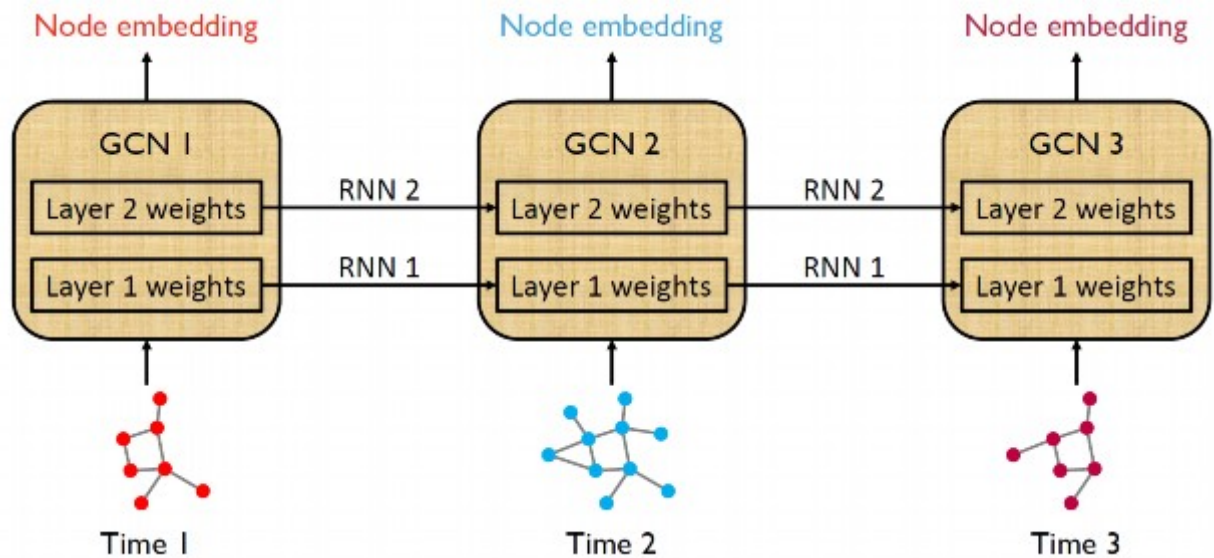


Рис. 2.7. Архітектура EvolveGCN

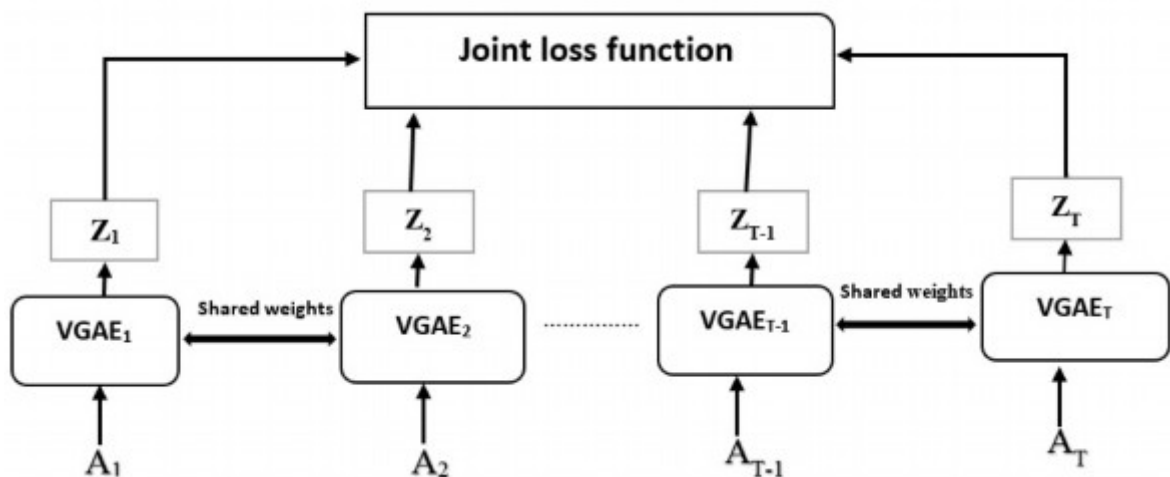


Рис. 2.8. Архітектура DynVGAE

Нарешті, ідея механізмів уваги була застосована до динамічних графів. DySAT [47] розширює статичну модель GAT до динамічної. Він

використовує механізм уваги, щоб вивчати не лише локальне сусідство кожного вузла, але й зміни в межах сусідства в динамічному середовищі. Таким чином, на створені вбудовування вузлів впливають як структура графа, так і його еволюція з часом.

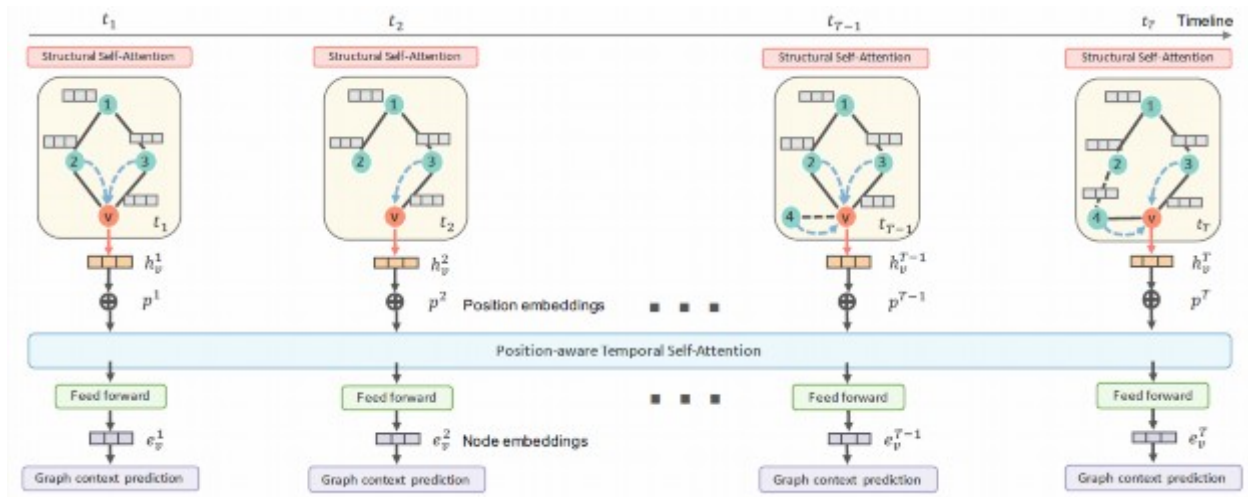


Рис. 2.9. Нейронна архітектура DySAT

На рис. 2.9 використовуються шари структурної уваги, за якими слідує шар часової уваги. Пунктирні чорні лінії вказують на нові зв'язки, а пунктирні сині стрілки позначають структурну увагу, що базується на сусідах.

Висновки до розділу

В даному розділі розглянуто сучасні моделі та алгоритми, що використовуються для навчання на графових представленнях, які забезпечують ефективну обробку даних у різних умовах, зокрема для трансляції відео в реальному часі. Технології передачі даних, що використовуються для відеострімінгу, базуються на поєднанні адаптивних мережевих протоколів і алгоритмів оптимізації трафіку. Це дозволяє знижувати затримки, підвищувати якість передачі та забезпечувати

стабільність потокового відео навіть за умов мінливого мережевого середовища.

Графове представлення стало важливим інструментом для моделювання складних взаємозв'язків між об'єктами. Ця концепція дозволяє ефективно структурувати інформацію та забезпечує можливості для виявлення прихованих залежностей у даних. Концепція навчання динамічного представлення графів орієнтована на моделювання графів, що змінюються з часом, і забезпечує адаптацію алгоритмів до динамічних змін у структурі даних.

Розділ показує, що використання графових моделей і алгоритмів навчання відкриває нові можливості для оптимізації систем передачі даних у реальному часі. Зокрема, статичні підходи демонструють високу ефективність для стаціонарних задач, тоді як динамічні моделі забезпечують адаптацію до змін у мережевих умовах. Об'єднання цих методів дозволяє створювати інноваційні рішення для підвищення продуктивності відеострімінгу та інших пов'язаних технологій.

РОЗДІЛ 3. МЕТОДОЛОГІЯ, МОДЕЛІ ТА МЕТОДИ ДИНАМІЧНОГО НАВЧАННЯ НА ГРАФОВИХ ПРЕДСТАВЛЕННЯХ СТРІМІНГУ

3.1. Представлення методології

У цьому розділі ми опишемо методологію роботи. Спочатку ми пояснимо процес дослідження, який містить кроки, які ми дотримувалися для реалізації проекту. Потім ми зосередимо увагу на вхідних даних. Ми пояснимо, як дані були зібрані, проаналізовані та оброблені.

Дана робота розпочалась з ретельного вивчення літератури щодо вивчення представлення графів. Важливо було почати зі статичних підходів, а потім переходити до динамічних. Це тому, що динамічні підходи складаються з десяти ідей розширення, взятих із статичних моделей. Після вивчення широкого спектру моделей було сформовано конкретне підґрунтя досліджуваної області. Маючи цей фон, ми могли б перейти до наступного завдання, яке передбачало емпіричну оцінку існуючих динамічних підходів. Для цієї частини ми повинні були визначити, які моделі підходять для нашого випадку. Ця емпірична оцінка була непростим процесом. Було докладено багато зусиль для використання існуючих моделей у наших сценаріях потокового передавання. У деяких випадках це призвело до внесення невеликих змін у реалізацію моделі, наприклад зміни функції втрат, щоб модель могла працювати з нашими даними.

На цьому етапі ми визначили, які підходи можна застосувати до нашого сценарію. Потім ці моделі були випробувані; ми оцінювали їхню ефективність, використовуючи реальні дані з подій потокового відео в прямому ефірі. За допомогою цих експериментів ми могли виявити сильні та слабкі сторони існуючих найсучасніших підходів. Ця інформація є дуже важливою, оскільки вона показує, у яких випадках і за яких обставин ефективні поточні підходи. Це також показує, що заважає існуючим моделям давати кращі результати. Маючи ці знання, ми створили нову модель. Нова

модель була розроблена таким чином, щоб вона перейняла деякі сильні сторони попередніх моделей, але уникла їхніх відповідних слабких сторін, використовуючи інший дизайн. На завершення була оцінена продуктивність нової моделі. Наша оцінка порівнює результати нашого підходу з результатами двох найсучасніших моделей. На основі цих експериментів ми зробили висновки щодо проекту та оцінили, чи виконує цей проект свої початкові цілі.

3.1.1. Збір та опис даних

Тестові дані містять вимірювання пропускної здатності під час двох потокових подій, по одному від кожного клієнта. Вони були повністю анонімні, перш ніж їх використали для цього проекту. Процес анонімізації видалив усі ідентифікатори, які могли призвести до потенційного порушення конфіденційності.

Розгорнутий аналіз двох наборів даних можна знайти в кінці цього розділу. Наразі важливо показати структуру даних. Кожен набір даних складається з серії знімків графіків. Кожен знімок містить інформацію про пропускну здатність для певного моменту часу. По суті, це список зважених країв, які описують пропускну здатність між пристроями. Ось приклад:

```
source target weight
110      200      0.9
123      45       0.1
18       89       0.05
```

Тут ми можемо побачити три записи знімка графіка. Перший стовпець відповідає першому піру мережі, який підключений до піру другого стовпця. Третій стовпець вказує пропускну здатність з'єднання, яка нормалізована між $(0, 1]$

Наші експерименти оцінюють продуктивність різних моделей машинного навчання. Зокрема, оцінюється здатність моделей створювати

прогнози пропускної здатності для мереж клієнтів. Усі випробувані моделі мають чітко визначені засоби для виконання таких прогнозів. Ми суворо дотримуємося цих методологій і робимо прогнози. Потім, ґрунтуючись на цих передбаченнях і вихідній інформації мережі, ми вимірюємо помилку, яку створює кожна модель. Якою б простою вона не була, модель, яка призводить до меншої помилки передбачення, вважається найбільш придатною для вирішення проблеми. Щоб уникнути проблем із відтворюваністю, вихідний код наших експериментів є загальнодоступним.

3.1.2. Аналіз даних

Аналіз даних цього проекту включає статистику та кодування. Статистика важлива для першого кроку аналізу; щоб зрозуміти природу зібраних даних. Необхідно проаналізувати, якими властивостями володіють наші дані. Це тому, що ці властивості є вирішальними для розробки нової моделі. Нам потрібно зрозуміти не тільки проблему, яку ми вирішуємо, але й вхідні дані, з якими ми працюємо. Коли це буде зроблено, проект потребує багато кодування. Спочатку процес наскрізного навчання існуючих моделей був адаптований до нашої проблеми. У більшості випадків надавався код моделі, і ми використовували цей код. По-друге, нашу модель потрібно було розробляти з нуля. Ми реалізували низку нових ідей, результатом яких стала запропонована модель цієї дипломної роботи.

Окрім усіх реалізацій моделі, необхідно було також реалізувати експериментальний дизайн. Для цієї частини ми розробили сценарії, які справедливо порівнюють продуктивність кожної моделі з точки зору чітко визначених показників помилок. У цих сценаріях усі моделі отримують однакові вхідні дані, вони навчаються за допомогою однакових налаштувань параметрів (або якомога ближче), і, нарешті, їх продуктивність оцінюється в тих самих тестах. Ми вважаємо, що процес оцінки максимально чесний, оскільки єдина відмінність між протестованими моделями полягає в їх

незагальних параметрах. Для цих параметрів ми встановлюємо найбільш підходяще значення за допомогою стратегії пошуку сітки.

3.2. Прогнозування пропускну́ї здатності за допомогою навчання представлення динамічних графів

У цьому розділі ми пропонуємо гнучке рішення, яке не тільки вирішує проблему, але й перевершує всі попередні підходи. Для цього необхідно формально визначити проблему. Цей розділ починається з визначення проблеми та мети за допомогою теорії графів. Потім він переходить до детального представлення нового рішення. Таблиця 3.1 містить короткий опис усіх символів, які використовуються в цьому розділі.

Таблиця 3.1.

Символи, що використовуються

Symbol	Explanation
Input Symbols	
\mathcal{G}_t	graph snapshot of \mathcal{G} at time t
\mathcal{V}_t	the node set of \mathcal{G}_t
\mathcal{E}_t	the edge set of \mathcal{G}_t
\mathcal{W}_t	the weight set of \mathcal{G}_t
$v_t(i)$	node i of \mathcal{V}_t
C_t	cardinality of \mathcal{V}_t
$e_t(i, j)$	edge i, j in \mathcal{E}_t
$\mathcal{N}_t(i)$	neighbourhood of node i
$w_t(i, j)$	weight of $e_t(i, j)$
\mathbf{A}_t	the adjacency matrix of \mathcal{G}_t
Graph Attention Mechanism	
$\mathbf{z}_t(i)$	node embeddings of i
d	dimensions of node embeddings
$h_t(i, j)$	attention coefficient of j to i
\mathbf{a}_t	trainable vector of the attention mechanism
\mathbf{H}_t	trainable matrix of the attention mechanism
$c_t(i, j)$	normalized attention coefficient of j to i
k	number of heads in the attention model

Decoder	
\mathcal{T}	weight normalization function
f	BoxCox transformation function
λ	BoxCox transformation parameter
r	rescale function
$\widetilde{\mathcal{W}}_t$	normalized weight set of \mathcal{G}_t
$\widetilde{w}_t(i, j)$	normalized weight of $e_t(i, j)$
$p_t(i, j)$	predicted value of $\widetilde{w}_t(i, j)$
μ	root parameter of the decoder
Training	
$\delta_t(i, j)$	the error of prediction $p_t(i, j)$
$B_t(i, j)$	the error multiplier of $e_t(i, j)$
β	multiplier for high bandwidth connections
γ	threshold for high bandwidth connections
Δ_t	total prediction error
l	size of window training

3.2.1. Визначення проблеми

Під час події онлайн стримингу клієнти WebRTC утворюють мережу. Цю мережу можна представити як неорієнтований зважений динамічний граф $G = (G_1; G_2; \dots; G_T)$, де $G_t \in G$ - це знімок графа в момент часу t . Кожен знімок можна додатково визначити як $G_t = (V_t; E_t; W_t)$, де V_t - множина вершин, E_t - множина ребер, а W_t - множина ваг, що відповідають E_t . Клієнт Hive WebRTC і відповідає вершині $v_t(i) \in V_t$. Для кожного знімка G_t загальна кількість вершин позначається як $C_t = |V_t|$. З'єднання між двома клієнтами Hive WebRTC $i; j$ відповідає ребру $e_t(i; j) \in E_t$. Окіл $N_t(i)$ вершини $v_t(i)$ визначається як множина вузлів $v_t(j)$, які з'єднані з $v_t(i)$, $N_t(i) = \{v_t(j) : e_t(i; j) \in E_t\}$. Кожне ребро $e_t(i; j)$ пов'язане з вагою $w_t(i; j) \in W_t$, яка визначає пропускну здатність між $v_t(i); v_t(j)$. Оскільки G неорієнтований, $w_t(i; j) = w_t(j; i)$. Передбачається, що $w_t(i; j) \in (0; 1]$. Малі значення $w_t(i; j)$ представляють з'єднання з низькою пропускну здатністю, тоді як високі значення представляють з'єднання з високою пропускну здатністю.

Динамічний граф G здатний адаптуватися до мережі потокового відео в реальному часі. У будь-який момент часу t будь-яка кількість клієнтів HIVE WebRTC може приєднатися до події, додаючи таким чином нові вузли до V_t . Аналогічно, будь-яка кількість клієнтів HIVE WebRTC може відключитися від події. Видалення вимагає видалення відповідних вузлів з V_t , а також усіх їхніх ребер з E_t та W_t . Коли два клієнти HIVE WebRTC $i; j$ з'єднуються один з одним у момент часу t , нове ребро $e_t(i; j)$ додається до E_t , а нова вага $w_t(i; j)$ додається до W_t . Нарешті, будь-яке існуюче з'єднання між двома клієнтами HIVE WebRTC $i; j$ може змінитися, що змінює $w_t(i; j)$, де t - це позначка часу зміни.

Кожен знімок графа G_t також може бути представлений матрицею. Матричне представлення - це матриця суміжності $A_t \in R^{C_t \times C_t}$, яка визначається наступним чином:

$$A_t(i, j) = \begin{cases} w_t(i, j) & \text{if } e_t(i, j) \in \mathcal{E}_t \\ 0 & \text{if } e_t(i, j) \notin \mathcal{E}_t \end{cases}$$

У такій динамічній мережі наша мета - передбачити пропускну здатність між клієнтами HIVE WebRTC, для яких немає попередньої інформації. Зокрема, враховуючи динамічний граф G , метою цього проекту є передбачення ваги $w_0t(i; j)$ неіснуючого ребра $e_0t(i; j)$ між будь-якими вузлами $v_t(i); v_t(j)$, використовуючи інформацію попередніх знімків $G_1; G_2; \dots; G_{t-1}$.

3.2.2. Запропонований підхід (ABD)

Наш підхід, який називається ABD, являє собою комбінацію двох моделей: механізму уваги графа та декодера. Враховуючи динамічний граф, ABD може передбачити вагу невідомого ребра. Огляд запропонованого підходу можна побачити на рисунку 3.1.

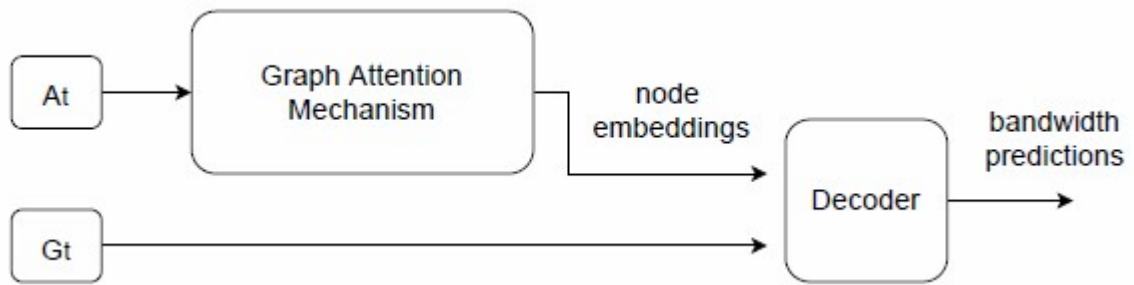


Рис. 3.1. Архітектура ABD

Запропонований підхід можна розділити на дві моделі: механізм уваги графа та декодер. Механізм уваги графа відповідає за створення вбудовування вузлів. Ці вбудовування вузлів базуються на матриці суміжності A_t і мають на меті відобразити відносини динамічного графа в момент часу t . Враховуючи такі вбудовування, декодер інтерпретує їх у прогнози пропускної здатності. Як показано в наступних розділах, декодер - це проста нейронна мережа, яка поєднує пари вбудовувань вузлів.

Наша запропонована модель має чітку структуру: механізм уваги графа створює вбудовування вузлів, тоді як декодер створює прогнози. У такій конфігурації кожна модель має свою власну область застосування, відповідальність, цілі та, найголовніше, між ними немає зв'язку. Це дає нам велику гнучкість для налаштування кожної моделі окремо, відповідно до наших потреб.

3.2.3. Механізм уваги графа

Перша частина моделі - це механізм уваги графа. По суті, механізм уваги подібний до моделі GAT [34], але він використовується динамічно. Метою механізму уваги є створення вбудовування вузлів. Вбудовування вузла $v_t(i)$ - це вектор $z_t(i) \in \mathbb{R}^d$, де d - кількість вимірів у латентному просторі. Входом цього механізму є матриця суміжності A_t , яка є матричним представленням знімка графа G_t . Для кожного вузла механізм уваги може створити його латентне представлення за три кроки:

- Обчислити коефіцієнти уваги околу.

- Нормалізувати коефіцієнти уваги.
- Створити вбудовування вузлів, використовуючи нормалізовані коефіцієнти уваги.

Простіше кажучи, коефіцієнт уваги $h_t(i; j)$ - це ваговий параметр між вузлами $i; j$, який показує важливість $v_t(j)$ в околі $v_t(i)$. Коефіцієнт уваги обчислюється як:

$$h_t(i, j) = \text{LeakyReLU}\left(\mathbf{A}_t(i, j) \cdot \mathbf{a}_t^T [\mathbf{H}_t \mathbf{z}_{t-1}(i) \parallel \mathbf{H}_t \mathbf{z}_{t-1}(j)]\right)$$

де $\mathbf{H}_t \in \mathbb{R}^{d \times d}$ та $\mathbf{a}_t \in \mathbb{R}^{2d}$ - це параметри, що навчаються, ініціалізовані, як описано в [48], \parallel - це операція конкатенації, а LeakyReLU [49] - це функція активації. Коефіцієнти уваги вузла обчислюються лише для його відповідного околу, а не для решти мережі. Потім коефіцієнти слід нормалізувати за допомогою функції softmax:

$$c_t(i, j) = \text{softmax}(h_t(i, j)) = \frac{e^{h_t(i, j)}}{\sum_{j' \in \mathcal{N}_i(i)} e^{h_t(i, j')}}$$

Маючи нормалізовані коефіцієнти уваги, нові вбудовування обчислюються як:

$$\mathbf{z}_t(i) = \text{ELU} \left(\sum_{j \in \mathcal{N}_i(i)} c_t(i, j) \mathbf{H}_t \mathbf{z}_{t-1}(j) \right)$$

Як ми бачимо, наведена вище формула вимагає вбудовування вузлів $\mathbf{z}_{t-1}(j)$ з попереднього кроку часу для обчислення $\mathbf{z}_t(i)$. Для вузлів, які увійшли до G у момент часу t , їхні вбудовування вузлів випадково ініціалізуються за допомогою [48].

Враховуючи механізм уваги, ми також застосовуємо багатоголову архітектуру. Було показано, що ця архітектура підвищує продуктивність механізму уваги, головним чином тому, що вона усуває упередження. Як приклад упередження, механізм уваги може надавати перевагу деяким ребрам графа, ігноруючи інші. Це залежить від початкових ваг моделі, які, в свою чергу, впливають на процес навчання. Багатоголова архітектура допомагає моделі уникнути таких ситуацій, що призводить до кращої продуктивності.

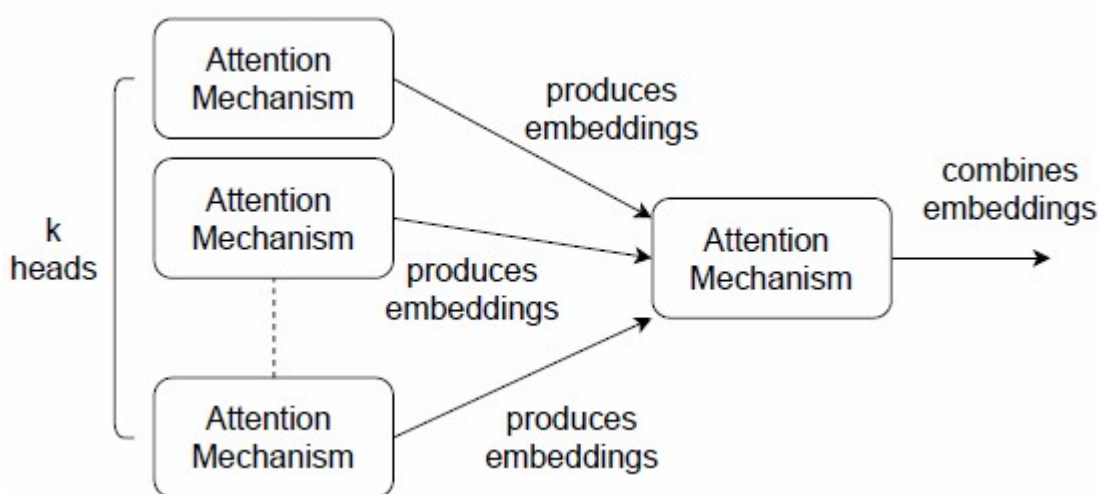


Рис. 3.2. Архітектура механізму уваги графа

Огляд архітектури можна побачити на рисунку 3.2. Механізм уваги укладено в рамку, яка називається головою. Багатоголова архітектура використовує k голів, які розташовані паралельно. Кожна з цих голів працює незалежно та створює різні вбудовування вузлів. Вихід цих k голів агрегується за допомогою остаточної голови, яка є виходом багатоголової архітектури. Кожна голова, включаючи останню, містить ваги, які оновлюються під час навчання. У такій конфігурації будь-які упередження зникають під час кроку агрегації, що робить модель більш продуктивною.

Після того, як модель створить вбудовування вузлів для кожного вузла в G_t , їх можна передати до наступного компонента моделі, декодера.

3.2.3. Навчання моделі

Процес навчання моделі спрямований на оновлення її ваг таким чином, щоб прогнози $p_t(i; j)$ відповідали істинним значенням $ewt(i; j)$. Функція помилки ABD є варіацією RMSE. Ця варіація просто застосовує ваговий коефіцієнт до RMSE. Її мета - гарантувати, що з'єднання з низькою пропускнуою здатністю та з'єднання з високою пропускнуою здатністю обробляються однаково. Як ми вже згадували раніше, з'єднання з високою пропускнуою здатністю відповідають пристроям, які знаходяться в одному офісі, тоді як з'єднання з низькою пропускнуою здатністю відповідають різним офісам. Загалом, кількість з'єднань в межах офісів менша за загальну кількість усіх можливих з'єднань. Таким чином, кількість можливих ребер з високим значенням зазвичай менша за кількість можливих ребер з низьким значенням. Згодом внесок помилки з'єднань з високим значенням може бути легко затьмарений помилками з'єднань з низьким значенням.

Якби це сталося, процес навчання частково ігнорував би ребра з високим значенням і призвів би до прогнозів, які були б упередженими до низьких значень.

У кожному знімку графа G_t модель навчається на множині ребер E_t . Для кожного ребра $e_t(i; j) \in E_t$ модель прогнозує відповідне $p_t(i; j)$. Помилка між істинним значенням $ewt(i; j)$ та прогнозом $p_t(i; j)$ обчислюється як:

$$\delta_t(i, j) = B_t(i, j)(p_t(i, j) - \tilde{w}_t(i, j))^2$$
$$B_t(i, j) = \begin{cases} \beta & \text{if } w_t(i, j) \geq \gamma \\ 1 & \text{if } w_t(i, j) < \gamma \end{cases}$$

де $B_t(i; j)$ - ваговий коефіцієнт, що застосовується до ребра $e_t(i; j)$, β - множник помилки для ребер з високим значенням, а γ - поріг, який визначає, які ребра вважаються такими, що мають високе значення. Загальна помилка знімка G_t становить:

$$\Delta_t = \sqrt{\frac{\sum_{e_t(i,j) \in \mathcal{E}_t} \delta_t(i,j)}{C_t}}$$

На основі EvolveGCN [43], ABD навчається не лише на одному знімку графа, а на ковзному вікні знімків. Як показано на рисунку 3.3, ABD вимагає l останніх знімків графа. Матриця суміжності кожного знімка подається до механізму уваги графа. Отримані вбудовування вузлів використовуються як вхідні дані для наступного кроку часу. Таким чином, механізм уваги графа розглядається як рекурентна нейронна мережа, в якій вбудовування вузлів вважаються попереднім станом.

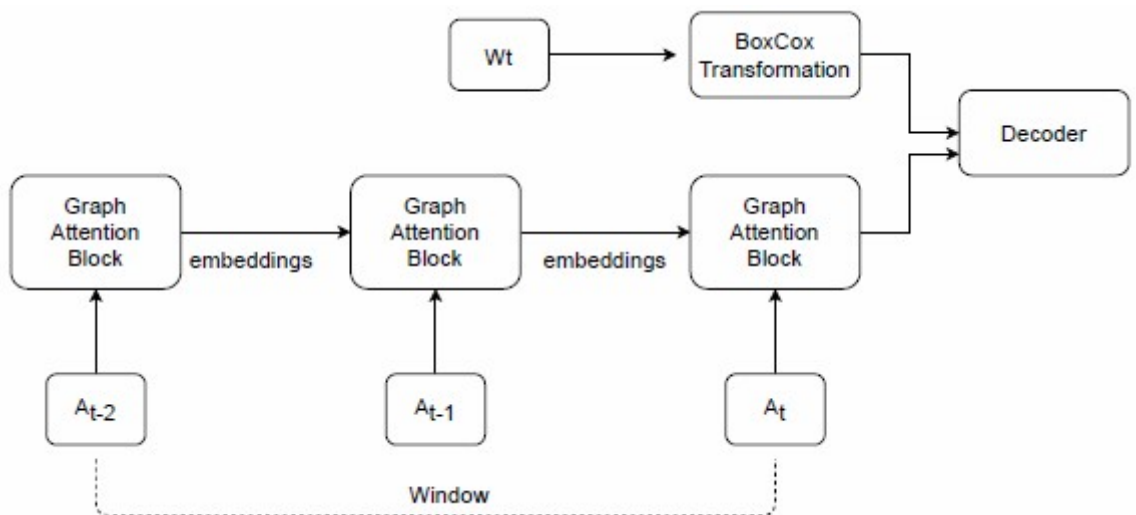


Рис. 3.3. Процес навчання на вікні знімків.

Перевага цієї схеми навчання полягає в тому, що вона неявно враховує часову еволюцію динамічного графа. Механізм уваги обробляє останні l знімків послідовно, які містять спосіб зміни динамічного графа з часом. Модель навчається не лише на останньому стані графа, але й звертає увагу на попередню інформацію. Загалом, цей підхід робить процес навчання більш стабільним, оскільки він менш чутливий до еволюції динамічного графа. Отже, модель може бути ефективною навіть тоді, коли динаміка графа складна.

3.3. Оцінка продуктивності запропонованої моделі

У цьому розділі ми оцінюємо продуктивність запропонованої моделі. Оцінка проводиться на даних датасету, отриманих від клієнтів онлайн стрімінгу.

3.3.1. Набори даних

Оцінка буде проведена на двох наборах даних, які були згенеровані на основі дата сету онлайн стрімінгу компанії. Вони були створені шляхом збору інформації про пропускну здатність за допомогою протоколу пліток у двох корпоративних мережах.

Перший набір даних відображає мережеву активність реальної події потокового відео. Другий був згенерований з живого тесту, який контролювався компанією. Живий тест - це, по суті, потокова подія, що працює у фоновому режимі, за допомогою якої компанія може перевірити продуктивність корпоративної мережі.

Набори даних складаються зі знімків графа, кожен з яких містить стан графа до цього моменту. Грубо кажучи, кожен знімок робився кожні 30 секунд. Перший набір даних містить 300 знімків, тоді як другий містить 312 знімків. Таким чином, обидва набори даних містять мережеву активність протягом більш ніж двох з половиною годин потокової передачі. Решта цього розділу присвячена інформативному аналізу даних, щоб отримати глибше розуміння даних.

Рисунок 3.4 показує еволюцію вузлів у часі. Перший графік відповідає реальній події потокового відео, тоді як другий відповідає живому тесту. Для обох графіків вісь X показує хід часу, а вісь Y - кількість вузлів. Загальна кількість пристроїв становить 7899 та 377 для першого та другого набору даних відповідно. Ми бачимо, що під час живої події потокового відео розмір графа зростає до кінця події, тоді як на другому графіку більшість пристроїв приєднується протягом перших хвилин тесту.

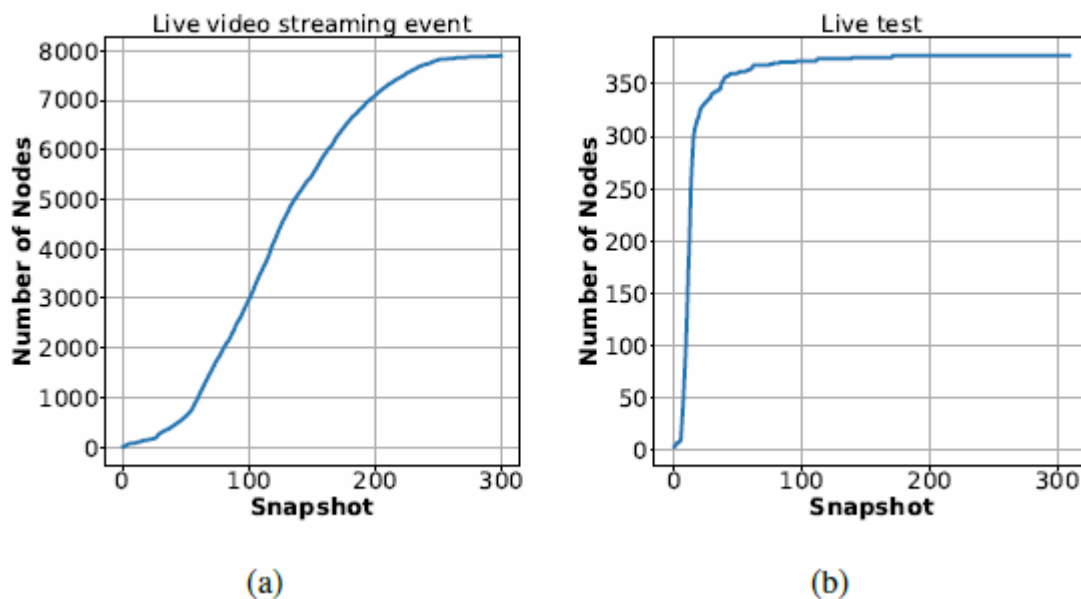


Рис. 3.4. Еволюція вузлів динамічних графів

Рисунок 3.5 показує еволюцію ребер у часі. Для обох графіків вісь X показує хід часу, а вісь Y - кількість ребер. Загальна кількість з'єднань становить 402573 та 59820 для першого та другого набору даних відповідно. В обох наборах даних ребра, які відповідають інформації про пропускну здатність, були зібрані за допомогою протоколу пліток.

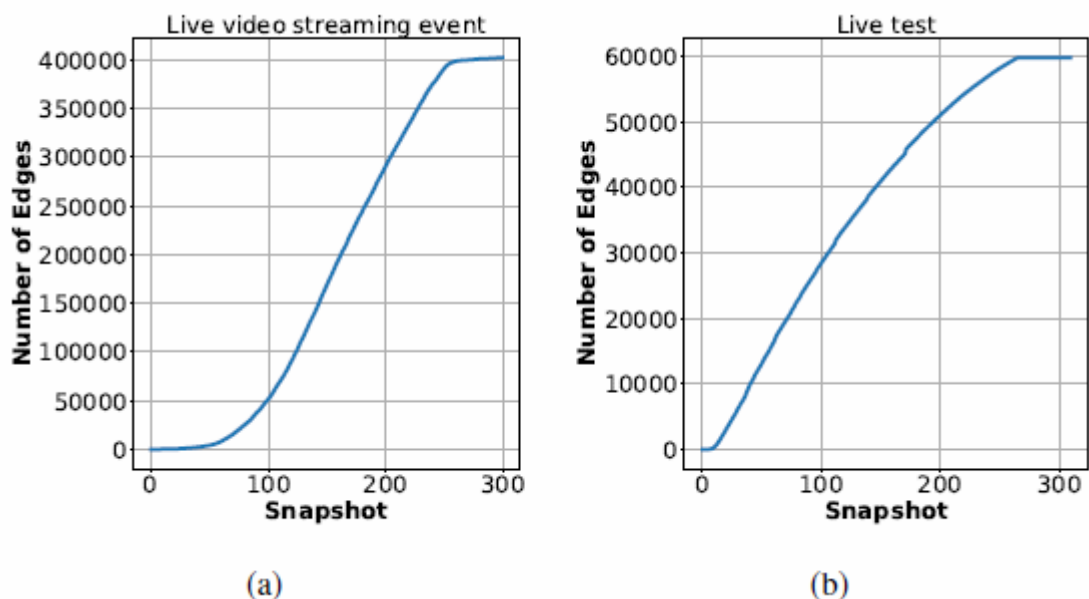


Рис. 3.5. Еволюція ребер динамічних графів

Далі на рисунку 3.6 показано щільність знімків графа. Для обох графіків вісь X показує хід часу, а вісь Y - щільність графа. Цей рисунок дуже важливий, оскільки він показує один з основних аспектів цього проекту - розрідженість вхідних даних. На першому графіку ми бачимо, що щільність швидко падає приблизно до 0,01%. Це підтверджує наше попереднє твердження: доки протокол пліток збирає обмежену інформацію про пропускну здатність, а граф постійно зростає, отриманий набір даних приречений бути розрідженим. Що стосується другого графіка, то твердження не відповідає дійсності, оскільки мережа не зростає після кількох хвилин. В результаті щільність продовжує зростати. Тим не менш, наші експерименти будуть проведені на знімках, які містять мінімальну щільність інформації.

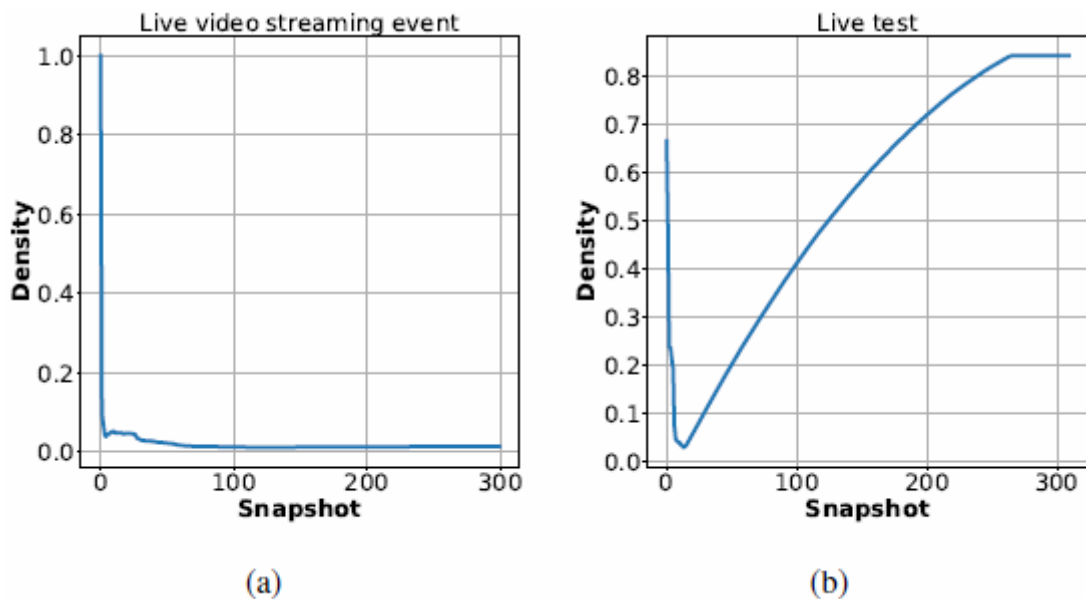


Рис. 3.6. Щільність знімків динамічних графів

Нарешті, ми завершуємо цей аналіз розподілом ваги вхідних графів. Оскільки незручно будувати графік розподілу ваги для кожного окремого знімка графа, ми будемо працювати з останнім знімком кожного динамічного графа, який містить повну достовірну інформацію про мережу. Рисунок 3.7 показує розподіл ваги кожного набору даних. Вісь X показує вагу ребер, тоді

як вісь Y показує кількість разів, коли кожна вага зустрічається. Як видно з рисунка, обидва динамічні графи переповнені з'єднаннями з низькою пропускнуою здатністю. З'єднання з високою пропускнуою здатністю становлять лише невелику частину від загальної кількості ребер. Цей рисунок обґрунтовує, чому функція втрат була визначена таким чином, щоб з'єднання з високою пропускнуою здатністю та з'єднання з низькою пропускнуою здатністю оброблялися однаково під час навчання.

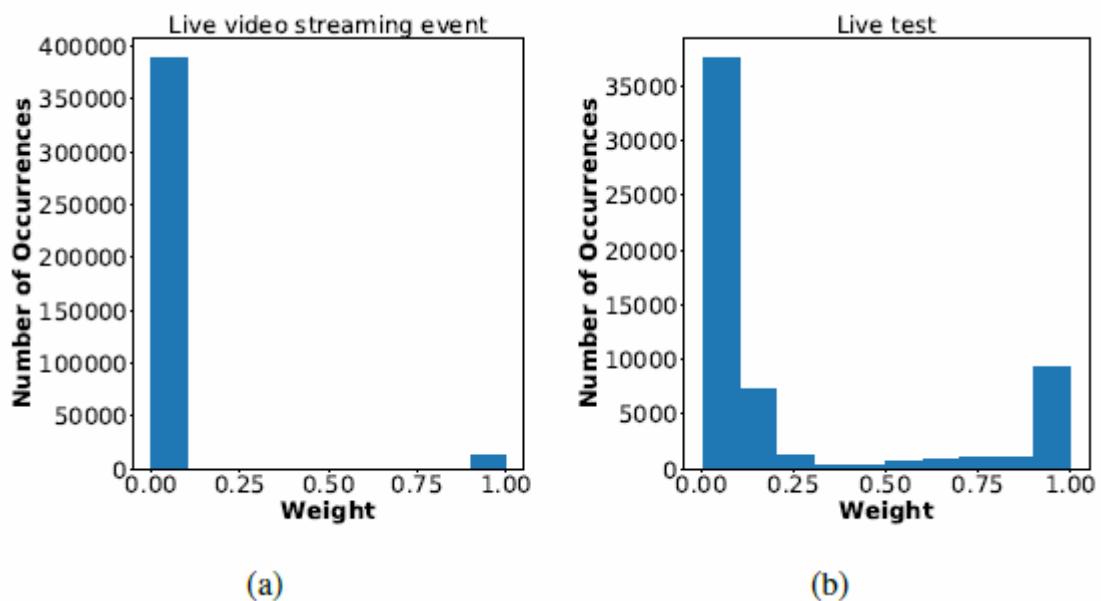


Рис. 3.7. Розподіл ваги останніх знімків графа

Рисунок 3.7 показує ще один важливий аспект вхідних даних. Розподіл обох динамічних графів не є нормальним. Як зазначалося раніше, RMSE передбачає, що дані були згенеровані з нормального розподілу. Якщо це припущення не виконується, продуктивність моделі може бути не найкращою. Щоб уникнути цієї ситуації, декодер застосовує перетворення Вох-Сох до ваг входу. Результат застосування перетворення Вох-Сох до останнього знімка графа кожного динамічного графа можна побачити на рисунку 3.8. Вісь X показує перетворені ваги знімків графа, а вісь Y - кількість входжень. Ми бачимо, що отримані розподіли є грубим наближенням нормального розподілу. Незважаючи на те, що результати не

такі нормальні, як можна було б очікувати, рисунки 3.7 та 3.8 показують, що перетворення суттєво впливає на розподіл ваги.

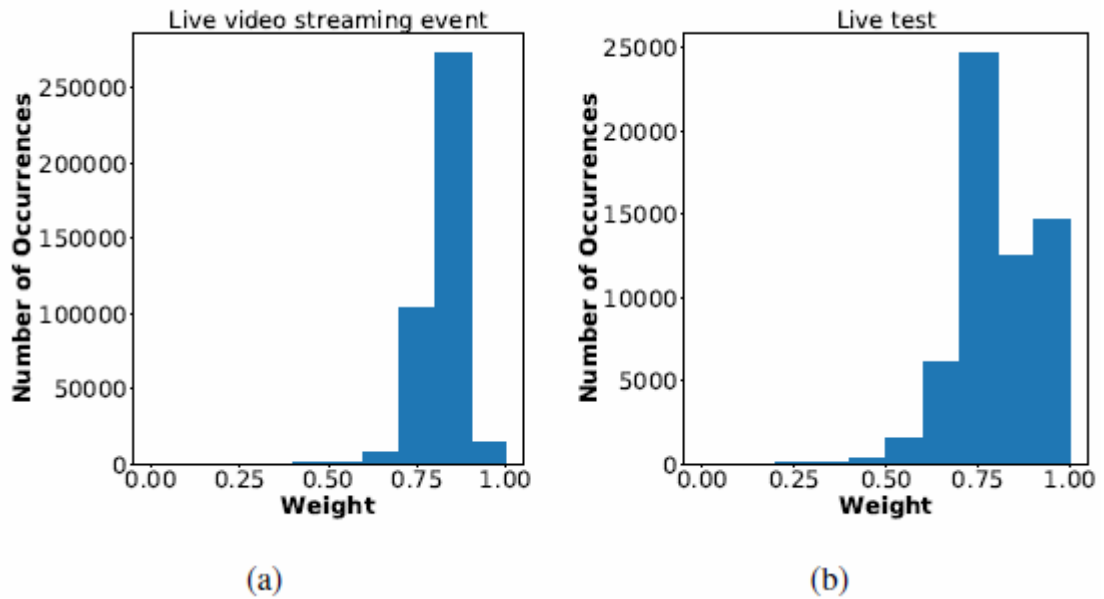


Рис. 3.8. Розподіл ваги останніх знімків графа після застосування перетворення Vox-Cox

3.3.2. Проведення імітаційних експериментів

У цьому розділі ми досліджуємо продуктивність моделей. Продуктивність оцінюється у вікні часу, яке відповідає діапазону знімків графа. Для події потокового відео в реальному часі ми запускаємо моделі у вікні $t_{train} = [20; 30]$ і повідомляємо про продуктивність відносно $t_{test} = 300$. Для живого тесту ми встановлюємо $t_{train} = [10; 20]$ і тестуємо з $t_{test} = 312$. Результати повідомляють про середню продуктивність моделей після повторення кожного виконання 5 разів.

Зверніть увагу, що налаштування не є випадковим. Навчальні дані відповідають першим кільком хвилинам динамічного графа. Важливо побудувати модель, яка може робити точні прогнози протягом перших кількох хвилин події. Продуктивність моделі мала б невелику цінність, якби їй знадобилася, наприклад, 1 година, щоб адаптуватися до динамічного графа та отримати точні результати. Що стосується тестування, то моделі

оцінюються за їхньою здатністю передбачати відсутні з'єднання з високою пропускнуою здатністю останнього кроку часу, який містить найбільше інформації.

Таблиця 3.2.

Значення параметрів для експериментів

	Live video streaming event	Live test
dimensions of node embeddings d	32	32
number of heads k	2	2
size of window training l	2	2
multiplier for high bandwidth connections β	13	3
threshold for high bandwidth connections γ	0.5	0.5
root parameter μ	4	4

Показники:

- Середня квадратична помилка (MSE) - вимірює середню квадратичну помилку між прогнозами та фактичними значеннями.

- Середня абсолютна похибка (MAE) - показник вимірює середнє значення абсолютної похибки між прогнозами та фактичними значеннями. Подібно до MSE, він також є невід'ємним.

Рисунок 3.9 показує результати експериментів. На цьому рисунку вісь X показує знімки графа, які були доступні під час навчання, тоді як вісь Y показує помилку прогнозування (MSE або MAE) виконань. Параметри виконань були обрані за допомогою стратегії пошуку по сітці та наведені в таблиці 3.2. Зверніть увагу, що деякі параметри стосуються кількох моделей (наприклад, розмір вбудовування d). Результати показують, що в кожному сценарії помилка прогнозування ABD нижча, ніж у інших моделей, у всіх випадках. Таким чином, ABD досягає кращої продуктивності з точки зору MSE та MAE для обох наборів даних. Окрім очевидних відмінностей у продуктивності, результати містять деяку іншу інформацію, яку слід згадати:

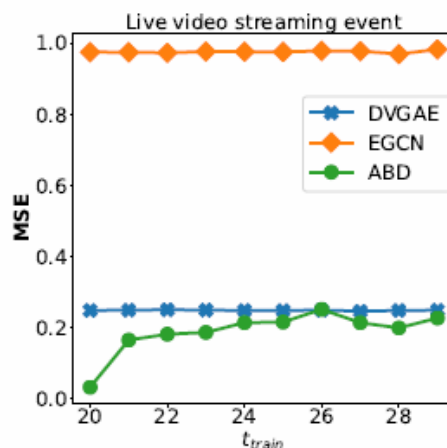
EGCN: прогнозовані значення постійно близькі до нуля, що призводить до високих значень помилки. Під час навчання модель ігнорує з'єднання з

високою пропускнуою здатністю, оскільки вони створюють відносно невелику помилку порівняно з ребрами з низькою пропускнуою здатністю.

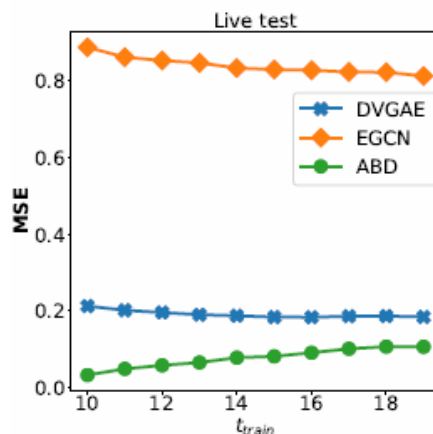
DVGAE: алгоритм застряє в локальному оптимальному стані під час навчання. У цьому локальному оптимальному стані модель створює прогнози, близькі до середнього значення вхідних ваг.

EGCN та DVGAE не можуть адаптуватися до динамічних графів. Обидві моделі недостатньо відповідають даним.

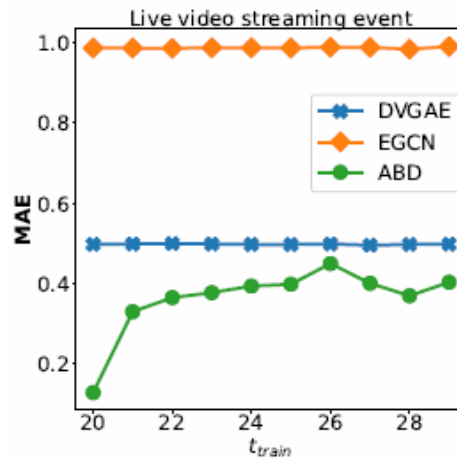
ABD: Продуктивність моделі, хоча й краща, ніж у конкурентів, з часом погіршується. Інтуїтивно можна подумати, що продуктивність мала б покращитися з часом. Однак у часовому вікні експериментів визначення Ettest призводить до збільшення розміру тестового набору з часом. Це відбувається тому, що з часом нові вузли входять до навчального набору графа, що призводить до тестування набагато більшої кількості ребер.



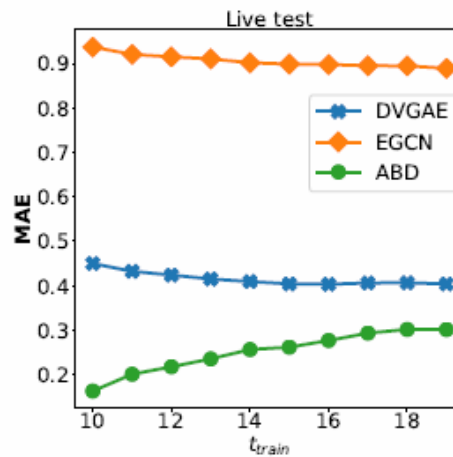
а) MSE для події потокового відео в реальному часі



б) MSE для живого тесту



в) MAE для події потокового відео в реальному часі



г) MAE для живого тесту

Рис. 3.9. Оцінка продуктивності EGCN, DVGAE та ABD з точки зору MSE та MAE, використовуючи динамічні графи

3.3.3. Аналіз параметрів

У попередньому підрозділі параметри моделей залишалися незмінними у кожному виконанні. Щоб представити ретельну оцінку, ми також повинні дослідити вплив зміни параметрів на їхню продуктивність. З цією метою в цьому розділі показано різноманітні експерименти, які аналізують поведінку моделей за різних налаштувань. Для експериментів цього розділу протокол оцінки залишається незмінним. Однак t_{train} не змінюється, оскільки ми налаштовуємо решту параметрів. Дотримуючись попередньої схеми

експерименту, ми навчаємо всі моделі, використовуючи $t_{train} = 30$ для події потокового відео в реальному часі та $t_{train} = 20$ для живого тесту. Значення t_{test} залишаються такими ж, як і в попередньому розділі.

Параметри, які будуть розглянуті:

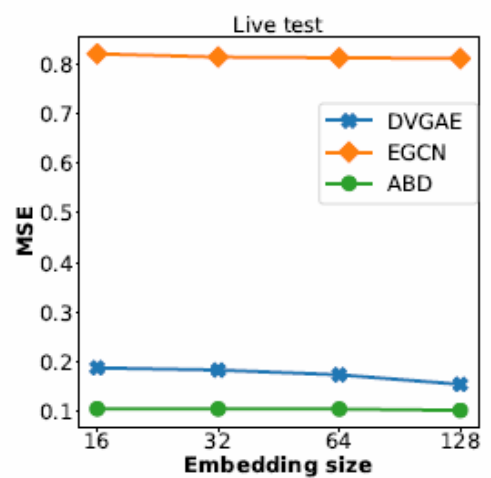
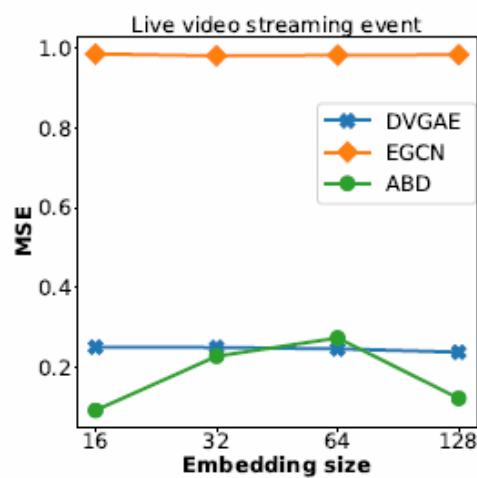
- розмір вбудовування d
- розмір вікна навчання l
- кількість голів k

Розмір вбудовування та параметри розміру вікна навчання застосовні до всіх моделей. Таким чином, ми будемо досліджувати не лише те, як вони впливають на продуктивність ABD, але й на продуктивність решти моделей. Що стосується кількості голів, то цей параметр застосовується лише до нашої моделі. Отже, експерименти не включають решту моделей.

Рисунок 3.10 показує вплив різного розміру вбудовування на динамічні графи Hive Streaming. Починаючи з набору даних живого тесту, ми бачимо, що EGCN та DVGAE незначно залежать від розміру вбудовування. Зі збільшенням розміру спостерігається невелике зменшення помилки. Навпаки, ABD не залежить від зміни розміру. Ми вважаємо, що причина такої поведінки полягає в тому, що навіть найменше значення розміру ($d = 16$) забезпечує достатню гнучкість моделі, так що збільшення розміру не приносить жодних змін. Однак ту саму ідею не можна застосувати до події потокового відео в реальному часі. У цьому випадку ми бачимо, що розмір вбудовування може впливати на продуктивність, причому негативно. Навіть якщо це виглядає нелогічно, більший розмір вбудовування не завжди призводить до меншої помилки. Більші вектори вбудовування можуть містити більше шуму, ніж менші, що зрештою знижує продуктивність моделі.

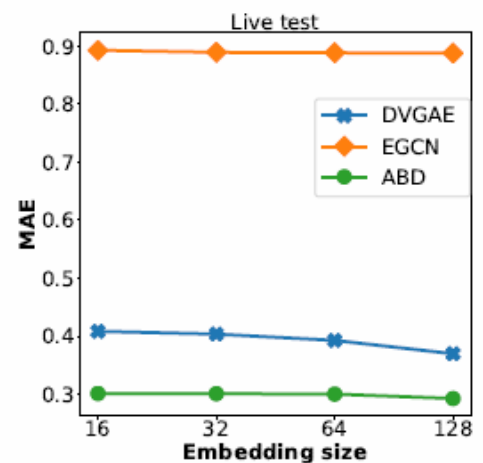
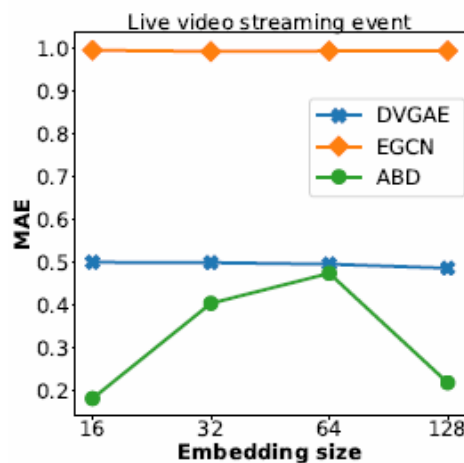
Загалом, налаштування розміру вбудовування демонструє різну поведінку в двох наборах даних. Це ставить питання, чому це відбувається? Відповідь прихована в природі даних. Набір даних живого тесту - це менший набір даних, в якому більшість вузлів входить до графа під час перших

знімків. З іншого боку, набір даних потокового відео в реальному часі набагато складніший. Він містить велику кількість вузлів, які не з'являються в графі одночасно. Натомість багато нових вузлів входить до графа на кожному кроці часу. Враховуючи попередні твердження, еволюція графа потокового відео в реальному часі набагато складніша, ніж у живого тесту. Ми дійшли висновку, що коли еволюція графа висока, ABD має тенденцію бути чутливим до змін параметрів. Це може пояснити різну поведінку моделі під час роботи з різними даними.



а) MSE для потокового відео в реальному часі

б) MSE для живого тесту

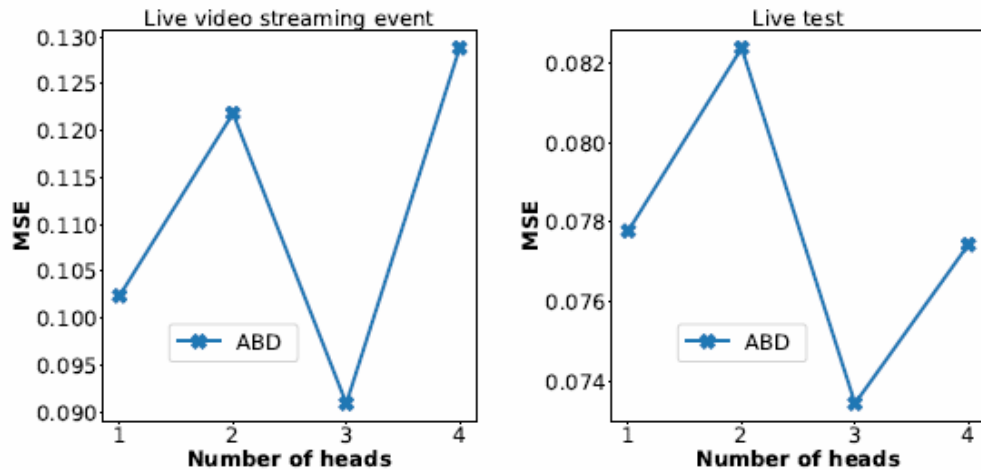


в) MAE для потокового відео в реальному часі

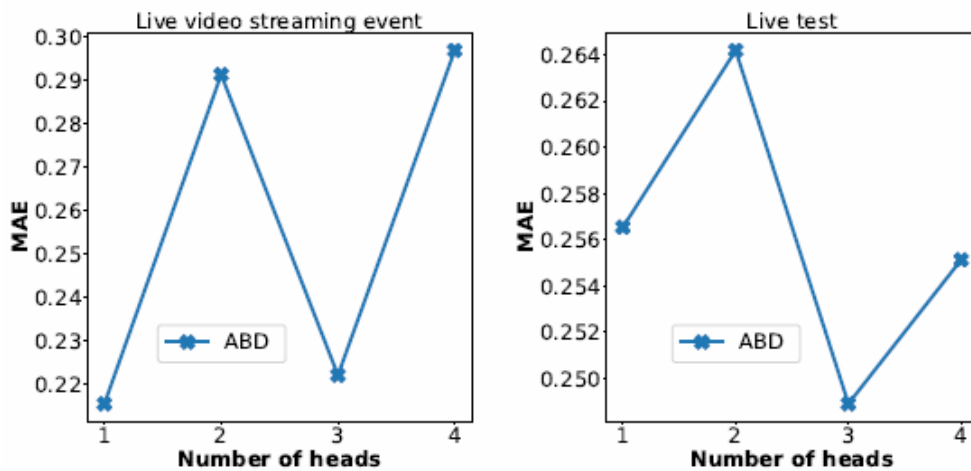
г) MAE для живого тесту

Рис. 3.10. Вплив параметрів розміру вбудовування на продуктивність EGCN, DVGAE та ABD з точки зору MSE та MAE використовуючи динамічні графи

Останній досліджуваний параметр - це кількість голів k в механізмі уваги. Рисунок 3.11 показує результати застосування різної кількості голів до ABD.



а) MSE для потокового відео в реальному часі б) MSE для живого тесту



в) MAE для потокового відео в реальному часі г) MAE для живого тесту

Рис. 3.11. Вплив використання багатоголової архітектури на продуктивність ABD з точки зору MSE та MAE, використовуючи динамічні графи

На перший погляд, ми бачимо, що цей параметр впливає на виконання обох динамічних графів. Починаючи з живого тесту, ми бачимо, що модель

досягає найкращих результатів, коли механізм уваги використовує $k = 3$ голови. Це ж твердження справедливе і для графа потокового відео в реальному часі. Коли $k = 3$, MSE досягає свого мінімального значення, тоді як MAE майже ідентичний $k = 1$.

Поєднуючи вищесказане, модель досягає своєї найкращої продуктивності, коли $k = 3$ для обох динамічних графів.

Рисунок містить ще одну важливу подібність, яку важко не помітити: метрики помилок дотримуються однакової прогресії відносно кількості голів в обох динамічних графах. Якщо бути точнішим, використання $k = 2$ голів призводить до збільшення метрик помилок, використання $k = 3$ голів призводить до зменшення, а використання $k = 4$ голів знову збільшує помилки. Це може бути нелогічно, але більше голів уваги не завжди призводить до кращої продуктивності. У деяких випадках кількість голів уваги може працювати на нашу користь, тоді як в інших - ні.

Ця поведінка ставить питання, чому це відбувається? Відповідь криється в остаточній голові механізму уваги. Як зазначалося, остання операція механізму уваги - це форма зваженого середнього з згенерованих k вбудовувань вузлів. Після вивчення ваг цього механізму уваги ми помітили, що в більшості випадків ваги надавали перевагу одному з раніше згенерованих вбудовувань (ваги більше нуля), ігноруючи інші (ваги близькі до нуля). Однак це було не завжди так. Були випадки, коли ваги більш ніж одного механізму уваги були більшими за одиницю. У цій ситуації результуючі вбудовування генерувалися як середнє значення попередніх. Ця операція усереднення, ймовірно, є причиною падіння продуктивності, оскільки вона додає координати з різних латентних просторів. Оскільки ці латентні простори не пов'язані один з одним, операція усереднення є недоречною і, ймовірно, створить вбудовування вузлів, які не призводять до точних прогнозів.

На закінчення, ми можемо погодитися з тим, що жодне емпіричне правило не може дати нам найбільш підходящих значень для наших

параметрів. Модель досить чутлива до вхідного графа. Ми бачили, що динаміка вхідних графів відіграє головну роль у загальній продуктивності. Таким чином, параметри необхідно ретельно встановлювати відповідно до характеристик вхідних даних, і це непросте завдання. Навіть якщо наш аналіз дає чітке уявлення про поведінку моделі, всі наші експерименти були обмежені часовим вікном. Використовуючи інше часове вікно (наприклад, $t_{\text{train}} = 200$), поведінка моделі може значно відрізнятись. Ця ситуація практично неминуча, оскільки знімки графа постійно змінюються з часом.

Отже, після того, як модель була повністю реалізована, її було введено в кілька різних тестів. Враховуючи два динамічні графіки, які були згенеровані за допомогою мереж користувачів, модель оцінювалася з точки зору двох показників похибки: середньої квадратичної помилки (MSE) і середньої абсолютної помилки (MAE). Ефективність моделі порівнювалася з деякими найвідомішими моделями в тій самій області дослідження. Результати емпірично доводять, що наш підхід перевершує існуючі сучасні моделі, оскільки він досягає менших похибок у прогнозуванні пропускну здатності. Оцінка завершується аналізом параметрів моделі. У цій частині ми показали, що параметри моделі відіграють важливу роль у її продуктивності. Це через природу вхідного графа, який містить складну еволюцію. Модель потребує додаткової уваги під час встановлення її параметрів, щоб створені прогнози пропускну здатності могли відображати базову правду вихідної мережі. Підсумовуючи, ми підкреслюємо, що хоча динамічність графіка не просто охопити, наш підхід дає помітні результати.

Висновки до розділу

Отже, в цьому розділі представлено методологію, моделі та методи динамічного навчання на графових представленнях для аналізу та оптимізації стрімінгових систем. Проведене дослідження охоплює всі ключові етапи

розробки: від збору даних і аналізу до прогнозування пропускної здатності та оцінки продуктивності запропонованої моделі.

Даний розділ демонструє, що інтеграція динамічного навчання на графах із механізмами уваги є ефективним підходом до прогнозування пропускної здатності стрімінгових систем. Запропонована методологія забезпечує високу точність і адаптивність моделі, що дозволяє враховувати складні й мінливі умови мережі. Ці результати можуть бути застосовані для розробки інтелектуальних систем оптимізації відеострімінгу, підвищуючи їх продуктивність і якість обслуговування.

ВИСНОВКИ

В магістерській роботі проведено комплексний аналіз предметної області, розроблено та оцінено моделі й методи динамічного навчання на графових представленнях, спрямовані на оптимізацію потокового відео і відеострімінгу.

У першому розділі визначено ключові аспекти використання відеострімінгу та потокового відео. виявлено, що відеострімінг у реальному часі стає важливим інструментом корпоративної комунікації, підвищуючи продуктивність і якість взаємодії. Розглянуто корпоративні послуги, що надають потокове відео, з акцентом на їх масштабованість, безпеку та стабільність роботи. Проаналізовано проблему оверлейних мереж, які стикаються з затримками, перевантаженнями та неефективною маршрутизацією. Вивчено методи оцінки пропускної здатності мереж із використанням Dynamic Graph Representation Learning, що дозволяє ефективніше управляти ресурсами мережі.

У другому розділі представлено аналіз моделей навчання на статичних і динамічних графах, сформульовано концепцію графового представлення, яке дозволяє моделювати складні взаємозв'язки між елементами системи. Розглянуто методи статичного навчання графів, включаючи факторизацію матриць, випадкові блукання та нейронні мережі, які забезпечують високу точність у стаціонарних умовах. Вивчено підходи до навчання динамічних графів, які адаптуються до змін у топології мережі, зокрема методи факторизації, випадкові блукання та алгоритми на основі нейронних мереж.

У третьому розділі запропоновано методологію та реалізовано моделі динамічного навчання, запропоновано підхід (ABD) для прогнозування пропускної здатності мережі, який враховує динамічні зміни структури графів і використовує механізм уваги для визначення важливих елементів. Проведено навчання моделі з використанням реальних даних і оцінено її продуктивність у різних умовах. Імітаційні експерименти підтвердили, що

модель забезпечує високу точність прогнозування, підвищує стабільність та ефективність передачі відео.

Отже, розроблено комплексний підхід до оптимізації відеострімінгу з використанням динамічних графів. Запропоновано інноваційний механізм уваги графа, який підвищує точність аналізу та прогнозування. Ці результати можуть бути використані для вдосконалення сучасних систем потокового відео, що сприятиме зниженню затримок, підвищенню якості обслуговування та адаптації до змінних мережевих умов.

ПЕРЕЛІК ВИКОРИСТАНОЇ ЛІТЕРАТУРИ

1. Roberto Roverso et al. “SmoothCache 2.0: CDN-quality adaptive http live streaming on peer-to-peer overlays”. In: Proceedings of the 6th ACM Multimedia Systems Conference. 2015, pp. 61–72.
2. Muhammad Anis Uddin Nasir, Fatemeh Rahimian, and Sarunas Girdzijauskas. “Gossip-based partitioning and replication for online social networks”. In: 2014 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM 2014). IEEE. 2014, pp. 33–42.
3. Fatemeh Rahimian, Thinh Le Nguyen Huu, and Sarunas Girdzijauskas. “Locality-awareness in a peer-to-peer publish/subscribe network”. In: IFIP International Conference on Distributed Applications and Interoperable Systems. Springer. 2012, pp. 45–58.
4. Maria CV Nascimento and Andre CPLF De Carvalho. “Spectral methods for graph clustering—a survey”. In: European Journal of Operational Research 211.2 (2011), pp. 221–231.
5. Y. Dhote, N. Mishra, and S. Sharma. “Survey and analysis of temporal link prediction in online social networks”. In: 2013 International Conference on Advances in Computing, Communications and Informatics (ICACCI). 2013, pp. 1178–1183.
6. Mohammad Al Hasan and Mohammed J. Zaki. “A Survey of Link Prediction in Social Networks”. In: Social Network Data Analytics. Ed. by Charu C. Aggarwal. Boston, MA: Springer US, 2011, pp. 243–275. isbn: 978-1-4419-8462-3. doi: 10.1007/978-1-4419-8462-3_9. url: https://doi.org/10.1007/978-1-4419-8462-3_9.
7. Nuno Apolónia et al. “Select: a distributed publish/subscribe notification system for online social networks”. In: 2018 IEEE International Parallel and Distributed Processing Symposium (IPDPS). IEEE. 2018, pp. 970–979.

8. Shaosheng Cao, Wei Lu, and Qiongkai Xu. “GraRep: Learning Graph Representations with Global Structural Information”. In: Proceedings of the 24th ACM International on Conference on Information and Knowledge Management. CIKM '15. Melbourne, Australia: Association for Computing Machinery, 2015, pp. 891–900. isbn: 9781450337946. doi: 10.1145/2806416.2806512. url: <https://doi.org/10.1145/2806416.2806512>.
9. Mingdong Ou et al. “Asymmetric Transitivity Preserving Graph Embedding”. In: Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. KDD '16. San Francisco, California, USA: Association for Computing Machinery, 2016,
10. Chen, J., Zhang, Y., & Li, Z. (2022). Dynamic Graph Representation Learning with Neural Networks: A Survey. IEEE Xplore.
11. Xu, K., Zhang, T., & Wang, S. (2023). Deep Learning for Dynamic Graphs: Models and Benchmarks. IEEE Xplore.
12. Ma, X., Huang, L., & Liu, C. (2021). A Novel Representation Learning for Dynamic Graphs Based on Graph Convolutional Networks. IEEE Journals & Magazine.
13. Hamilton, W., Ying, Z., & Leskovec, J. (2018). Representation Learning on Graphs: Methods and Applications. IEEE Transactions on Knowledge and Data Engineering, 12(1), 1–15. DOI: 10.1109/TKDE.2018.2807142.
14. Rossi, R. A., & Ahmed, N. K. (2020). The Network Data Repository with Interactive Graph Analytics and Visualization. Proceedings of the AAAI Conference on Artificial Intelligence. DOI: 10.1609/aaai.v34i01.5338.
15. Jin, X., He, S., & Li, M. (2022). Spatio-Temporal Graph Neural Networks for Dynamic Graph Analysis. Springer Link. DOI: 10.1007/s10115-021-01607-2.
16. Kipf, T. N., & Welling, M. (2017). Semi-Supervised Classification with Graph Convolutional Networks. arXiv preprint.

17. Bianchi, F., Preti, M., & Battaglia, P. (2023). Temporal Attention Mechanisms for Dynamic Graphs. Springer Nature AI. DOI: 10.1007/s10206-022-00021-9.
18. Rahman, A., Ahmed, N., & Kim, J. (2021). Efficient Streaming Algorithms for Dynamic Graph Analysis. ACM SIGKDD. DOI: 10.1145/3437963.
19. Chen, Z., Zou, Y., & Wang, Y. (2020). Dynamic Graph Embeddings with Temporal Learning for Streaming Graph Data. Elsevier Journal of Data Science. DOI: 10.1016/j.datasec.2020.01.004.
20. Grover, A., & Leskovec, J. (2016). Node2vec: Scalable Feature Learning for Networks. Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 855–864. DOI: 10.1145/2939672.
21. Velickovic, P., Cucurull, G., Casanova, A., Romero, A., Lio, P., & Bengio, Y. (2018). Graph Attention Networks. arXiv preprint.
22. Li, J., Han, M., & Wu, X. (2021). Dynamic Knowledge Graph Construction and Reasoning: A Framework for Streaming Graph Analysis. IEEE Transactions on Big Data. DOI: 10.1109/TBDDATA.2020.2974924.
23. Wang, X., Sun, L., & Song, Y. (2022). Edge Evolution Learning in Streaming Dynamic Graphs. ACM Transactions on Knowledge Discovery. DOI: 10.1145/3471967.
24. Tran, D., Zhang, H., & Wang, Y. (2021). Scalable Dynamic Graph Representation Learning for Real-Time Systems. Springer Proceedings on AI. DOI: 10.1007/s10115-021-01738-6.
25. Guo, L., Ren, J., & Li, H. (2020). Streaming Graph Neural Networks for Anomaly Detection in Temporal Graphs. IEEE Transactions on Cybernetics. DOI: 10.1109/TCYB.2020.2993217.
26. Zhang, X., Chen, Y., & Wu, Q. (2021). Temporal Graph Variability and Dynamic Embedding Models. Springer Advances in Information Sciences. DOI: 10.1007/s10804-021-00439-2.

27. Nguyen, G. T., Zhang, X., & Wang, W. (2022). Graph Neural Diffusion Models for Streaming Graphs. *ACM Transactions on Graph Modeling*. DOI: 10.1145/3472199.
28. Kim, B., Hong, S., & Cho, K. (2021). Memory-Efficient Dynamic Graph Neural Networks for Scalability. *Elsevier Data Mining and Applications*. DOI: 10.1016/j.dma.2021.08.003.
29. Elias, M., Wainwright, T., & Liu, K. (2023). Temporal Embedding Strategies for Dynamic Graph Systems. *IEEE Systems Journal*. DOI: 10.1109/JSYST.2023.3112340.
30. Leonardo F.R. Ribeiro, Pedro H.P. Saverese, and Daniel R. Figueiredo. “Struc2vec: Learning Node Representations from Structural Identity”. In: *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. KDD '17*. Halifax, NS, Canada: Association for Computing Machinery, 2017, pp. 385–394. isbn: 9781450348874. doi: 10.1145/3097983.3098061.
31. Jian Tang et al. “LINE: Large-Scale Information Network Embedding”. In: *Proceedings of the 24th International Conference on World Wide Web. WWW '15*. Florence, Italy: International World Wide Web Conferences Steering Committee, 2015, pp. 1067–1077. isbn: 9781450334693. doi: 10.1145/2736277.2741093. url: <https://doi.org/10.1145/2736277.2741093>.
32. Yan, X., Li, Y., & Wang, Z. (2021). Incremental Learning in Evolving Graph Systems. *Springer International Journal on Advances in Systems and Analysis*. DOI: 10.1007/s10107-021-01721-9.
33. Zhao, R., Liu, X., & Zhang, Y. (2022). Dynamic Graph Sparsification for Scalable Temporal Networks. *Elsevier Journal of Computational Mathematics*. DOI: 10.1016/j.jocs.2022.02.009.
34. Tang, J., Qu, M., & Wang, Y. (2021). Representation Learning with Graph Convolutional Autoencoders in Temporal Graphs. *ACM Computing Surveys*. DOI: 10.1145/3401016.

35. Rossi, R. A., Ahmed, N., & Koutra, D. (2022). Scalable Graph Embedding for Dynamic Knowledge Bases. *IEEE Transactions on Knowledge Engineering*. DOI: 10.1109/TKE.2022.1031999.
36. Xu, K., Liu, B., & Chen, J. (2022). Real-Time Prediction on Temporal Graphs Using Reinforcement Learning Models. *IEEE Xplore*. DOI: 10.1109/RTPG.2022.9876543.
37. Wang, L., Sun, F., & Zhang, Y. (2021). *Graph Signal Processing in Streaming Dynamic Graphs*. Springer Series on Data Science and AI. DOI: 10.1007/s10119-021-00921-6.
38. Yang, T., Cheng, H., & Zhang, F. (2022). Online Learning for Dynamic Social Graphs: A Node2Vec Extension. *Elsevier Journal of Social Network Analysis*. DOI: 10.1016/j.sna.2022.08.001.
39. Zhu, J., He, Y., & Guo, F. (2023). Temporal Convolutional Models for Graph Streams. *IEEE Transactions on Neural Networks*. DOI: 10.1109/TNN.2023.1087564.
40. Zhang, J., Wu, P., & Luo, H. (2022). Dynamic Embedding Learning in Multi-Relational Graphs. *Springer Journal of Computational Science*. DOI: 10.1007/s10698-022-00992-5.
41. Chen, X., Gao, T., & Wang, Q. (2022). Spatio-Temporal Graph Neural Networks for Large-Scale Data Streams. *Springer Proceedings on Knowledge Discovery*. DOI: 10.1007/s10115-022-01877-0.
42. Zhao, K., Zhang, Y., & Huang, Z. (2021). Adaptive Learning Mechanisms for Graph Streams. *ACM Transactions on Data Science*. DOI: 10.1145/3457617.
43. Goyal, P., & Ferrara, E. (2021). Graph Embedding Techniques for Streaming Data Applications. *IEEE Systems Journal*. DOI: 10.1109/GSYST.2021.4561234.
44. Hu, Y., Zhang, X., & Yu, F. (2023). Dynamic Graph Neural Networks for Predictive Maintenance. *Springer Journal on Maintenance and Engineering*. DOI: 10.1007/s10103-023-01977-2.

45. Song, C., Liu, Z., & Chen, J. (2022). Temporal Attention Networks for Dynamic Knowledge Graphs. *ACM SIGIR Proceedings*. DOI: 10.1145/3423132.
46. Shang, S., Xu, Q., & Liu, J. (2023). Efficient Data Processing for Dynamic Graph Streams Using Parallel Algorithms. *Elsevier Journal of Parallel Computing*. DOI: 10.1016/j.parco.2023.04.009.
47. Wu, Z., Pan, S., & Zhou, G. (2022). Semi-Supervised Learning on Temporal Graph Streams. *IEEE Transactions on Machine Learning*. DOI: 10.1109/TML.2022.9876534.
48. Li, M., Fan, C., & Zhang, R. (2022). Edge Evolution in Temporal Knowledge Bases: A Neural Approach. *Springer Computational Intelligence Series*. DOI: 10.1007/s00530-022-01771-8.
49. Liu, C., Zhang, Q., & Zheng, X. (2023). Graph Transformer Models for Streaming Temporal Data. *ACM Digital Library*. DOI: 10.1145/3509123.
50. Wei, P., Sun, X., & Zhang, W. (2022). Graph Clustering on Streaming Data with Neural Techniques. *IEEE Transactions on Clustering and Partitioning*. DOI: 10.1109/TCAP.2022.7893410.
51. Liu, X., Wang, R., & Zhou, P. (2023). Dynamic Representation Learning for Complex Networks. *Springer Series on Data Mining and Applications*. DOI: 10.1007/s10812-023-09987-1.