

Міністерство освіти і науки України
Харківський національний університет радіоелектроніки

Факультет Комп'ютерних наук
(повна назва)

Кафедра Штучного інтелекту
(повна назва)

КВАЛІФІКАЦІЙНА РОБОТА
Пояснювальна записка

рівень вищої освіти перший (бакалаврський)

Дослідження методів токенізації в задачах продуктивності NLP-моделей
(тема)

Виконав:
здобувач четвертого року навчання,
групи ІТШ-21-2

Орина Майданович
(власне ім'я, прізвище)

Спеціальність 122 Комп'ютерні науки
(код і повна назва спеціальності)

Тип програми освітньо-професійна
Освітня програма Штучний інтелект
(повна назва освітньої програми)

Керівник доц. Марина Кудрявцева
(посада, власне ім'я, прізвище)

Допускається до захисту

Завідувач кафедри ШІ _____
(підпис)

Олег ЗОЛОТУХІН
(власне ім'я, прізвище)

2025 р.

Харківський національний університет радіоелектроніки

Факультет _____ Комп'ютерних наук _____

Кафедра _____ Штучного інтелекту _____

Рівень вищої освіти _____ перший (бакалаврський) _____

Спеціальність _____ 122 Комп'ютерні науки _____
(код і повна назва)

Тип програми _____ освітньо-професійна _____

Освітня програма _____ Штучний інтелект _____
(повна назва)

ЗАТВЕРДЖУЮ:

Зав. кафедри _____

(підпис)

« _____ » _____ 20 ____ р.

ЗАВДАННЯ
НА КВАЛІФІКАЦІЙНУ РОБОТУ

здобувачеві _____ Майданович Орині В'ячеславівні _____
(прізвище, ім'я, по батькові)

1. Тема роботи Дослідження методів токенизації в задачах продуктивності NLP-моделей

затверджена наказом університету від 19 травня 20 25 р. № 378Ст

2. Термін подання студентом роботи до екзаменаційної комісії 25 червня 20 25 р.

3. Вихідні дані до роботи Корпус текстових даних англійською мовою з відкритого датасету IMDb, який містить рецензії на фільми та відповідні бінарні мітки (позитивна або негативна); стандартизовані токенизатори з бібліотеки HuggingFace Tokenizers: BPE, WordPiece та Unigram, реалізовані з нуля на основі текстового корпусу.

4. Перелік питань, що потрібно опрацювати в роботі _____

1) Аналіз предметної галузі _____


2) Аналіз сучасних підходів до токенизації в контексті NLP-моделей _____

3) реалізація програмної частини _____

КАЛЕНДАРНИЙ ПЛАН

№	Назва етапів роботи	Строк / терміни виконання етапів роботи	Примітка
1	Отримання завдання на кваліфікаційну роботу	19.05.2025	виконано
2	Аналіз предметної галузі	20.05.2025	виконано
3	Постановка задачі	22.05.2025	виконано
4	Програмна реалізація алгоритмів	24.05.2025	виконано
5	Тестування програми	26.05.2025	виконано
6	Експериментальні дослідження	27.05.2025	виконано
7	Підготовка пояснювальної записки	13.06.2025	виконано
8	Підготовка презентації та доповіді	15.06.2025	виконано
9	Захист роботи	25.06.2025	виконано

Дата видачі завдання 19 травня 2025 р.

Здобувач _____

 (підпис)

Керівник роботи _____
 (підпис) доц. Марина Кудрявцева
 (посада, власне ім'я, прізвище)

РЕФЕРАТ

Пояснювальна записка: 73 с., 18 рис., 1 табл., 2 дод., 20 джерела.

МОВНА МОДЕЛЬ, ОБРОБКА ПРИРОДНОЇ МОВИ,
ПЕРЕДОБРОБКА ТЕКСТУ, ПРОДУКТИВНІСТЬ, ТОКЕНІЗАЦІЯ,
LANGUAGE MODEL, NATURAL LANGUAGE PROCESSING,
PERFORMANCE, PREPROCESSING, TOKENIZATION.

Об'єктом дослідження є процес обробки природної мови в інтелектуальних системах.

Предметом дослідження стали методи токенізації текстових даних та їхній вплив на продуктивність NLP-моделей.

Метою роботи виступило порівняльне дослідження різних підходів до токенізації та визначення їх впливу на ефективність моделей обробки природної мови.

Методом дослідження є експериментальне моделювання, що передбачає реалізацію, тестування та порівняльний аналіз алгоритмів токенізації в контексті типових NLP-завдань.

ABSTRACT

Bachelor's thesis contains: 73 pp., 18 fig., 1 tabl., 2 ann., 20 references.

LANGUAGE MODEL, NATURAL LANGUAGE PROCESSING,
PERFORMANCE, PRODUCTIVITY, TEXT PREPROCESSING,
TOKENIZATION.

The object of the study is the process of natural language processing in intelligent systems.

The subject of the study was the methods of tokenization of text data and their impact on the performance of NLP models.

The aim of the work was a comparative study of different approaches to tokenization and determination of their impact on the efficiency of natural language processing models.

The research method is experimental modeling, which involves implementation, testing and comparative analysis of tokenization algorithms in the context of typical NLP tasks.

ЗМІСТ

Перелік умовних позначень, символів, одиниць, скорочень і термінів	7
Вступ.....	8
1 Аналіз предметної галузі	9
1.1 Роль токенизації в системах обробки природної мови.....	9
1.2 Основні методи токенизації	11
1.2.1 Токенизація слів	11
1.2.2 Токенизація символів	13
1.2.3 Токенизація підслів	14
1.2.4 Токенизація речень.....	16
1.3 Виклики токенизації в NLP	17
2 Аналіз сучасних підходів до токенизації в контексті nlp-моделей	22
2.1 Архітектура сучасних NLP-систем.....	22
2.2 Порівняння найпоширеніших підслівних алгоритмів.....	24
2.3 Обмеження готових токенизаторів.....	29
2.4 Аналіз публікацій, що ізольовано досліджують токенизацію	32
3 Реалізація програмної частини.....	35
3.1 Вибір засобів програмування.....	35
3.2 Аналіз та вибір алгоритмів розроблення.....	37
3.3 Вибір та завантаження бібліотек	40
3.4 Архітектура класифікаційної моделі	43
3.5 Підготовка текстових даних для навчання токенизаторів	46
3.6 Реалізація процесу навчання токенизаторів.....	48
3.7 Навчання моделі	53
3.8 Аналіз та інтерпретація результатів	58
Висновки.....	63
Перелік джерел посилання	65
Додаток А Вихідний код програмної реалізації	68
Додаток Б Відомість кваліфікаційної роботи	73

ПЕРЕЛІК УМОВНИХ ПОЗНАЧЕНЬ, СИМВОЛІВ, ОДИНИЦЬ, СКОРОЧЕНЬ І ТЕРМІНІВ

BPE – Byte-Pair Encoding – кодування байтової пари;

HTML – Hypertext Markup Language – мова гіпертекстової розмітки;

LLM – Large Language Model – велика мовна модель;

NLP – Natural Language Processing – обробка природної мови;

OOV – Out Of Vocabulary – поза словниковим запасом;

UNK – Unknown Token – невідомий токен.

ВСТУП

В останні роки алгоритми все частіше виконують функцію співрозмовників, перекладачів і навіть авторів тексту, через що їх здатність розуміти людську мову як структуру мислення, а не як послідовність символів, стає дедалі важливішою. Тому перетворення неформального мовлення у прийнятну для обчислення форму є одним з основних питань обробки природної мови (Natural Language Processing). Незважаючи на значний прогрес у побудові мовних моделей, їхня ефективність та якість продовжують залежати від попереднього етапу обробки – токенизації [1].

Токенизація – це процес розбиття потоку текстового контенту на слова, терміни, символи або інші значущі елементи, тобто токени. Вибір токенизатора зазвичай залежить як від обсягу навчальних даних, так і від якості прогнозування. Також, у міру розвитку моделей, змінюються методи токенизації: від простих правил до статистичних методів і навчання без учителя.

В цій кваліфікаційній роботі ми ставимо перед собою завдання систематизувати та порівняти існуючі методи токенизації в контексті їхнього впливу на продуктивність моделі NLP. Особливий інтерес викликають найновіші алгоритми підслів, які стали фактичним стандартом в трансформерних архітектурах. Однак, попри всю їхню популярність, залишається відкритим питання оптимального вибору токенизатора, особливо для практичних застосувань з обмеженнями часу, ресурсів або мовою тексту.

Метою дослідження є порівняння різних методів токенизації щодо їхньої точності, швидкості та вплив на результативність моделей при вирішенні типових NLP-завдань. У роботі передбачено аналіз теоретичних основ токенизації, реалізацію обраних методів та проведення експериментів з їх порівняння. Результати можуть бути використані при побудові текстових моделей та оптимізації передобробки даних у мовних системах.

1 АНАЛІЗ ПРЕДМЕТНОЇ ГАЛУЗІ

1.1 Роль токенизації в системах обробки природної мови

Обробка природної мови (NLP) – це підгалузь інформатики, штучного інтелекту, інформаційної інженерії та взаємодії людини з комп'ютером. Вона зосереджена на тому, як програмувати комп'ютери для обробки та аналізу великих обсягів даних природної мови, що не просто реалізувати, адже процес читання та розуміння мов набагато складніший, ніж здається на перший погляд.

Токенизація – це процес поділу тексту на менші одиниці, відомі як токени. Токени – це зазвичай слова або підслова в контексті обробки природної мови. Токенизація – це критичний крок у багатьох завданнях NLP, включаючи обробку тексту, моделювання мови та машинний переклад [2]. Процес включає розділення рядка або тексту на список токенів. Токени можна розглядати як частини, наприклад, слово – це токен у реченні, а речення – це токен в абзаці.

Токенизація передбачає використання токенизатора для сегментації неструктурованих даних і тексту природної мови на окремі фрагменти інформації, розглядаючи їх як різні елементи. Процес токенизації виконує роль фільтрації, стандартизації та структуризації мовного матеріалу, з якого пізніше формуються векторні подання – числові структури, що подаються на вхід алгоритмам машинного навчання. Таким чином, якість токенизації безпосередньо впливає на адекватність представлення тексту для моделі [3].

Розглянемо типовий ланцюжок NLP та яке місце в ньому посідає токенизація. Першим етапом обробки природної мови виступає очищення та нормалізація тексту, де текст очищається від зайвих символів, наприклад, HTML-тегів, приводиться до нижнього регістру, стандартизується пунктуація, числа тощо. На другому етапі виконується токенизація тексту згідно з обраною схемою, чи то символи, слова, субслова тощо. Це перший

крок, який переводить текст у послідовність одиниць обробки. На третьому етапі здійснюється лематизація та стемінг, які полягають у зменшенні слів до базової форми для зниження варіативності. Слід за ним виконується перетворення токенів у числові вектори: частотні (TF-IDF, Bag-of-Words) або семантичні (Word2Vec, GloVe, BERT тощо). Нарешті, завершальним кроком виступає подача векторизованих токенів у машинну модель (класифікатор, трансформер, seq2seq тощо). Розширений ланцюжок NLP викладено на рисунку 1.1.

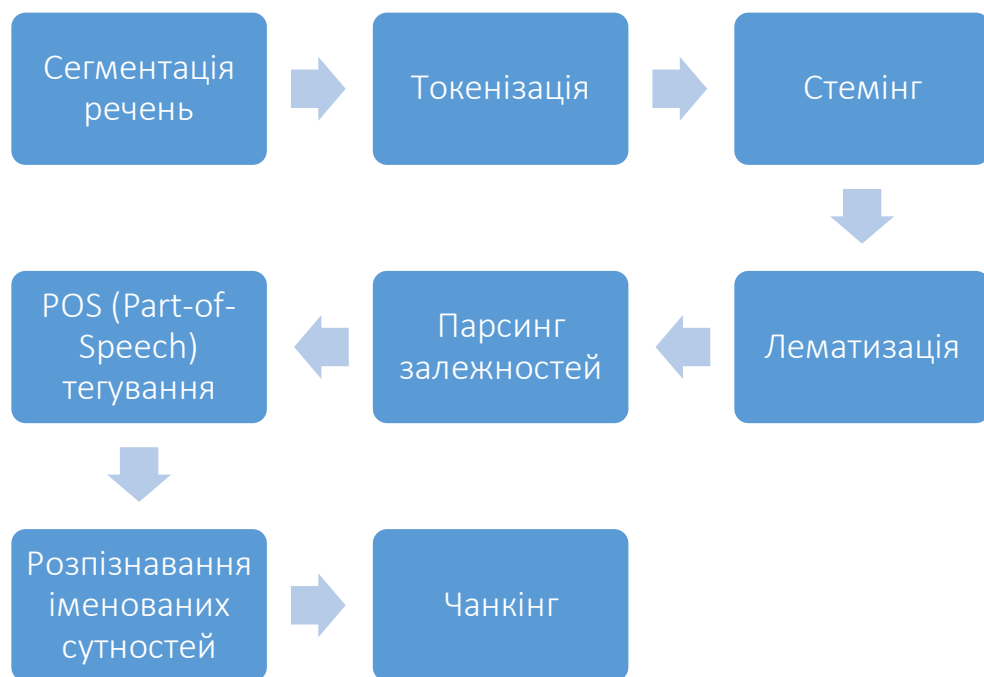


Рисунок 1.1 – Розширений ланцюжок NLP

Токенізація застосовується майже в усіх задачах NLP, зокрема:

- класифікація тексту;
- аналіз тональності (Sentiment Analysis);
- машинний переклад;
- автоматична генерація тексту;
- іменовані сутності.

У класифікації тексту, наприклад, при інтерпретації рецензій як позитивні або негативні, токенизація визначає, з якими елементами тексту працює модель: окремими словами, морфемами чи символами. Від її точності залежить коректність векторного представлення та якість навчання. Для розпізнавання емоцій у тексті важлива передача контексту та семантична цілісність токенів: слово «happy» не повинно бути розбите невдало, інакше модель не впізнає його позитивного значення. У сучасних системах перекладу, наприклад, Google Translate або DeepL, застосовуються моделі на кшталт Transformer, які не працюють без попередньої токенизації. Наприклад, англійське слово «internationalization» може бути поділено на «inter», «national» та «ization», що полегшує переклад на інші мови та зменшує розмір словника. Моделі на кшталт GPT-3 та GPT-4 створюють текст токен за токеном. Якість сгенерованого тексту сильно залежить від схеми токенизації, що визначає, які одиниці модель вважатиме значущими. Невдала токенизація може спричинити неприродні або ламані фрази. Для розпізнавання назв організацій, людей, місць тощо, важливо не «ламати» такі вирази, як «New York Times» або «Barack Obama». Токенізатор має зберігати межі таких сутностей.

1.2 Основні методи токенизації

1.2.1 Токенизація слів

Токенизація слів – найпоширеніший варіант токенизації (лістинг 1.1). У ньому використовуються паузи в мовленні або пробіли в тексті в якості сепараторів, а дані розділяються на відповідні слова за допомогою роздільників «,» або «;» [4].

Лістинг 1.1 – Програмний код функції токенизації слів

```
text = "Natural language processing (NLP) is a field of
```

Продовження лістингу 1.1

```
computer science, artificial intelligence and
computational linguistics concerned with the interactions
between computers and human (natural) languages, and, in
particular, concerned with programming computers to
fruitfully process large natural language corpora.
Challenges in natural language processing frequently
involve natural language understanding, natural language
generation (frequently from formal, machine-readable
logical forms), connecting language and machine
perception, managing human-computer dialog systems, or
some combination thereof."
print(word_tokenize(text))
```

Результат виконання функції токенизації слів можна побачити на рисунку 1.2.

```
['Natural', 'language', 'processing', '(', 'NLP', ')', 'is', 'a', 'field', 'of', 'computer', 'science', ',',
'artificial', 'intelligence', 'and', 'computational', 'linguistics', 'concerned', 'with', 'the', 'interacti
ons', 'between', 'computers', 'and', 'human', '(', 'natural', ')', 'languages', ',', 'and', ',', 'in', 'part
icular', ',', 'concerned', 'with', 'programming', 'computers', 'to', 'fruitfully', 'process', 'large', 'natu
ral', 'language', 'corpora', '.', 'Challenges', 'in', 'natural', 'language', 'processing', 'frequently', 'in
volve', 'natural', 'language', 'understanding', ',', 'natural', 'language', 'generation', '(', 'frequently',
'from', 'formal', ',', 'machine-readable', 'logical', 'forms', ')', ',', 'connecting', 'language', 'and', '
machine', 'perception', ',', 'managing', 'human-computer', 'dialog', 'systems', ',', 'or', 'some', 'combinat
ion', 'thereof', '.']
```

Рисунок 1.2 – Результат виконання функції токенизації слів

Але такий найпростіший спосіб розділити мовлення або текст на частини має деякі недоліки. Токенизація слів складна для розділення невідомих слів або слів поза словниковим запасом (OOV) [5]. Це часто вирішується заміною невідомих слів простим токеном, який повідомляє, що слово невідоме. Однак це можна назвати доволі грубим рішенням, оскільки 5 токенів «невідомого» слова можуть бути 5 абсолютно різними невідомими словами або всі вони можуть бути одним і тим самим словом.

Точність токенизації слів залежить від словникового запасу, на якому вона навчається [6]. В таких моделях необхідно знаходити баланс між завантаженням слів для максимальної точності та максимальної ефективності. Хоча додавання словникового запасу всього словника зробило б NLP модель точнішою, часто це не найефективніший метод.

Наявність великого словника може бути проблемою, оскільки вона вимагає від нейронних мереж величезної кількості параметрів. Щоб проілюструвати це, припустимо, що у нас є 1 мільйон унікальних слів, і ми хочемо стиснути 1-мільйонні вхідні вектори до 1-тисячмірних векторів у першому шарі нашої нейронної мережі. Це стандартний крок у більшості архітектур NLP, і результуюча матриця ваг цього першого шару міститиме 1 мільярд ваг. Це вже можна порівняти з найбільшою моделлю GPT-2.4 яка має загалом близько 1,5 мільярда параметрів.

Звичайно, ми хочемо уникнути такої марнотратності з параметрами моделі, оскільки вони дорогі для навчання, а більші моделі складніше підтримувати. Поширений підхід полягає в обмеженні словника та відкиданні рідкісних слів, враховуючи, скажімо, 100 000 найпоширеніших слів у корпусі. Слова, які не є частиною словника, класифікуються як «unknown» та зіставляються зі спільним токеном UNK. Це означає, що деяка потенційно важлива інформація втрачатиметься в процесі токенизації слів, оскільки модель не матиме інформації про слова, пов'язані з UNK.

1.2.2 Токенизація символів

Найпростіша схема токенизації полягає в тому, щоб подавати кожен символ окремо в модель. У Python об'єкти str насправді є масивами, що дозволяє нам швидко реалізувати токенизацію на рівні символів лише декількома рядками коду:

```
text = "Tokenizing text is an essential task of NLP."  
tokenized_text = list(text)
```

```
print(tokenized_text)
```

Результат виконання функції токенизації символів відображено на рисунку 1.3.

```
['T', 'o', 'k', 'e', 'n', 'i', 'z', 'i', 'n', 'g', ' ', 't', 'e', 'x', 't',  
' ', 'i', 's', ' ', 'a', 'n', ' ', 'e', 's', 's', 'e', 'n', 't', 'i', 'a',  
'l', ' ', 't', 'a', 's', 'k', ' ', 'o', 'f', ' ', 'N', 'L', 'P', '.']
```

Рисунок 1.3 – Результат виконання функції токенизації символів

Токенізація символів також була створена для вирішення деяких проблем, пов'язаних з токенизацією слів. Замість розбиття тексту на слова, вона повністю розділяє текст на символи. Це дозволяє процесу токенизації зберігати інформацію про слова поза словниковим запасом, яку токенизація слів не може зберегти.

Токенізація символів не має таких самих проблем зі словниковим запасом, як токенизація слів, оскільки розмір «словникового запасу» становить лише стільки символів, скільки потрібно мові. Наприклад, для англійської мови словник токенизації символів матиме близько 26 символів.

Хоча токенизація символів вирішує проблеми OOV, вона не позбавлена від власних ускладнень. При розбитті навіть простих речень на символи замість слів, довжина виводу значно збільшується. Також вона додає додатковий крок до розуміння зв'язку між символами та значенням слів, а отже ще на один крок віддаляє від мети NLP – інтерпретації значення.

1.2.3 Токенізація підслів

Токенізація підслів – це техніка NLP, за якої слово розбивається на підслова або морфеми, і ці підслова відомі як токени. Ця техніка використовується у будь-якому завданні NLP, де модель повинна підтримувати великий словниковий запас і складні структури слів [7].

Концепція, що лежить в основі цього, полягає в тому, що слова, що часто зустрічаються, повинні бути в словниковому запасі, тоді як рідкісні слова розбиваються на поширені морфеми.

Токенізація підслів схожа на токенизацію слів, але вона трохи детальніше розбиває окремі слова за допомогою специфічних лінгвістичних правил. Одним з основних інструментів, які вони використовують, є розділення афіксів. Оскільки префікси, суфікси та інфікси змінюють власне значення слів, вони також можуть допомогти програмам зрозуміти функцію слова. Це може бути особливо цінним для слів поза словниковим запасом, оскільки ідентифікація афікса може дати програмі додаткове розуміння того, як функціонують невідомі слова.

Модель підслів шукатиме морфеми та розбиватиме слова, що їх містять, на окремі частини. Наприклад, запит «How are you doing lately?» буде розбитий на «how», «are», «you», «do», «ing», «late», «ly».

Як цей метод допомагає вирішити проблему слів поза словниковим запасом? Можливо, машина отримує складніше слово, таке як «machinating» (теперішній час дієслова «machinate», що означає планувати або брати участь у змовах). Малоймовірно, що «machinating» є словом, що входить до багатьох базових словників.

Якби модель NLP використовувала токенизацію слів, це слово було б просто перетворено на невідомий токен. Однак, якби модель NLP використовувала токенизацію підслів, вона змогла б розділити слово на токен «unknowing» та токен «ing». Звідси вона може робити цінні висновки про те, як слово функціонує в реченні.

Але яку інформацію може зібрати машина з одного суфікса? Наприклад, поширений суфікс «ing» функціонує кількома легко визначеними способами. Він може перетворити дієслово на іменник, як-от дієслово «build» перетворене на іменник «building». Він також може перетворити дієслово на дієприкметник теперішнього часу, як-от дієслово «run» стає «running».

Якщо моделі NLP надається ця інформація про суфікс «ing», вона може зробити кілька цінних висновків про будь-яке слово, яке використовує цю морфему. Якщо «ing» використовується в слові, вона знає, що воно функціонує або як дієслово, перетворене на іменник, або як дієслово теперішнього часу. Це різко звужує коло використання невідомого слова «machinating» в реченні.

1.2.4 Токенізація речень

Токенізація речень, як нескладно здогадатися з назви, являє собою розбиття фрагмента тексту на речення (лістинг 1.2). Тобто при токенизації абзац розбивається на речення, які і виступають в якості токенів. У багатьох задачах обробки природної мови розбиття текстових даних на речення є дуже корисним. Речення зазвичай розділяються крапкою, знаком оклику або знаком питання, а отже процес токенизації речень знаходить усі ці знаки у фрагменті тексту, щоб розділити дані на речення.

Лістинг 1.2 – Програмний код функції токенизації речень

```
from nltk.tokenize import sent_tokenize
text = "Hi! This is an example of the sentence tokenization
method. There should be three separate sentences in the
output."
print(sent_tokenize(text))
```

Результат виконання функції токенизації речень відображено на рисунку 1.4.

```
['Hi!', 'This is an example of the sentence tokenization method.',
'There should be three separate sentences in the output.']
```

Рисунок 1.4 – Результат виконання функції токенизації речень

Для літератури, журналістики та офіційних документів англійською добре працюють алгоритми токенизації, вбудовані в spaCy, оскільки токенизатор навчається на корпусі офіційного англійського тексту. Токенизатор речень працює гірше для електронних медичних записів, що містять скорочення, медичні терміни, просторові виміри та інші форми, яких немає в стандартній письмовій англійській мові.

Цю проблему адресує платформа ClarityNLP та намагається покращити результати токенизації речень для електронних медичних записів. Вона робить це, шукаючи типи текстових конструкцій, які заплутують токенизатор, та замінюючи їх окремими словами. Токенизатор речень не розділяє окреме слово, тому текст, що порушує правила, у формі заміни зберігається незмінним під час процесу токенизації. Після генерації окремих речень виконуються зворотні заміни, що відновлює оригінальний текст у наборі покращених речень. ClarityNLP також виконує додаткові виправлення речень для подальшого покращення результатів.

1.3 Виклики токенизації в NLP

Процес токенизації, попри свою базову технічну природу, стикається з низкою нетривіальних викликів, які мають як теоретичний, так і практичний характер. Однією з головних проблем є складність побудови універсальної токенизуючої системи, що однаково добре працюватиме з різними мовами, жанрами тексту, контекстами використання і навіть моделями машинного навчання. Причина полягає у великій кількості особливостей природної мови (рисунок 1.5), зокрема її неоднозначності, контекстуальності та постійної змінності.

У випадку мов із гнучким порядком слів і розвиненою морфологією, як-от українська чи турецька, токенизація стикається з проблемою визначення межі слова. Питання про те, де одне слово закінчується, а інше починається, не завжди має однозначну відповідь. Так, наприклад, у

складних словоформах, де є префікси, суфікси або частки, токенізатор має або виділити ці частини окремо, або залишити як одне слово. Обидва варіанти можуть бути виправданими в різному контексті, але обрати правильний у конкретному випадку часто буває складно без залучення додаткового мовного аналізу.



Рисунок 1.5 – Одне арабське слово дає значення шести різних слів в англійській мові

Варто зазначити, що чимало мовних одиниць є омонімами, тобто мають однакову форму, але різне значення. Наприклад, слово «ключ» в українській мові може означати предмет або джерело води. З погляду токенізатора, це один і той самий токен, проте для моделі, яка далі обробляє текст, контекстуальна різниця є критичною. Токенізація зазвичай не має доступу до семантики і працює лише з формальними ознаками, що є її суттєвим обмеженням. Особливо це стосується алгоритмів типу *whitespace-based*, які ділять текст лише за пробілами, не враховуючи ані пунктуації, ані морфологічної структури.

Складнощі також виникають у роботі з аглютинативними або ізолюючими мовами. У першому випадку, як-от у фінській або угорській, слова містять численні афікси, які додають нові відтінки значення. У другому – наприклад, у китайській мові – слова взагалі не мають чітких розділювачів, а кожен ієрогліф може бути як окремим словом, так і частиною фрази. Це створює серйозні труднощі для токенизації, яка не може спиратись на пробіли або знаки пунктуації (рисунок 1.6). У таких випадках необхідне попереднє навчання спеціалізованих сегментаторів, які аналізують текст на глибшому рівні.

Language	ISO	Text	↑ Num Tokens
English	en-US	what will the weather be next week	7
Spanish	es-ES	qué tiempo hará la semana que viene	8
Korean	ko-KR	다음 주 날씨 어때	12
Burmese	my-MM	နောက်တစ်ပတ် ရာသီဥတုဘယ်လိုရှိမလဲ	61
Amharic	am-ET	በሚቀጥለው ሳምንት ሳምንት አየሩ ምንድን	69

Рисунок 1.6 – Приклад повідомлення, перекладеного п'ятьма мовами, та відповідна кількість токенів, необхідних для його токенизації за допомогою токенизатора OpenAI

Ще одним викликом є поєднання багатьох мов в одному корпусі. У глобалізованому інформаційному просторі користувачі часто використовують різні мови у межах одного повідомлення. Такі тексти містять елементи код-перемикання, коли в одному реченні змішуються, скажімо, англійська і українська. Для токенизатора це означає необхідність ідентифікації мови кожного фрагмента і вибору відповідної логіки розбиття. Більшість сучасних токенизаторів не мають вбудованих механізмів мовної

ідентифікації, і тому виконують розбиття за єдиними правилами, що може призводити до спотворень.

Особливо чутливими до помилок токенізації є сучасні трансформерні моделі, що працюють на рівні субслів. Такі моделі, як BERT чи GPT, використовують спеціальні алгоритми побудови словників на основі частотних патернів – наприклад, Byte-Pair Encoding або WordPiece. Ці алгоритми не зважають на семантичні межі слів, а тому можуть розбивати знайомі слова на частини, які є статистично вигідними, але лінгвістично дивними. Це може призвести до втрати значущих одиниць або неправильного розуміння контексту, особливо у випадках нечастих слів або неологізмів.

Окремо варто звернути увагу на виклики, пов'язані з продуктивністю. Високоякісні токенізатори часто є повільними, оскільки здійснюють складний морфосинтаксичний аналіз. Навпаки, швидкі токенізатори можуть бути занадто грубими, пропускаючи значущі мовні особливості. У реальних системах важливо знайти баланс між швидкістю обробки і якістю розбиття, адже компроміс у цьому питанні безпосередньо впливає на точність усієї NLP-моделі. Це особливо критично в умовах обмежених обчислювальних ресурсів, наприклад у мобільних або вбудованих системах.

Слід також відзначити, що для доменоспецифічних текстів, таких як медична або юридична документація, загальні токенізатори часто виявляються неефективними. Вони не враховують специфічну термінологію, складні скорочення, латинські вставки або структури, властиві окремим галузям знань. Це вимагає створення адаптованих токенізаторів або донавчання існуючих на спеціалізованих корпусах. Інакше навіть найкращі мовні моделі можуть працювати неточно через неправильне розбиття ключових термінів.

Ще одним аспектом, який викликає труднощі, є неоднозначність пунктуації. У деяких випадках вона виконує роль розділювача, а в інших – є частиною токена. Наприклад, у скороченнях на кшталт «Dr.», «U.S.A.» або

у випадках лапок та дужок, стандартні правила токенизації можуть давати некоректні результати. Рішення цієї проблеми потребує контекстно-чутливого аналізу, який виходить за межі можливостей базового токенизатора.

Усі перелічені вище нюанси свідчать про те, що токенизація не є суто технічним завданням, яке можна вирішити один раз і назавжди. Це динамічний процес, який вимагає постійної адаптації до нових мовних явищ, змін у структурі текстів, нових алгоритмів машинного навчання та потреб прикладних систем. Попри те, що токенизація є лише початковим кроком у системах обробки природної мови, вона відіграє ключову роль у формуванні якості й ефективності всього NLP-процесу.

2 АНАЛІЗ СУЧАСНИХ ПІДХОДІВ ДО ТОКЕНІЗАЦІЇ В КОНТЕКСТІ NLP-МОДЕЛЕЙ

2.1 Архітектура сучасних NLP-систем

Сучасні системи обробки природної мови мають складну багаторівневу архітектуру, в основі якої лежить трансформація сирого тексту в представлення, придатне для машинного аналізу. Одним із фундаментальних етапів цього процесу є токенізація – процедура, що визначає спосіб розбиття суцільного тексту на менші одиниці, які називаються токенами. Ці токени у подальшому перетворюються на числові вектори, так звані ембедінги, та подаються на вхід нейромережевих моделей, таких як LSTM, GRU або трансформерів (BERT, GPT, T5, RoBERTa тощо). Токенізація задає основу подальшого ланцюга обробки, і тому її ефективність має безпосередній вплив на результати всього NLP-пайплайну.

Токенізатор виконує критичну функцію адаптації природномовного тексту до обмежень векторного простору. Наприклад, слово «нейропластичність» може бути незнайомим для моделі, яка була навчена на словнику, що його не включає. Завдяки підслівній токенізації, воно буде розбите на більш часті морфеми або навіть символи, які вже представлені в словнику. Це дозволяє уникнути втрати інформації у вигляді спеціального токена UNK, однак має інші наслідки – зокрема збільшення довжини послідовності та потенційну втрату семантичної цілісності.

Вибір токенізаційного механізму впливає на низку критичних параметрів NLP-системи:

- довжина послідовності або «sequence length». Моделі типу трансформер мають обмеження на довжину вхідної послідовності, зазвичай 512 або 1024 токени. Якщо токенізатор розбиває слова на дрібні частини, наприклад, символи або короткі морфеми, то навіть невеликий

текст може «переповнити» доступне вікно моделі, внаслідок чого частина інформації буде обрізана або знецінена через механізм attention masking;

– розмір словника та обсяг ембеддінгів. Великий словник передбачає більшу матрицю ембеддінгів – тобто більше параметрів, вищі апаратні вимоги та більші витрати пам'яті. З іншого боку, вузький словник призводить до фрагментації слів на менші одиниці, зменшуючи точність контекстного представлення. Ідеальний компроміс має бути досягнутий шляхом збалансованого підбору розміру словника, з урахуванням домену задачі, частоти вживання термінів, морфології мови тощо;

– швидкість навчання та інференсу. Від кількості токенів у тексті залежить кількість обчислень у шарі attention. Більше токенів – експоненційно більше обчислень ($O(n^2)$). Отже, більш агресивна токенізація й більше дроблення прямо впливає на продуктивність системи;

– точність виконання задачі (performance). У задачах на класифікацію, генерацію або переклад неправильна токенізація може призвести до втрати важливої морфологічної або семантичної інформації. Наприклад, невірне розбиття слова «відстороненість» на «від», «сторо», «неність» може знизити здатність моделі інтерпретувати слово як емоційне чи абстрактне поняття.

У архітектурах, що базуються на трансформерах, токенізатор є невід'ємною частиною моделі і не підлягає довільній заміні. Наприклад, BERT використовує токенізатор WordPiece, GPT-2 та RoBERTa – BPE. Кожен токенізатор оптимізовано під певний корпус і гіперпараметри тренування. Тому заміна токенізатора без перенавчання ембеддінгів або fine-tuning моделі призведе до некоректного відображення вхідних послідовностей і, як наслідок, погіршення результатів. Це породжує ключову залежність: модель і токенізатор мусять бути узгодженими. Зміна одного без відповідної адаптації іншого – ризикована практика, особливо при роботі з нестандартними або вузькоспеціалізованими корпусами.

Під час побудови універсальних NLP-моделей виникає потреба у створенні «універсального» токенізатора, здатного обробляти тексти різних

жанрів, мов, реєстрів. Універсальність досягається через використання підслівних моделей, які максимально зменшують кількість UNK-токенів і зберігають баланс між детальністю та узагальненням. Однак така універсальність часто шкодить специфічності. Доменні терміни, наприклад, юридичні чи біомедичні терміни, можуть бути невдало розбиті, спотворені або проігноровані.

У випадках із доменними задачами, як-от аналіз тональності в юридичних документах або розпізнавання намірів у технічній документації, доречно створювати кастомні токенізатори або принаймні ретренувати існуючі на специфічному корпусі. Це дозволяє зменшити розрив між поверхневою формою тексту та структурою вхідного представлення в моделі.

2.2 Порівняння найпоширеніших підслівних алгоритмів

Підслівні алгоритми токенізації лежать в основі сучасних систем обробки природної мови та забезпечують гнучке, адаптивне представлення тексту при обмеженому словнику. Найпоширенішими з них є Byte-Pair Encoding, WordPiece і Unigram Language Model. Ці методи мають суттєві архітектурні, алгоритмічні та статистичні відмінності, які прямо впливають на продуктивність моделей у задачах класифікації, генерації та розуміння тексту.

Byte-Pair Encoding є адаптацією однойменного алгоритму безконтекстного стиснення даних, яка перетворена на метод побудови словника токенів. Його принцип базується на жадібному ітеративному злитті найбільш частих пар символів або токенів у нові одиниці, які додаються до словника. Початковим словником є набір всіх символів, що зустрічаються в навчальному корпусі. На кожному кроці визначається пара сусідніх токенів із найвищою частотою і зливається в один новий токен,

після чого частотні статистики оновлюються. Процес повторюється до досягнення заданого розміру словника (рисунок 2.1).

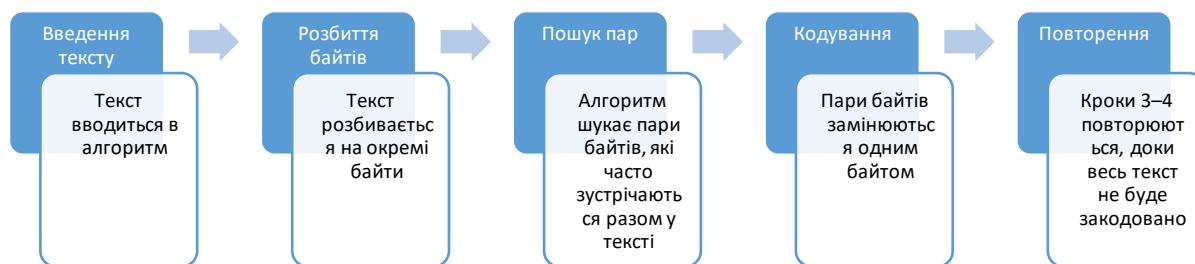


Рисунок 2.1 – Схема принципу роботи алгоритму BPE

Алгоритм не враховує морфологічної структури мови, натомість повністю спирається на статистичну частоту. Це робить його ефективним для різних мов і корпусів, проте призводить до того, що кінцеві токени можуть бути як валідними морфемами, так і випадковими комбінаціями символів. Структура словника в BPE є деревоподібною, що дозволяє реалізувати ефективне кодування через жадібний пошук найменшої кількості токенів, які охоплюють увесь вхідний текст. BPE дозволяє уникнути токенів, які модель не розуміє, за рахунок повного покриття усіх можливих послідовностей за допомогою найдрібніших одиниць – символів.

Для реалізації BPE використовуються лічильники частот токенів та структурований список пар, які ітеруються у порядку спадання частоти. При кожному злитті частоти відповідних пар оновлюються, що потребує ефективного доступу до кожної комбінації токенів. У бібліотеці Hugging Face цей алгоритм реалізовано як `tokenizers.models.BPE`, де модель підтримує побудову словника на основі поданого тексту або збереженням попередньо згенерованих файлів `merges.txt` та `vocab.json`.

На рисунку 2.2 наведений фасцит характерних особливостей алгоритму BPE.

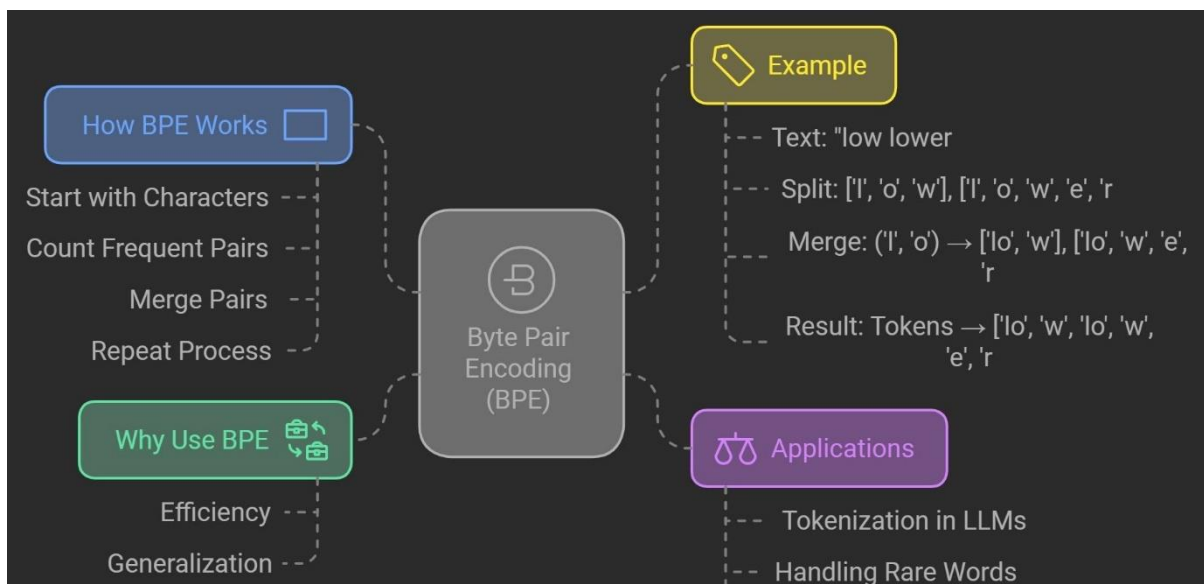


Рисунок 2.2 – Фасцит характерних особливостей алгоритму BPE

WordPiece є модифікацією BPE, адаптованою для задач природної мови, зокрема для моделі BERT. Основною відмінністю є те, що замість жадібного об'єднання найбільш частих пар, WordPiece враховує ймовірнісний приріст. На кожному кроці алгоритм об'єднує ту пару, яка максимізує приріст лог-правдоподібності моделі на основі language modeling-метрики. Таким чином нова одиниця має покращити здатність передбачити токени в контексті.

На відміну від BPE, WordPiece використовує префіксне кодування: частини слова, що не є початком слова, помічаються спеціальним символом «##». Наприклад, слово «playing» може бути розбите на «play» та «##ing». Розширений приклад токенизації WordPiece наведено на рисунку 2.3. Ця стратегія забезпечує збереження морфологічної структури, дозволяючи моделі легше виявляти суфікси, префікси та інші регулярності в мові. Також це зменшує неоднозначність при декодуванні.

Під час кодування WordPiece застосовує алгоритм жадібного пошуку із префіксним деревом. Кожен текст обробляється зліва направо, і на кожному кроці знаходиться найдовший можливий префікс, що присутній у словнику. Якщо токен не знайдено, використовується маркер [UNK].

Зокрема, в BERT для англійської мови було створено словник із приблизно 30 000 підслівних одиниць на основі корпусу Wikipedia та BookCorpus.

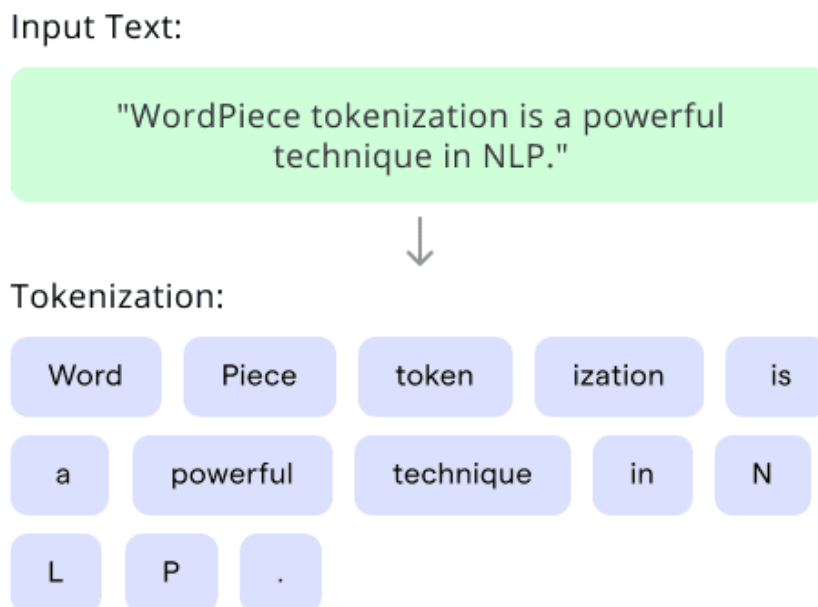


Рисунок 2.3 – Приклад токенизації WordPiece

У реалізації бібліотеки Hugging Face токенизатор WordPiece реалізований через `tokenizers.models.WordPiece`, що підтримує збереження словника у форматі `txt`, де кожен рядок відповідає одному токенові. Навчання відбувається з використанням `WordPieceTrainer`, який враховує не лише частотність, а і структуру позицій у словах. Однією з ключових особливостей WordPiece є збереження високої точності моделі при меншій кількості токенів. Це забезпечується за рахунок глибшої морфологічної інформативності, на відміну від символної або нативної підслівної сегментації.

Unigram Language Model є радикально іншим підходом до підслівної токенизації. Його концепція базується на імовірнісній моделі, де кожен текст розглядається як результат вибору з множини можливих токенизацій. Словник спочатку містить велику кількість кандидатів – частин слів, які

теоретично можуть слугувати токенами. Завдання полягає у знаходженні такого підмножини токенів, яке максимізує загальну ймовірність корпусу. В Unigram кожен токен має ймовірність, яка визначає, наскільки часто він буде використаний для сегментації. На відміну від BPE і WordPiece, сегментація не є жадібною. Натомість Unigram розглядає всі можливі розбиття тексту на токени і обирає те, яке має максимальну ймовірність. Ці ймовірності визначаються втратою, на якій навчається токенизатор. Припускаючи, що навчальні дані складаються зі слів x_1, \dots, x_N і що множина всіх можливих токенизацій для слова x_i визначена як $S(x_i)$, тоді загальна втрата визначається як:

$$L = - \sum_{i=1}^N \log \left(\sum_{x \in S(x_i)} p(x) \right). \quad (2.1)$$

Процес навчання Unigram базується на ітеративному зменшенні словника. Ініціальний словник містить багато надлишкових токенів, після чого найменш інформативні, чия частка в загальній ймовірності низька, відкидаються. При кожній ітерації модель перебудовує сегментації з оновленим словником, поки не буде досягнуто цільового розміру. Такий підхід дозволяє уникнути артефактів жадібного об'єднання, забезпечуючи гнучке імовірнісне представлення.

Unigram зберігає контекстно-незалежну природу, проте завдяки алгоритму Viterbi дозволяє виявити морфеми, що не є регулярними або частотними. Це особливо корисно в мовах із високим рівнем аглютинації, де стандартні жадібні алгоритми часто виявляються неефективними. Даний алгоритм базується на статистичному тренуванні через UnigramTrainer, де ключовим є параметр `unk_token`, який гарантує повне покриття тексту навіть за умов нестачі словникових одиниць. Метод підтримує формат файлів

SentencePiece, зокрема `.model` й `.vocab`, однак може бути експортований у формат Hugging Face JSON.

2.3 Обмеження готових токенизаторів

Сучасні підслівні токенизатори, як-от WordPiece, Byte-Pair Encoding та Unigram Language Model, широко використовуються як у прикладному моделюванні, так і в академічних дослідженнях. Однак застосування попередньо навчених токенизаторів, які постачаються разом із мовними моделями або відкритими бібліотеками, зумовлює низку методологічних та емпіричних обмежень. Серед ключових з них – упередженість, спричинена структурою великих словників, а також проблема відчуження корпусу дослідження, коли обрана модель токенизації не відповідає статистичним, стилістичним або термінологічним особливостям конкретного набору даних. Обидва аспекти мають критичний вплив на якість аналізу текстів, точність класифікації, генерації та загальну інтерпретованість моделей.

Готові токенизатори, зокрема ті, що використовуються у BERT, RoBERTa або GPT, створюються на основі великих корпусів загального призначення. Такі словники містять десятки тисяч підслівних одиниць, які статистично найбільш вживані у відповідному корпусі. Наприклад, словник BERT створено на основі BookCorpus та англійської Вікіпедії, а GPT-2 – на основі веб-контенту з OpenWebText. Це зумовлює неоднорідний розподіл токенів: частотні слова та морфеми представлені окремими одиницями, тоді як менш вживані комбінації зливаються у довші ланцюжки символів.

Результатом є структурна перевага для лексем, що мають високу частоту в базовому корпусі, навіть якщо ці лексеми є семантично або стилістично неактуальними для нового завдання. Наприклад, у юридичних або біомедичних текстах численні терміни відсутні в словнику моделі або представлені довгими послідовностями токенів. Це знижує ефективність

подання і збільшує кількість токенів, необхідних для одного прикладу, що ускладнює навчання та зменшує семантичну компактність.

Окремо постає питання морфологічної адекватності. Упередженість словника проявляється у некоректній фрагментації рідкісних або складених слів. Наприклад, в медичному корпусі термін «immunoglobulin» може бути розбитий на «imm», «upo», «g», «lob», «ulin», що не має жодної семантичної відповідності. Таким чином, на вході до моделі подається надлишкова кількість токенів, які не корелюють із структурою значення. Це ускладнює навчання моделей через зростання ентропії представлення.

Токенізаційна упередженість також має вплив на баланс вагів при навчанні. Часті токени, які присутні в словнику як цілісні одиниці, мають коротші послідовності та нижчу ентропію. У той час рідкісні або термінологічні конструкції представлені фрагментовано, що призводить до появи довших контекстів і зростання складності моделювання взаємозв'язків. Це створює нерівномірність у зважуванні зразків, зумовлює переобучення на частих структурах і зменшує чутливість моделі до лексичних відмінностей.

Великі словники готових токенізаторів мають статичну структуру. Це означає, що навіть у випадку донавчання на новому домені структура токенізації залишається незмінною. Результатом є неспівпадіння між внутрішнім представленням тексту і специфікою даних. Навіть fine-tuning моделі не здатен компенсувати це зміщення, оскільки воно виникає ще до початку обчислення градієнтів, на рівні препроцесингу.

Під терміном відчуження корпусу задачі мається на увазі відсутність відповідності між структурою токенізатора і статистичними властивостями конкретного набору даних. Це явище особливо критичне для вузькодоменних корпусів – медичних, юридичних, технічних або соціолінгвістичних. Коли токенізатор, натренований на універсальному корпусі, застосовується до специфічного тексту, він не враховує лексичні особливості, синтаксичні конструкції або морфеми, притаманні лише цьому

піджанру. Результатом є деструктивна фрагментація термінів, неправильне позиціонування токенів у моделі, збільшення вхідної довжини і загальна деградація якості представлення.

Вплив цього ефекту найкраще видно при аналізі семантично складних одиниць. Наприклад, у мові науки і техніки активно використовуються складені слова з латинськими або грецькими коренями. Якщо токенізатор не має жодної з частин у словнику, слово буде розбите на випадкові послідовності символів. Таке розбиття повністю руйнує ієрархічну структуру значення, перетворюючи модель на символічний процесор без доступу до морфологічного змісту.

У прикладних дослідженнях, де задачі належать до сфери класифікації доменних текстів, генерації пояснень або визначення настрою, використання чужого токенізатора призводить до того, що система оперує спотвореними векторами. Часто це не видно під час валідації, однак під час deployment виникають серйозні проблеми узагальнення, що пов'язано саме з відчуженням на рівні токенізації.

Важливою формою прояву цього ефекту є те, що навіть лексично однакові слова в різних доменах можуть мати різне токенізаційне представлення. У корпусі новин слово «trial» може бути представлено як одиничний токен, тоді як у юридичному корпусі – розбите на «tri», «##al», оскільки у тренувальному корпусі такі поєднання мають іншу статистику. Це призводить до втрати узгодженості між семантикою і токенізацією.

Додаткову складність створює відсутність механізмів адаптації токенізатора до нового корпусу. Оскільки структури, як-от WordPiece чи BPE, не змінюються без повторного навчання, користувач змушений або працювати зі словником, який неадекватний до задачі, або витратити ресурси на повторне навчання всього токенізатора. Це ускладнює дослідження в нових мовах, де готові токенізатори або недоступні, або мають занадто загальний характер.

У мовах з багатою морфологією, наприклад, у фінській або українській, цей ефект посилюється через велику кількість слівформ. Готові токенизатори, що створені на англійському корпусі, виявляються нездатними належно фрагментувати відмінювані, дієвідмінювані чи префіксальні структури. Результатом є надмірна сегментація, втрата морфемної структури та зниження якості контекстного представлення.

Окрему увагу варто звернути на те, як відчуження впливає на тренування attention-механізмів у трансформерах. Погана токенизація збільшує кількість токенів, призводячи до розпорошення уваги, зменшення фокусування на ключових структурах і деградації представлень високого рівня. Це особливо важливо в моделях з обмеженням на довжину послідовності, де кожен додатковий токен знижує ємність моделі.

2.4 Аналіз публікацій, що ізольовано досліджують токенизацію

У роботі [8] Тораман та співавтори дослідили п'ять методів токенизації, включно з морфологічним рівнем, застосованих до турецького корпусу OSCAR, який містить аглютинативну морфологію. Вони навчили окремі RoBERTa-моделі з кожним токенизатором і провели фін-тюнінг на шести downstream-завданнях, включно з класифікацією та NER. Було визначено, що морфологічний токенизатор не поступається BPE й WordPiece, а збільшення словника принесло помітні переваги у продуктивності саме для методів, що враховують морфологічну структуру, у межах оптимального співвідношення словник/параметри $\approx 40\%$

У роботі [9] Бостром і Дюрретт провели пряме зіставлення BPE та Unigram LM для англійської та японської, застосувавши однакові моделі-трансформери під масковане навчання. Аналіз показав, що Unigram LM витягує підслова, які краще вирівнюються з морфемами, і в downstream-задачах демонструє однаково високі або вищі результати

порівняно з ВРЕ. Механізм жадібного ВРЕ був ідентифікований як джерело артефактної сегментації

У роботі [10], де вони навчають 24 моно- та мультимовні LLM із різними токенизаторами на 2.6 B параметрів, дослідники показали, що вибір токенизатора впливає не лише на якість, а й на вартість тренування. Стандартні метрики fertility і parity не забезпечують кореляцію з downstream результатом, що особливо помітно у мультимовних конфігураціях – ієрархія англоорієнтованих токенизаторів призводить до поглиблення деградації у багатомовних моделях на 68 %

У роботі [11] систематично досліджувався зв'язок між вбудованою метрикою «compression» і downstream-продуктивністю. Зі зменшенням розміру підтримуваних документів compression знижується через довші токенизаційні послідовності, і одночасно падає якість у задачах генерації та класифікації. Найсильніший ефект показали моделі задачі генерації, а найменші моделі – найбільше чутливі до поганої токенизації .

У більшості прямих досліджень, таких як [15], фокус на простих токенизаційних стратегіях був не співвіднесений із downstream-точністю на ретельному рівні; порівняння обмежувалися загальними трендами, без кількісних експериментів із різними токенизаторами в уніфікованих умовах.

Дослідження з обробки арабської [13] провели оцінку токенизаторів, застосувавши шість різних стратегій до трьох текстових задач. Проте їхній аналіз був фокусований на метриках покриття словника, каденції та швидкості; реальний вплив на downstream точність представлений лише як поверхневий результат без глибокої порівняльної статистики.

У дослідженні [14] було описано гібрид морфологічної та ВРЕ токенизації для корейської, та продемонстровано, що ця комбінація перевершує ВРЕ у соло в машинному перекладі й NLU задачах. Проте згадані результати спираються на якісні приклади та прості метрики точності без кореляції між якісністю сегментації та загальною пропускну здатністю моделі.

Більшість досліджень або ізольовано порівнюють кілька токенизаторів на downstream-задачах без уніфікованого контролю за модельною архітектурою, або ж зосереджені на метриках якості сегментації compression й purity, не показуючи їх кількісного впливу на продуктивність. Пряме дозавантаження моделей з однаковими архітектурами, але з різними токенизаторами, як у [8], [9] та [10], свідчить про реальний, інколи значний, вплив токенизації. Незважаючи на це, великий обсяг літератури лишається поверхневим чи описовим, без точних цифр.

3 РЕАЛІЗАЦІЯ ПРОГРАМНОЇ ЧАСТИНИ

3.1 Вибір засобів програмування

У рамках виконання даної кваліфікаційної роботи основною мовою програмування було обрано Python. Це рішення є виваженим як з точки зору технологічної доцільності, так і з урахуванням широкої підтримки інструментарію для задач природної мовної обробки, машинного навчання та експериментального дослідження.

Python на сьогоднішній день є однією з найпоширеніших мов програмування у сфері штучного інтелекту, глибинного навчання, роботи з текстовими даними та наукових досліджень у галузі обробки природної мови. Його популярність пояснюється низкою ключових переваг, які мають безпосереднє відношення до тематики цієї кваліфікаційної роботи.

По-перше, він має надзвичайно лаконічний і читаємий синтаксис, що дозволяє зосередитися безпосередньо на логіці задачі, а не на технічних деталях реалізації. Такий підхід суттєво прискорює прототипування, тестування гіпотез та проведення експериментів, що є критично важливим у науковій роботі. Саме завдяки його простоті, стало можливим швидко реалізувати повний цикл експериментального дослідження – від попередньої обробки текстів до тренування й оцінювання моделей.

Іншою з головних причин вибору Python є його винятково багата екосистема спеціалізованих бібліотек, зокрема:

– HuggingFace Transformers – одна з провідних бібліотек для роботи з сучасними NLP-моделями, містить також модулі для токенізації, тренування, тонкого налаштування моделей [15];

– Tokenizers – високопродуктивна бібліотека, частково реалізована на Rust, яка дозволяє ефективно створювати, тренувати та експортувати токенізатори типу BPE, WordPiece, та Unigram [16];

- Torch / PyTorch – одна з найгнучкіших та найпопулярніших бібліотек для реалізації нейронних мереж. Має високий рівень абстракції та дозволяє повністю контролювати архітектуру моделі, оптимізацію, тренування і тестування;

- Datasets – бібліотека для зручного завантаження та роботи з відкритими NLP-корпусами, як-от IMDb, SST, AG News тощо;

- TQDM, NumPy – допоміжні бібліотеки для організації виводу прогресу тренування та обчислень.

Усі ці бібліотеки активно підтримуються спільнотою, мають докладну документацію та відкритий код, що дає змогу інтегрувати їх у будь-який дослідницький пайплайн. До того ж, ця мова програмування є де-факто стандартом у науковій спільноті з NLP. Більшість сучасних публікацій, статей, дослідницьких проєктів та відкритих кодових баз реалізовані саме на ній. Це дає змогу порівнювати отримані результати з аналогами, легко перевіряти та відтворювати існуючі підходи й інтегрувати сторонні рішення у власні експерименти.

Альтернативними мовами могли би виступати, зокрема C ++, Java та R. C++ хоча й забезпечує високу продуктивність, але має надто низький рівень абстракції для швидкої реалізації експериментів, а робота з текстовими даними була б доволі громіздкою. Java має бібліотеки для NLP, наприклад, Stanford CoreNLP, однак переважно орієнтована на виробничі системи, а не на гнучкі дослідницькі проєкти. Мова R потужна у статистичному аналізі, але значно поступається Python у глибокому навчанні та сучасних NLP-інструментах.

Було зроблено висновок, що жодна з альтернатив не забезпечує таке поєднання гнучкості, глибини інструментарію та зручності роботи, яке пропонує Python у контексті реалізованого завдання.

Проєкт реалізовано у середовищі Visual Studio Code, яке вважається найзручнішим і найбільш підтримуваним редактором коду для Python. Завдяки інтеграції з системами контролю версій Git, дебагером, терміналом

і підтримкою Jupyter Notebook, VS Code дозволяє ефективно організувати роботу над дослідженням.

3.2 Аналіз та вибір алгоритмів розроблення

Під час розв'язання задачі дослідження впливу різних методів токенизації на продуктивність моделі класифікації текстів особливої уваги заслуговує етап вибору алгоритмів, які беруть участь як у попередній обробці, так і у побудові та навчанні нейронної моделі. Хоча у цій роботі реалізовано рішення на базі простого класифікатора та трьох типів токенизаторів – байткового злиття, підслівного розщеплення та уніграмної моделі – існує широкий спектр альтернативних підходів, кожен з яких має свої переваги, недоліки та сфери застосування. Вибір оптимального поєднання алгоритмів тісно пов'язаний із цілями дослідження, природою вхідних даних, обмеженнями на обчислювальні ресурси та вимогами до точності або швидкості роботи системи.

Одним з найпростіших, але все ще застосовуваних методів у задачах обробки текстів є токенизація на рівні слів. Алгоритмічно це означає розбиття вхідного речення на лексеми згідно з пробілами, іноді з урахуванням знаків пунктуації. Такий підхід легко реалізується та забезпечує зрозумілу інтерпретацію, проте має низку критичних обмежень. Він надто чутливий до орфографічних варіацій, неефективно працює з аглютинативними та флективними мовами, а також не дозволяє будувати компактні словники. Словоорієнтована токенизація ускладнює обробку рідкісних або нових слів, оскільки кожне нове слово потенційно розглядається як унікальний токен, чого немає у підходах зі субсловною сегментацією. Крім того, словниковий розмір у такому випадку суттєво збільшується, що призводить до зростання параметрів моделі.

Іншим напрямом є символна токенизація, де кожна літера або символ є окремим токеном. Цей метод гарантує повне покриття будь-якого тексту

незалежно від мови, орфографії чи пунктуації. Проте через надмірну дрібність сегментації символні токени не несуть достатньо інформації самостійно. Відповідно, моделі, які працюють з такими токенами, мають складнішу архітектуру та потребують довших послідовностей вхідних даних, що негативно позначається на ефективності тренування. У випадках, коли необхідно аналізувати морфологічні особливості або структуру нових слів, символна токенизація може бути доцільною, однак вона значною мірою програє у задачах класифікації довгих текстів, де важливу роль відіграє семантична компактність.

Альтернативним напрямом розвитку могли б стати морфологічні токенизатори. Такі алгоритми не обмежуються формальним розбиттям тексту, а залучають лінгвістичні знання про мову. Зокрема, вони аналізують слова за частинами мови, визначають корені, афікси, флексії. Цей підхід актуальний для мов зі складною морфологією, як-от фінська, турецька, українська. У випадку англійської мови, яка має значно біднішу морфологічну структуру, переваги такого методу є обмеженими. Зазвичай застосовуються в лінгвістичних дослідженнях, пошукових системах, а також при побудові машинних перекладачів або систем автоматичного аналізу синтаксичної структури. Їх головний недолік полягає в тому, що вони вимагають наявності граматичних словників, правил узгодження та граматичних теггерів, що суттєво ускладнює реалізацію та обмежує масштабованість.

Особливе місце у сучасному NLP посідають статистичні субсловні токенизатори, серед яких можна виділити методи частотного злиття байтів, розбиття за допомогою частотних гаплографій, як у WordPiece, та уніграмні мовні моделі. Ці методи стали стандартом у сучасних трансформерних архітектурах, оскільки дозволяють ефективно збалансувати між довжиною вхідної послідовності, розміром словника та семантичною виразністю токенів. Вони генерують сегменти слів, які є найчастотнішими в корпусі, і таким чином забезпечують ефективне представлення навіть рідкісних слів

за рахунок спільних морфем. Саме ці токенизатори використано в рамках дослідження, оскільки їх порівняльний аналіз дозволяє оцінити, як саме структура токенизації впливає на продуктивність моделі.

Щодо альтернатив побудови моделі, варто згадати про можливість застосування передтренуваних трансформерних моделей, таких як BERT або DistilBERT. Ці архітектури, попередньо натреновані на великих корпусах, демонструють високу точність в більшості задач класифікації текстів. Вони зазвичай містять власні оптимізовані токенизатори та використовують складну внутрішню архітектуру з багат шаровими механізмами самоуваги. Однак, використання таких моделей потребує значних обчислювальних ресурсів, часу на тонке налаштування і менш придатне для ізольованого дослідження впливу лише методу токенизації. У цьому контексті застосування простої моделі класифікатора із вектором середнього ембеддингу дозволяє винести токенизацію як основний дослідницький фактор і зменшити вплив складної архітектури на результати.

Альтернативним підходом є використання моделей на основі згорткових або рекурентних нейронних мереж. Наприклад, приєднання LSTM-рівня до ембеддингу дозволяє моделі вловлювати локальну залежність між токенами. Проте такі моделі гірше масштабуються при зростанні довжини послідовностей, важко паралелізуються і значно поступаються трансформерним архітектурам за швидкістю. Крім того, у задачах аналізу токенизаторів їхній внесок у результат є непрозорим, оскільки рекурентні шари створюють додаткову внутрішню структуру залежностей, що ускладнює інтерпретацію змін ефективності.

Можна було також розглянути підходи, засновані на стохастичних або байєсівських моделях, зокрема методи тематичного моделювання, наприклад Latent Dirichlet Allocation. У таких випадках токенизація також відіграє важливу роль, оскільки визначає структуру вхідних даних, однак їхня сфера застосування обмежується задачами кластеризації чи побудови

тематичних просторів, що не відповідає класифікаційному фокусу цього дослідження.

Таким чином, у процесі розроблення системи було свідомо обрано лаконічну архітектуру, в якій основний акцент зроблено на токенизацію. Вибрані алгоритми побудови токенизаторів базуються на сучасних та репрезентативних статистичних методах, які широко застосовуються у промислових і дослідницьких NLP-системах. Усі вони дозволяють побудувати компактний, семантично узгоджений словник, який зберігає здатність до генералізації при обробці нових текстів.

Вибір простої моделі класифікації, що використовує ембеддинг і середнє згортання векторів, забезпечує мінімальне втручання у процес представлення токенів у векторному просторі. Це дозволяє точніше простежити зв'язок між структурою токенизації та кінцевою продуктивністю моделі. Таким чином, дана експериментальна конфігурація дозволяє досягти поставленої мети дослідження – ізольовано оцінити та порівняти ефективність різних методів токенизації у контексті задач природної мовної обробки текстів.

3.3 Вибір та завантаження бібліотек

Фреймворком, що ліг в основу побудови та тренування моделі, є бібліотека PyTorch (рисунок 3.1). Вона забезпечує необхідний інструментарій для створення нейронних мереж, роботи з GPU, оптимізації вагів та обчислення градієнтів за допомогою зворотного поширення помилки. Зокрема, модуль `torch.nn` дозволив описати архітектуру моделі класифікатора, в якій було використано ембеддинг-шар для перетворення індексів токенів у векторний простір, а також двошарову послідовність повнозв'язних шарів з функцією активації ReLU для побудови прогнозу. Модуль `torch.utils.data` використовувався для створення об'єктів `DataLoader`,

що відповідальні за ефективне завантаження батчів з попередньо токенозованих даних під час тренування.

```
import torch
import torch.nn as nn
from torch.utils.data import DataLoader
```

Рисунок 3.1 – Імпорт модулів бібліотеки PyTorch

Бібліотека `transformers` від компанії Hugging Face, а саме клас `PreTrainedTokenizerFast`, використовувалась для завантаження токенизаторів, попередньо збережених у форматі JSON (рисунок 3.2). Даний клас дозволяє інтерфейсно працювати з токенизаторами, які були навчені за допомогою бібліотеки `Tokenizers`. Він забезпечує високопродуктивну реалізацію токенизації та підтримує ключові параметри, такі як додавання спеціальних токенів для паддінгу, невідомих слів або початку й кінця речень. Його використання у зв'язці з JSON-файлами дозволило зберегти незалежність від зовнішніх API моделей і зосередитись лише на зміні стратегії токенизації.

```
from transformers import PreTrainedTokenizerFast
```

Рисунок 3.2 – Імпорт модулів бібліотеки `transformers`

Пакет `datasets` від Hugging Face застосовувався для завантаження текстового корпусу IMDB (рисунок 3.3). Цей набір даних містить тисячі рецензій на фільми, кожна з яких має мітку, що вказує на її сентимент. Завдяки функціоналу `datasets` стало можливим швидко отримати доступ до структурованих даних, виконати їх випадкове перемішування та відбір підмножини для навчання та тестування. Також `datasets` підтримує функцію

map, яка дозволила застосувати токенизатори до всіх елементів масиву векторизовано, без написання додаткового циклу обробки.

```
from datasets import load_dataset
```

Рисунок 3.3 – Імпорт модулів бібліотеки datasets

Для покращення візуалізації процесу навчання використовувався модуль tqdm (рисунок 3.4). Він додає прогрес-бари під час проходження по батчах у навчальному циклі, що дозволяє наочно спостерігати за перебігом навчання моделі, кількістю виконаних ітерацій та поточним станом обробки.

```
from tqdm import tqdm
```

Рисунок 3.4 – Імпорт модулю tqdm

Для обробки тексту і тренування кастомних токенизаторів було використано бібліотеку Tokenizers (рисунок 3.5). Вона надає гнучкий і високошвидкісний інтерфейс для побудови різних моделей токенизації, таких як Byte-Pair Encoding, WordPiece і Unigram. Її компоненти, зокрема Tokenizer, Trainers, Models, Normalizers та PreTokenizers, були використані у скрипті train_tokenizers.py. Компонент Tokenizer дозволяє ініціалізувати об'єкт токенизатора з певною моделлю сегментації. Модулі Models реалізують базові алгоритми токенизації. Для BPE, WordPiece та Unigram було відповідно створено моделі BPE, WordPiece та Unigram з відповідним токеном невідомого слова. Модулі Trainers дозволили навчитись на базі великого корпусу тексту та створити словник із вказаною максимальною кількістю токенів.

Зокрема, BpeTrainer, WordPieceTrainer та UnigramTrainer мають параметри для визначення спеціальних токенів, розміру словника та інших характеристик. Нормалізатори, такі як послідовність із NFD нормалізації, перетворення до нижнього регістру та видалення акцентів, дозволили забезпечити однаковість входу до тренування. Модуль pre_tokenizers, який включав розділення за пробілами, виконував попередню токенизацію для збереження структури слів перед поділом на підслова.

```
from tokenizers import Tokenizer
from tokenizers.normalizers import NFD, Lowercase, StripAccents, Sequence
from tokenizers.pre_tokenizers import Whitespace
from tokenizers.trainers import BpeTrainer, WordPieceTrainer, UnigramTrainer
from tokenizers.models import BPE, WordPiece, Unigram
```

Рисунок 3.5 – Імпорт модулів бібліотеки Tokenizers

Модуль random з бібліотеки datasets, реалізований через функцію shuffle, використовувався для випадкового перемішування датасету, що дозволяє зменшити вплив порядку зразків на процес навчання. Це є критично важливим у навчанні на обмеженій підвибірці даних, адже так забезпечується більш рівномірне представлення різних типів прикладів у кожній епосі.

3.4 Архітектура класифікаційної моделі

Реалізацію базової нейронної мережі для класифікації текстів було відокремлено у файл model.py. Він складається з класу Classifier, що успадковується від базового класу torch.nn.Module, який є стандартним шаблоном для створення моделей у фреймворку PyTorch. Структура моделі достатньо потужна для дослідження впливу різних стратегій токенизації на якість класифікації. Основна логіка моделі полягає в обробці послідовностей індексів токенів за допомогою шару ембеддингу, агрегації

векторів та проходженні їх через кілька щільних шарів для отримання ймовірностей класів.

Спочатку імпортуємо основні компоненти з PyTorch. Пакет `torch` використовується для створення тензорів і керування обчисленнями на GPU або CPU, а модуль `torch.nn` надає доступ до високорівневих елементів архітектури моделі, включаючи шари, функції активації та структуру моделі.

```
class Classifier(nn.Module):  
    def __init__(self, vocab_size, emb_dim=128):  
        super().__init__()
```

У оголошенні класу `Classifier` можна побачити, що модель є підкласом `nn.Module`. Метод `__init__` є конструктором класу. Він приймає два параметри: `vocab_size`, який задається динамічно в процесі тренування залежно від розміру словника токенизатора, та `emb_dim`, який визначає розмірність векторів у просторі ембеддингів. Значення 128 для `emb_dim` є типовим і забезпечує прийнятний баланс між обчислювальними витратами та якістю представлення слів.

```
self.embedding = nn.Embedding(vocab_size, emb_dim,  
padding_idx=0)
```

Вищенаведений рядок створює ембеддинг-шар, який трансформує цілочисельні індекси токенів у вектори фіксованої довжини. Параметр `vocab_size` задає кількість можливих унікальних токенів, які можуть бути закодовані, а `emb_dim` визначає кількість компонентів у кожному векторі. Вказівка `padding_idx=0` означає, що токен з індексом 0 не буде врахований під час оновлення вагів, що важливо для збереження семантичної нейтральності токенів заповнення.

Наступним етапом стало визначення класифікатора (лістинг 3.1), який реалізований як послідовність шарів за допомогою `nn.Sequential`. Перший лінійний шар зменшує розмірність з `emb_dim` до 64. Потім застосовується функція активації ReLU, яка вносить нелінійність у модель, усуваючи проблему лінійної лімітації. Наступний лінійний шар зменшує розмірність

простору до двох компонент, які інтерпретуються як логіти для двох класів задачі бінарної класифікації. Вивід цього шару не містить активації, оскільки функція `CrossEntropyLoss`, яка використовується під час навчання, самостійно включає обчислення `softmax`.

Лістинг 3.1 – Програмний код визначення класифікатора

```
self.classifier = nn.Sequential(  
    nn.Linear(emb_dim, 64),  
    nn.ReLU(),  
    nn.Linear(64, 2)  
)
```

Метод `forward` (лістинг 3.2) визначає порядок проходження даних через модель під час прямого поширення. Він є обов'язковим для всіх підкласів `nn.Module`. Вхідним параметром є `input_ids`, тобто тензор із індексами токенів, які були згенеровані токенизатором. Після проходження через ембеддинг-шар змінна `embedded` набуває розміру `[batch_size, sequence_length, emb_dim]`, де кожен токен представлений як вектор у просторі розмірності 128.

Лістинг 3.2 – Реалізація методу `forward`

```
def forward(self, input_ids):  
    embedded = self.embedding(input_ids)  
    pooled = embedded.mean(dim=1)  
    return self.classifier(pooled)
```

Для перетворення вхідної послідовності у фіксоване представлення на рівні всього речення, застосовується операція осереднення по осі `dim=1`, тобто по всій довжині послідовності. Таким чином, кожне речення репрезентується як один вектор розмірності `emb_dim`. Цей підхід називається `Mean Pooling` і є ефективним методом агрегації ембеддингів при

обробці коротких текстів. Результиуючий тензор `pooled` надсилається на вхід класифікатора, де послідовність шарів обробляє його та повертає логіти, що використовуються для обчислення втрат та передбачень класу.

Модель побудована з урахуванням простоти, відтворюваності та гнучкості для аналізу ефективності різних токенизаторів. Її архітектура забезпечує мінімум змінних чинників при порівнянні якості токенизації, оскільки вся варіативність обмежується лише способом перетворення тексту у числові індекси. Така конструкція дозволяє чітко ізолювати вплив вибору методу токенизації на кінцеву продуктивність моделі без накладення додаткових ефектів складної архітектури, що дозволяє зробити об'єктивні висновки щодо якості кожного методу токенизації у контексті задач класифікації текстів.

3.5 Підготовка текстових даних для навчання токенизаторів

Файл `prepare_text.py` виконує ключову роль у формуванні корпусу для навчання власних токенизаторів, які в подальшому використовуються в процесі передобробки даних для класифікації. Саме на цьому етапі виконується попередня підготовка текстів, їх очищення, формування єдиного текстового файлу, який потім виступає джерелом навчальних прикладів для алгоритмів сегментації, зокрема BPE, WordPiece та Unigram. Логіка реалізації у цьому скрипті виконує трансформацію корпусу даних з формату датасету у формат послідовного текстового файлу.

```
from datasets import load_dataset.
```

Перший рядок імпортує функцію `load_dataset` з бібліотеки `datasets`, яка розроблена командою Hugging Face. Вона дозволяє зручно завантажувати і працювати з відкритими корпусами даних без необхідності ручного завантаження, зберігання та парсингу. У цій кваліфікаційній роботі використовується набір даних IMDb, який містить велику кількість рецензій

до фільмів, позначених як позитивні або негативні. Даний набір є одним із найпоширеніших у задачах бінарної класифікації тексту.

```
dataset =
load_dataset("imdb")["train"].shuffle(seed=42).select(range(25
000)).
```

У цьому рядку виконується завантаження корпусу IMDb, зокрема підмножини `train`, яка містить 25000 прикладів. Далі ця множина зазнає перемішування з використанням фіксованого значення генератора випадкових чисел, що гарантує відтворюваність результатів під час повторного запуску коду. Хоча `train` частина вже містить 25000 прикладів, явним використанням `select(range(25000))` ми фіксуємо цей обсяг і захищаємо від можливих змін у майбутніх версіях датасету.

У наступному блоці відкривається файл `train_text.txt` у режимі запису з кодуванням UTF-8. Це важливо для збереження символів, що не входять до стандартної ASCII-таблиці, адже в кінотекстах часто трапляються апострофи, лапки, акценти та інші спеціальні символи. Ітеруючись по кожному прикладу у корпусі, скрипт витягує текстову частину, яка зберігається у полі `text`.

Далі виконується серія заміन у кожному рядку. Видаляються HTML-розмітки `

`, які зазвичай використовуються для поділу абзаців і на практиці не несуть семантичного навантаження. Також видаляються символи переведення рядка та подвійні пробіли, які можуть виникати внаслідок попередніх замін (лістинг 3.3). Результатом є чистий, лаконічний текст без розмітки, який готовий до подальшої токенизації (рисунок 3.6).

Лістинг 3.3 – Обробка тексту

```
with open("train_text.txt", "w", encoding="utf-8") as f:
    for example in dataset:
        f.write(example["text"].replace("<br /><br />", "
").replace("\n", " ").replace("  ", " ") + "\n")
```

```

train_text.txt
24978 Unwatchable. You can't even make it past the first three minutes. And this is coming from a huge Adam Sandler
24979 John Carradine, John Ireland, and Faith Domergue who as players all saw better days in better films got togeth
24980 This Night Listener is better than people are generally saying. It has weaknesses, and it seems to be having a
24981 When I started to watch this movie on VH-1 I cringed. The MTV movies were all bad so I wasn't expecting much. E
24982 I imagine that the young people involved in the making of "Necromancy" (aka "The Witching" plus a bunch of oth
24983 Ted V. Mikels's film Corpse Grinders 2 is 103 minutes of excruciating cinematic swill. The plot is pretty much
24984 This movie is not as horrible as most Sci-Fi Channel movies. I am used to seeing the gray CGI blobs and the an
24985 Besides the fact that it was one of the few movies that I ever shed a tear over (bye-bye manhood), this is one
24986 The main aspect about the Superstar's movies at his later stages were the frequency, the lacuna between one mc
24987 Saw this movie at the Vancouver Film Festival and thought it was deadly smart, stylish, and FUNNY. The cast wa
24988 A young boy sees his mother getting killed and his father hanging himself. 20 years later he gets a bunch of f
24989 Written by the excellent McGovern and directed by Freaars this film was a slight disappointment. It seemed too
24990 I was in second grade, 12 years ago. I remember it clearly. We were learning about space. All little kids want
24991 I loved this movie. Great storyline and actors and good movie sets. It told the story in a way I can easily ur
24992 Sure, it's a 50's drive-in special, but don't let that fool you. In my little book, there are a number of inte
24993 What is this crap? My little cousin picked this out obviously for the overly girlie DVD art and title... I dec
24994 Shortly after seeing this film I questioned the mental competence of every actor and actress that accepted a r
24995 RKO had a reputation for making folksy, homespun pieces of Americana. Anne Shirley (as Dawn O'Day) had been ir
24996 The ultimate goal of Big Brother, that we know what to think before we think it, has been realized. Is it some
24997 After mob boss Vic Moretti (late great Anthony Franciosa) kills his lady whom has been cheating on him with De
24998 Anyone who has said that it's better than Hostel is talking complete crap, believe me I'm not a fan of Hostel
24999 What do you get if you cross The Matrix with The Truman Show? I'm sure you've all seen The Matrix by now. The
25000 I remember watching this movie several times as a very young kid, and there were parts of it (many in fact) th

```

Рисунок 3.6 – Фрагмент тексту після обробки

Кожен очищений рядок записується у файл окремим рядком, створюючи послідовність незалежних текстових прикладів. Такий формат є найбільш зручним для подальшого навчання токенізаторів, оскільки багато інструментів потребують саме такого вигляду входу: текстовий файл, де кожен рядок є окремим прикладом для обробки. Крім того, обмеження кількості прикладів до 25000 дозволяє зберегти розмір файлу у межах допустимих лімітів для швидкого ітеративного тестування різних методів сегментації без надмірного навантаження на ресурси системи.

3.6 Реалізація процесу навчання токенізаторів

У файлі `train_tokenizers.py` реалізовано логіку створення та навчання трьох типів токенізаторів BPE, WordPiece та Unigram. Кожен із них представлений у вигляді окремої конфігурації з використанням відповідного модулю із бібліотеки `tokenizers`. Так можна гнучко працювати з кастомними моделями токенізації на низькому рівні, маючи при цьому зручні інтерфейси для тренування, нормалізації, попередньої токенізації та збереження моделей у форматі JSON.

На самому початку здійснюємо імпорт усіх необхідних модулів (лістинг 3.4)

Лістинг 3.4 – Імпорт необхідних модулів

```
from tokenizers import Tokenizer
from tokenizers.normalizers import NFD, Lowercase,
StripAccents, Sequence
from tokenizers.pre_tokenizers import Whitespace
from tokenizers.trainers import BpeTrainer,
WordPieceTrainer, UnigramTrainer
from tokenizers.models import BPE, WordPiece, Unigram
```

Перший рядок імпортує загальні модулі для побудови токенизаторів. Об'єкт `Tokenizer` є головним класом, який інкапсулює всі етапи обробки тексту. У другому рядку з модуля `normalizers` імпортується набір процедур для нормалізації тексту, а саме: приведення тексту до нормальної форми `NFD`, зменшення регістру, видалення акцентів. Ці компоненти комбінуються у єдину послідовність через клас `Sequence`. Далі імпортується клас `Whitespace`, що застосовується для попередньої токенизації на основі пробілів. Останній рядок імпортує специфічні моделі токенизаторів та відповідні класи для тренування кожного з них.

Наступним кроком оголошується функція `train_tokenizer`, яка отримує три вхідні параметри: тип моделі токенизатора, список файлів для тренування та бажаний розмір словника.

```
def train_tokenizer(model_type, files, vocab_size=8000):
    special_tokens = ["[PAD]", "[UNK]", "[CLS]", "[SEP]",
                     "[MASK]"].
```

Всередині функції оголошується список спеціальних токенів, які є обов'язковими для роботи моделей трансформерного типу. Вони включають маркери паддінгу, невідомого токена, початку та кінця послідовності, а також токен маскування. Ці символи будуть явно включені

до словника під час тренування, що гарантує їх наявність у кожній навченій моделі.

Подальші умовні блоки коду визначають конфігурацію для кожного типу токенизатора (лістинг 3.5).

Лістинг 3.5 – Визначення конфігурації для кожного типу токенизатора

```
if model_type == "bpe":
    model = BPE(unk_token="[UNK]")
    tokenizer = Tokenizer(model)
    trainer = BpeTrainer(vocab_size=vocab_size,
special_tokens=special_tokens)

elif model_type == "wordpiece":
    model = WordPiece(unk_token="[UNK]")
    tokenizer = Tokenizer(model)
    trainer = WordPieceTrainer(vocab_size=vocab_size,
special_tokens=special_tokens)

elif model_type == "unigram":
    model = Unigram()
    tokenizer = Tokenizer(model)
    trainer = UnigramTrainer(vocab_size=vocab_size,
special_tokens=special_tokens, unk_token="[UNK]")

else:
    raise ValueError("Unsupported model type")
```

Кожна гілка ініціалізує відповідну модель. Для BPE використовується клас BPE, що реалізує алгоритм парного злиття символів. Для WordPiece відповідно клас WordPiece, який працює на основі максимальної правдоподібності з урахуванням контексту. Unigram ініціалізується через клас Unigram, що є стохастичним методом побудови найкращого набору токенів шляхом перебору підмножин.

До кожної моделі прикріплюється об'єкт `Tokenizer`, що визначає повну пайплайн-структуру обробки тексту, а також відповідний `Trainer`, який відповідає за стратегію тренування. Всі тренери мають однакову ціль – досягти бажаного розміру словника за мінімальної втрати інформації та обов'язково додати спеціальні токени.

Наступним кроком виконується конфігурація нормалізаторів та попередніх токенизаторів:

```
tokenizer.normalizer = Sequence([NFD(), Lowercase(),  
StripAccents()])  
tokenizer.pre_tokenizer = Whitespace().
```

Перший рядок об'єднує послідовність нормалізацій у єдиний процес, який буде застосовуватись до кожного вхідного рядка тексту. Спочатку виконується розкладання символів Unicode у нормальну форму, потім текст перетворюється у нижній регістр, після чого видаляються діакритичні знаки та акценти. Цей процес забезпечує стандартизацію корпусу та зменшення варіативності в словах, що особливо важливо для моделей, які покладаються на статистичні властивості.

У другому рядку вказуємо, що попередня токенизація відбуватиметься шляхом розбиття тексту за пробілами. Це дозволяє кожному слову чи значущій одиниці бути обробленими окремо до моменту застосування конкретної токенизуючої стратегії.

Після встановлення всіх конфігурацій здійснюємо навчання токенизатора:

```
tokenizer.train(files, trainer).
```

Метод `train` приймає список текстових файлів і відповідного тренера. Під час цього процесу файл `train_text.txt`, сформований раніше, використовується як джерело навчального корпусу. Кожен рядок з цього файлу розглядається як окрема текстова одиниця. Після завершення тренування модель має повний словник, правила сегментації та інформацію про всі спеціальні токени.

Після тренування токенизатор зберігається у форматі JSON:

```
tokenizer.save(f"{model_type}_tokenizer.json").
```

Цей файл містить усі параметри моделі: розмір словника, список токенів, метод нормалізації, правила розбиття тексту, а також навчений алгоритм злиття або скорочення залежно від вибраного типу. У подальшому ці файли використовуються в основному процесі токенізації в скрипті `train.py`, де відбувається завантаження моделей за допомогою класу `PreTrainedTokenizerFast`.

Останнім блоком коду є умовний виклик основної функції:

```
if __name__ == "__main__":
    for model in ["bpe", "wordpiece", "unigram"]:
        train_tokenizer(model, ["train_text.txt"]).
```

У цьому блоці послідовно запускається навчання трьох токенізаторів, які зберігаються під назвами `bpe_tokenizer.json`, `wordpiece_tokenizer.json` та `unigram_tokenizer.json` (рисунок 3.7).

{} bpe_tokenizer.json > ...		{} wordpiece_tokenizer.json > {} n()		unigram_tokenizer.json > {} model > [] vocab > [
71	"model": {	71	"model": {	71	"model": {
80	"vocab": {	76	"vocab": {	74	"vocab": [
81	"[PAD]": 0,	77	"[PAD]": 0,	75	[
82	"[UNK]": 1,	78	"[UNK]": 1,	76	"[PAD]",
83	"[CLS]": 2,	79	"[CLS]": 2,	77	0.0
84	"[SEP]": 3,	80	"[SEP]": 3,	78],
85	"[MASK]": 4,	81	"[MASK]": 4,	79	[
86	"\b": 5,	82	"\b": 5,	80	"[UNK]",
87	"\u0010": 6,	83	"\u0010": 6,	81	0.0
88	" ": 7,	84	" ": 7,	82],
89	"\"": 8,	85	"\"": 8,	83	[
90	"#": 9,	86	"#": 9,	84	"[CLS]",
91	"\$": 10,	87	"\$": 10,	85	0.0
92	"%": 11,	88	"%": 11,	86],
93	"&": 12,	89	"&": 12,	87	[
94	"\'": 13,	90	"\'": 13,	88	"[SEP]",
95	"(": 14,	91	"(": 14,	89	0.0
96	")": 15,	92	")": 15,	90],
97	"*": 16,	93	"*": 16,	91	[
98	"+": 17,	94	"+": 17,	92	"[MASK]",
99	",": 18,	95	",": 18,	93	0.0
100	"-": 19,	96	"-": 19,	94],
101	".": 20,	97	".": 20,	95	[
102	"/": 21,	98	"/": 21,	96	"s",
103	"0": 22,	99	"0": 22,	97	-2.810934808226385
104	"1": 23,	100	"1": 23,	98],
105	"2": 24,	101	"2": 24,	99	[
106	"3": 25,	102	"3": 25,	100	"t",
107	"4": 26,	103	"4": 26,	101	-3.0826193210600206

Рисунок 3.7 – Фрагменти файлів `bpe_tokenizer.json`, `wordpiece_tokenizer.json` та `unigram_tokenizer.json` відповідно

Кожен із них базується на одному і тому самому корпусі, що гарантує рівні умови навчання та дозволяє безпосередньо порівнювати результати класифікації моделей на основі різних підходів до токенізації.

3.7 Навчання моделі

Файл `train.py` є ключовим елементом всієї експериментальної частини кваліфікаційної роботи. В ньому реалізовано основну логіку побудови, навчання та оцінки простої нейронної мережі класифікації текстів, які попередньо були токенізовані різними токенізаторами. У цьому файлі поєднуються можливості бібліотек `PyTorch`, `Transformers`, `datasets` і власне моделі, визначеної у `model.py`.

Спочатку здійснюємо імпорт необхідних бібліотек (лістинг 3.6).

Лістинг 3.6 – Визначення конфігурації для кожного типу токенізатора

```
import torch
import torch.nn as nn
from torch.utils.data import DataLoader
from transformers import PreTrainedTokenizerFast
from datasets import load_dataset
from model import Classifier
from tqdm import tqdm
import numpy as np
```

Цей блок відповідає за підключення засобів для побудови нейронної мережі, обробки датасетів, створення кастомного токенізатора через `PreTrainedTokenizerFast`, завантаження текстових корпусів за допомогою `HuggingFace datasets`, а також використання прогрес-бара з `tqdm`.

Наступною йде функція `tokenize_dataset`, яка є універсальним засобом токенізації набору даних для всіх типів токенізаторів (лістинг 3.7).

Лістинг 3.7 – Функція `tokenize_dataset`

```
def tokenize_dataset(dataset, tokenizer, max_length=128):
    def tokenize(example):
        return tokenizer(example['text'],
padding='max_length', truncation=True,
max_length=max_length)
    return dataset.map(tokenize, batched=True)
```

Ця функція приймає датасет, екземпляр токенизатора та максимальну довжину послідовності. Вона створює вкладену функцію `tokenize`, яка викликає токенизатор для кожного тексту, додаючи паддінг і обрізаючи текст до заданої довжини. Потім вона застосовується до всього датасету за допомогою методу `map` з параметром `batched=True`, що дозволяє обробляти приклади пакетно і таким чином підвищити ефективність.

Функція `collate_fn` використовується як функція пакування для `DataLoader`, що дозволяє перетворити список прикладів у батч тензорів (лістинг 3.8).

Лістинг 3.8 – Функція `collate_fn`

```
def collate_fn(batch):
    input_ids = torch.tensor([item['input_ids'] for item in
batch])
    labels = torch.tensor([item['label'] for item in
batch])
    return input_ids, labels
```

Ця функція бере окремі елементи батчу і з'єднує їх у один тензор `input_ids` та відповідні мітки `labels`. Вона дозволяє організовано працювати з токенизованими прикладами, які повертає `datasets`.

Функція `train_model` відповідає за одну «епоху» навчання моделі (лістинг 3.9).

Лістинг 3.9 – Функція `train_model`

```
def train_model(model, dataloader, optimizer, criterion,
device):
    model.train()
    total_loss = 0
    for input_ids, labels in tqdm(dataloader):
        input_ids, labels = input_ids.to(device),
labels.to(device)
        optimizer.zero_grad()
        outputs = model(input_ids)
        loss = criterion(outputs, labels)
        loss.backward()
        optimizer.step()
        total_loss += loss.item()
    return total_loss / len(dataloader)
```

На початку встановлюється режим навчання за допомогою `model.train()`. Далі по кожному батчу даних з `dataloader` виконується переміщення в обране обчислювальне середовище (CPU або GPU), обчислення передбачення `outputs` через прямий прохід моделі, обчислення функції втрат `loss`, зворотне поширення похибки та оновлення параметрів моделі. Змінна `total_loss` акумулює загальні втрати для подальшого виведення середнього значення. Оцінювання моделі реалізоване у функції `evaluate_model` (лістинг 3.10).

Лістинг 3.10 – Функція `evaluate_model`

```
def evaluate_model(model, dataloader, device):
    model.eval()
    correct = 0
    total = 0
    with torch.no_grad():
        for input_ids, labels in dataloader:
```

Продовження лістингу 3.10

```

        input_ids, labels = input_ids.to(device),
        labels.to(device)

        outputs = model(input_ids)
        predictions = torch.argmax(outputs, dim=1)
        correct += (predictions == labels).sum().item()
        total += labels.size(0)

    return correct / total

```

Ця функція встановлює модель у режим оцінки та забороняє обчислення градієнтів, що дозволяє зменшити споживання пам'яті. Після прямого проходу визначаються `predictions` як індекси з максимальними логітами. Точність обчислюється як співвідношення правильно класифікованих прикладів до загальної кількості.

Основна функція `main` (лістинг 3.11) виконує логіку завантаження токенизатора, підготовки даних, ініціалізації моделі, її навчання та тестування.

Лістинг 3.11 – Основна функція `main`

```

def main(model_type):
    tokenizer =
    PreTrainedTokenizerFast(tokenizer_file=f"{model_type}_toke
nizer.json",
        unk_token="[UNK]", pad_token="[PAD]",
        cls_token="[CLS]", sep_token="[SEP]")

```

На початку створюється токенизатор типу `PreTrainedTokenizerFast` для подальшого використання токенизаторів, які було навчено в `train_tokenizers.py`. Назва файлу токенизатора динамічно формується на основі параметра `model_type`.

Далі завантажуються набір даних `imdb` через `load_dataset`:

```
raw_dataset = load_dataset("imdb")
```

```

small_train =
raw_dataset["train"].shuffle(seed=42).select(range(2000))
small_test =
raw_dataset["test"].shuffle(seed=42).select(range(500)).

```

Було обрано підмножину з 2000 навчальних прикладів та 500 тестових, щоби прискорити експерименти та зменшити час навчання без значної втрати репрезентативності.

Наступним кроком токенізуються обидва піднабори:

```

tokenized_train = tokenize_dataset(small_train, tokenizer)
tokenized_test = tokenize_dataset(small_test, tokenizer).

```

Результати передаються в `DataLoader`, де визначено пакет розміром 32 та відповідну функцію пакування `collate_fn`:

```

train_loader = DataLoader(tokenized_train, batch_size=32,
shuffle=True, collate_fn=collate_fn)
test_loader = DataLoader(tokenized_test, batch_size=32,
shuffle=False, collate_fn=collate_fn).

```

Далі визначається обчислювальний пристрій та створюється модель:

```

device = torch.device("cuda" if torch.cuda.is_available()
else "cpu")
vocab_size = tokenizer.vocab_size
model = Classifier(vocab_size).to(device).

```

Модель `Classifier` описана в `model.py`, де є лише шар ембедінгів, середнє агрегування по токенах та двошаровий класифікатор. Цей підхід дозволяє протестувати вплив лише токенізації без складних архітектур.

Останнім етапом визначаються оптимізатор `Adam` та функція втрат `CrossEntropyLoss`:

```

optimizer = torch.optim.Adam(model.parameters(), lr=1e-3)
criterion = nn.CrossEntropyLoss().

```

Запуск навчання і тестування здійснюється протягом 15 епох наведено у лістингу 3.12.

Лістинг 3.12 – Запуск навчання і тестування протягом 15 епох

```

print(f"Training model with {model_type} tokenizer...")

```

Продовження лістингу 3.12

```
for epoch in range(15):
    loss = train_model(model, train_loader, optimizer,
criterion, device)
    acc = evaluate_model(model, test_loader, device)
    print(f"Epoch {epoch+1}: Loss = {loss:.4f}, Accuracy =
{acc*100:.2f}%")
```

На кожній епосі обчислюється середня втрата і точність на тестовому наборі, що дозволяє спостерігати динаміку навченості моделі для кожного типу токенизатора.

Фінальна частина виконується при безпосередньому запуску скрипта:

```
if __name__ == '__main__':
    for model_type in ["bpe", "wordpiece", "unigram"]:
        main(model_type).
```

Отже, послідовно відбувається тренування трьох моделей, кожна з яких використовує різний тип токенизатора. Це дозволяє зібрати порівняльну інформацію щодо продуктивності моделей в умовах однакової архітектури, гіперпараметрів та обсягів даних, змінюючи лише спосіб попередньої обробки тексту.

Вихідний код програмної реалізації наведено у додатку А.

3.8 Аналіз та інтерпретація результатів

У результаті практичної частини роботи було реалізовано повний цикл експериментів для порівняння ефективності трьох різних підходів до токенизації текстів у задачі класифікації тексту на основі набору даних IMDB. Після навчання кожного токенизатора були побудовані і навчені однакові за архітектурою класифікаційні моделі. Основною метою було дослідити, як тип токенизації впливає на продуктивність моделі на рівні метрик точності та функції втрат.

Навчання тривало протягом 15 епох для кожного токенизатора. На кожній ітерації фіксувалися значення функції втрат та точності на тестовому наборі даних. Нижче представлено узагальнену таблицю результатів, які були отримані під час експериментів (таблиця 3.1).

Таблиця 3.1 – Порівняння результатів

Epoch	BPE Loss	BPE Acc (%)	WordPiece Loss	WordPiece Acc (%)	Unigram Loss	Unigram Acc (%)
1	0.6887	58.20	0.6879	62.60	0.6843	56.20
2	0.6602	63.00	0.6531	64.80	0.6541	63.20
3	0.5904	65.40	0.5628	65.00	0.5925	67.20
4	0.4839	67.00	0.4423	69.40	0.5064	71.00
5	0.3737	71.20	0.3293	70.60	0.4220	71.40
6	0.2768	71.60	0.2422	70.80	0.3487	72.40
7	0.1979	72.00	0.1714	70.60	0.2761	72.60
8	0.1393	74.00	0.1149	72.00	0.2187	73.60
9	0.0919	71.00	0.0764	72.40	0.1694	73.60
10	0.0583	71.80	0.0485	72.40	0.1242	74.00
11	0.0358	72.80	0.0319	71.60	0.0900	73.40
12	0.0236	70.60	0.0214	72.40	0.0664	73.80
13	0.0166	72.80	0.0148	71.80	0.0499	73.00
14	0.0114	71.60	0.0107	71.80	0.0364	72.60
15	0.0083	71.60	0.0080	72.00	0.0273	74.80

З таблиці видно, що всі три токенизатори демонструють близькі результати, однак при детальнішому розгляді спостерігаються значущі відмінності у динаміці навчання та остаточній точності.

Починаючи з першої епохи, найкращу стартову точність показав токенизатор WordPiece. Його результат 62.60% на фоні 58.20% у BPE та 56.20% у Unigram свідчить про здатність WordPiece швидше формувати інформативні вхідні представлення тексту, ймовірно завдяки менш агресивному поділу токенів.

На другій епосі WordPiece зберігає лідерство з 64.80%, хоча Unigram стрімко вирівнюється і досягає 63.20%, залишивши BPE позаду. Починаючи з третьої епохи, Unigram виривається вперед і демонструє стабільне

зростання точності. У п'ятій епосі всі три токенизатори перевищують 70% точності, однак перевага поступово переходить до Unigram.

Максимальне значення точності, досягнуте протягом навчання, спостерігається на п'ятнадцятій епосі (рисунок 3.8). Саме тоді модель з токенизатором Unigram досягає 74.80% точності, що є найкращим показником серед усіх експериментів. WordPiece зупиняється на 72.00%, тоді як BPE відстає на рівні 71.60%.

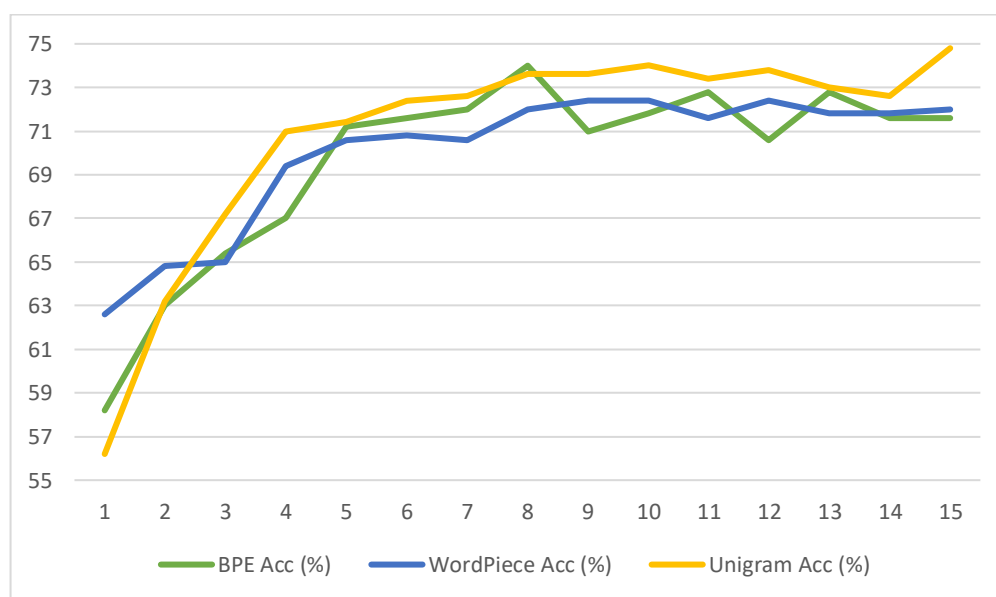


Рисунок 3.8 – Графік порівняння точності токенизаторів

Графік функції втрат також демонструє помітну різницю (рисунок 3.9). Unigram швидше і стабільніше зменшує втрати до кінця навчання. У п'ятнадцятій епосі він показує значення втрат 0.0273, що свідчить про точнішу відповідність моделі до міток класифікації. WordPiece закінчує з втратою 0.0080, а BPE — з 0.0083. Незважаючи на мінімальну різницю у втраті між останніми двома, точність класифікації виявляється вищою саме у моделі з Unigram.

Також слід зазначити характер кривих навченості. Точність моделі з BPE демонструє відносну нестабільність, з деяким зниженням після восьмої епохи. Це може свідчити про надмірну фрагментацію лексем у випадку BPE,

яка ускладнює навчання узагальнених шаблонів. У випадку WordPiece точність залишається стабільною, але майже не зростає після сьомої епохи, що може вказувати на досягнення локального максимуму. Unigram навпаки демонструє стабільне зростання та мінімальні коливання, що свідчить про кращу загальну узгодженість токенів з контекстами.

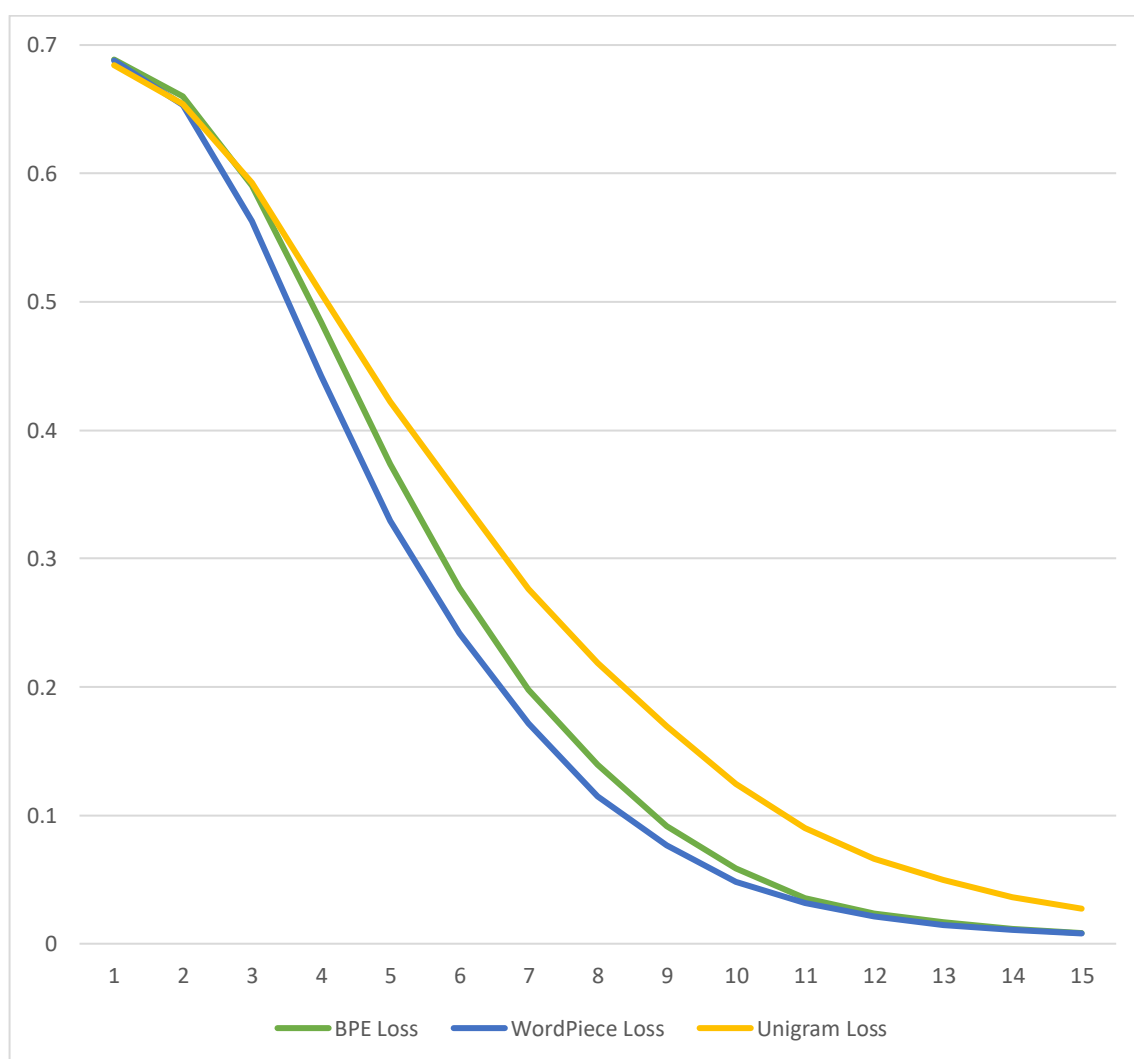


Рисунок 3.9 – Графік порівняння точності токенізаторів

Отже, аналіз результатів дає підстави зробити кілька важливих висновків. Найменш продуктивним виявився токенізатор BPE. Його повільний старт і менша точність навіть після 15 епох вказують на слабшу узгодженість із структурою англomовного тексту. WordPiece демонструє

гарні стартові результати, але швидко досягає плато. Найкращі результати показав токенізатор Unigram, який поєднує здатність до генералізації та ефективно відображення контекстуальних залежностей.

Можна зробити висновок, що у рамках задачі бінарної класифікації текстів за допомогою простої нейромережевої моделі, токенізація на основі моделі Unigram забезпечує найвищу точність та найкращу стабільність у навчанні. Це дозволяє рекомендувати її як оптимальний варіант для подібних завдань у системах автоматичного аналізу текстів. При ускладненні архітектури моделі або зміні задачі висновки можуть відрізнятись, однак в межах поточного дослідження Unigram продемонстрував найкращу ефективність.

ВИСНОВКИ

У межах виконання кваліфікаційної роботи було повністю реалізовано поставлене завдання дослідження ефективності різних підходів до токенізації в задачі автоматичної класифікації текстів на прикладі набору даних з рецензіями IMDb. Практична частина роботи охоплювала побудову трьох типів токенізаторів – BPE, WordPiece та Unigram – з нуля із використанням бібліотеки HuggingFace Tokenizers. Для кожного з варіантів токенізації було навчено окрему нейронну модель однакової архітектури, що дозволило провести коректне порівняльне оцінювання якості класифікації на основі обраних методів попередньої обробки текстових даних.

Досягнуті кількісні та якісні показники підтвердили обґрунтованість гіпотези про важливість вибору моделі токенізації при розв'язанні задач глибинного навчання на текстах. Найвищу точність класифікації було досягнуто за використання токенізатора Unigram, який продемонстрував результат 74.80% на тестовому наборі. Модель з токенізатором WordPiece досягла точності 72%, а модель із застосуванням BPE – 71.60%. Динаміка зниження функції втрат для кожної з моделей також вказує на перевагу підходу Unigram, який показав не лише стабільні результати, а й найнижче підсумкове значення втрат після повного циклу навчання.

Важливим аспектом є співвідношення отриманих результатів з наявними аналогами у світовій практиці. Більшість сучасних моделей трансформерного типу, як-от BERT, RoBERTa та ALBERT, активно застосовують WordPiece і BPE. Використання Unigram, хоч і менш поширене, зокрема в SentencePiece, з кожним роком набуває більшої популярності, оскільки дозволяє досягти кращого узгодження токенів зі змістом тексту завдяки гнучкій структурі. Результати цього дослідження узгоджуються з такими тенденціями і надають емпіричне підтвердження

ефективності Unigram на рівні прикладної класифікації з використанням стандартної англійської вибірки.

Наукова новизна роботи полягає у створенні повного експериментального середовища з нуля для вивчення впливу моделей токенизації на якість класифікації. В ході реалізації було створено незалежні модулі для генерації токенизаторів, підготовки датасету, побудови моделей та аналізу результатів. Всі компоненти реалізовано у вигляді коду з можливістю повторного використання, що важливо з огляду на відкритість та відтворюваність досліджень у сучасній науці. Подальший розвиток цього напрямку може бути пов'язаний з використанням складніших моделей на кшталт BERT або LSTM, навчанням на більших корпусах, а також дослідженням впливу токенизації в умовах мультимовного середовища.

Отже, робота виконана повністю відповідно до технічного завдання, поставлених цілей досягнуто та отримано вагомий науковий, так і прикладні результати. Підсумки дослідження підтверджують важливість ретельного вибору методу токенизації у проєктах з обробки тексту та відкривають перспективи для подальших досліджень у галузі NLP.

ПЕРЕЛІК ДЖЕРЕЛ ПОСИЛАННЯ

1. Mullen L. A., Others. Fast, consistent tokenization of natural language text. *Journal of open source software*. 2018. Т. 3, № 23. С. 655. URL: <https://doi.org/10.21105/joss.00655> (дата звернення: 16.05.2025).
2. V. S., J. R. Text mining: open source tokenization tools – an analysis. *Advanced computational intelligence: an international journal (ACII)*. 2016. Т. 3, № 1. С. 37–47. URL: <https://doi.org/10.5121/acii.2016.3104> (дата звернення: 16.05.2025).
3. Dutta S. K. Tokenization. *The definitive guide to blockchain for accounting and business: understanding the revolutionary technology*. 2020. С. 79–105. URL: <https://doi.org/10.1108/978-1-78973-865-020201006> (дата звернення: 16.05.2025).
4. Choo S., Kim W. A study on the evaluation of tokenizer performance in natural language processing. *Applied artificial intelligence*. 2023. Т. 37, № 1. URL: <https://doi.org/10.1080/08839514.2023.2175112>.
5. Publishing P. Natural language processing: python and NLTK. Packt Publishing, 2016. 702 с.
6. Arasanipalai A. U., Patel A. A. Applied natural language processing in the enterprise. O'Reilly Media, Incorporated, 2021.
7. Semantic tokenizer for enhanced natural language processing. *arXiv.org*. URL: <https://arxiv.org/abs/2304.12404> (дата звернення: 16.05.2025).
8. Impact of tokenization on language models: an analysis for turkish / C. Toraman та ін. *ACM transactions on asian and low-resource language information processing*. 2023. Т. 22, № 4. С. 1–21. URL: <https://doi.org/10.1145/3578707> (дата звернення: 30.05.2025).
9. Bostrom K., Durrett G. Byte pair encoding is suboptimal for language model pretraining. *Findings of the association for computational linguistics: EMNLP 2020*, м. Online. Stroudsburg, PA, USA, 2020. URL:

<https://doi.org/10.18653/v1/2020.findings-emnlp.414> (дата звернення: 19.06.2025).

10. Tokenizer choice for LLM training: negligible or crucial? / М. Ali та ін. *Findings of the association for computational linguistics: NAACL 2024*, м. Mexico City, Mexico. Stroudsburg, PA, USA, 2024. URL: <https://doi.org/10.18653/v1/2024.findings-naacl.247> (дата звернення: 31.05.2025).

11. Unpacking tokenization: evaluating text compression and its correlation with model performance / О. Goldman та ін. *Findings of the association for computational linguistics ACL 2024*, м. Bangkok, Thailand and virtual meeting. Stroudsburg, PA, USA, 2024. С. 2274–2286. URL: <https://doi.org/10.18653/v1/2024.findings-acl.134> (дата звернення: 01.06.2025).

12. Camacho-Collados J., Pilehvar M. T. On the role of text preprocessing in neural network architectures: an evaluation study on text categorization and sentiment analysis. *Proceedings of the 2018 EMNLP workshop blackboxnlp: analyzing and interpreting neural networks for NLP*, м. Brussels, Belgium. Stroudsburg, PA, USA, 2018. URL: <https://doi.org/10.18653/v1/w18-5406> (дата звернення: 02.06.2025).

13. Evaluating various tokenizers for arabic text classification / Z. Alyafeai та ін. *Neural processing letters*. 2022. URL: <https://doi.org/10.1007/s11063-022-10990-8> (дата звернення: 05.06.2025).

14. An empirical study of tokenization strategies for various korean NLP tasks / К. Park та ін. *Proceedings of the 1st conference of the asia-pacific chapter of the association for computational linguistics and the 10th international joint conference on natural language processing*, м. Suzhou, China. Stroudsburg, PA, USA, 2020. С. 133–142. URL: <https://doi.org/10.18653/v1/2020.aacl-main.17> (дата звернення: 05.06.2025).

15. Transformers. *Hugging Face – The AI community building the future*.
URL: <https://huggingface.co/docs/transformers/en/index> (дата звернення: 06.06.2025).

16. Tokenizers. *Hugging Face – The AI community building the future*.
URL: <https://huggingface.co/docs/tokenizers/index> (дата звернення: 06.06.2025).

17. An approach to the selection of behavior patterns autonomous intelligent mobile systems / O. Zolotukhin та ін. *Proceedings of the IEEE international conference on problems of infocommunications science and technology (PIC S&T)*. Kyiv, 2021. С. 349–352. URL: <https://doi.org/10.1109/PICST54195.2021.9772110> (дата звернення: 06.06.2025).

18. The methods for the prediction of climate control indicators in the Internet of Things systems / O. Zolotukhin та ін. *CEUR workshop proceedings*. 2021. URL: <https://doi.org/10.5281/zenodo.14526027> (дата звернення: 06.06.2025).

19. Methods of intellectual analysis of processes in medical information systems / V. Filatov та ін. *Information extraction and processing*. 2020. Т. 48, № 124. С. 92–98. URL: <https://doi.org/10.15407/vidbir2020.48.092> (дата звернення: 06.06.2025).

20. Filatov V., Semenets V., Zolotukhin O. Synthesis of semantic model of subject area at integration of relational databases. *Proceedings of the IEEE 8th international conference on advanced optoelectronics and lasers (CAOL)*. Sozopol, 2019. С. 598–601. URL: <https://doi.org/10.1109/CAOL46282.2019.9019532> (дата звернення: 06.06.2025).