

**ДЕРЖАВНИЙ УНІВЕРСИТЕТ  
ІНФОРМАЦІЙНО-КОМУНІКАЦІЙНИХ ТЕХНОЛОГІЙ  
НАВЧАЛЬНО-НАУКОВИЙ ІНСТИТУТ  
ІНФОРМАЦІЙНИХ ТЕХНОЛОГІЙ  
КАФЕДРА ІНЖЕНЕРІЇ ПРОГРАМНОГО ЗАБЕЗПЕЧЕННЯ**

**КВАЛІФІКАЦІЙНА РОБОТА**

на тему: «Покращення трансформерних методів автоматичної оцінки рівня володіння мовою шляхом введення адаптивного семантико-граматичного коефіцієнта»

на здобуття освітнього ступеня магістра  
зі спеціальності 121 Інженерія програмного забезпечення  
освітньо-професійної програми «Інженерія програмного забезпечення»

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

\_\_\_\_\_ Олександр МАТЮШКО  
(підпис)

Виконав: здобувач вищої освіти групи ПДМ-61  
Олександр МАТЮШКО

Керівник: \_\_\_\_\_ Тимур ДОВЖЕНКО  
канд. техн. наук

Рецензент: \_\_\_\_\_

**ДЕРЖАВНИЙ УНІВЕРСИТЕТ  
ІНФОРМАЦІЙНО-КОМУНІКАЦІЙНИХ ТЕХНОЛОГІЙ**

**Навчально-науковий інститут інформаційних технологій**

Кафедра Інженерії програмного забезпечення

Ступінь вищої освіти Магістр

Спеціальність 121 Інженерія програмного забезпечення

Освітньо-професійна програма «Інженерія програмного забезпечення»

**ЗАТВЕРДЖУЮ**

Завідувач кафедри

Інженерії програмного забезпечення

\_\_\_\_\_ Ірина ЗАМРІЙ

«\_\_\_\_\_» \_\_\_\_\_ 2025 р.

**ЗАВДАННЯ  
НА КВАЛІФІКАЦІЙНУ РОБОТУ**

Матюшко Олександрю Вадимовичу

1. Тема кваліфікаційної роботи: «Покращення трансформерних методів автоматичної оцінки рівня володіння мовою шляхом введення адаптивного семантико-граматичного коефіцієнта»

керівник кваліфікаційної роботи Тимур ДОВЖЕНКО, канд. техн. наук,

затверджені наказом Державного університету інформаційно-комунікаційних технологій від «30» жовтня 2025 р. № 467.

2. Строк подання кваліфікаційної роботи «19» грудня 2025 р.

3. Вихідні дані до кваліфікаційної роботи: науково-технічна література, опис моделей, що використовуються як база, датасети та мовні матеріали, вимоги до точності та функціональності оцінювальної делі, параметри для розробки адаптивного семантико-граматичного коефіцієнта, вимоги до експериментального дослідження.

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

1. Дослідження існуючих трансферних моделей.

2. Аналіз сучасних обмежень трансформерних методів оцінювання та потреба в їх покращенні.

3. Розробити та теоретично обґрунтувати новий метод оцінювання на основі адаптивного семантико-граматичного коефіцієнта (АСГК), який балансує вагу семантики та граматики залежно від CEFR-рівня і типу мовного завдання.

5. Перелік ілюстративного матеріалу: *презентація*

1. Актуальність розробки адаптивного семантико-граматичного коефіцієнта.
2. Існуючі моделі та їх неточність оцінювання.
3. Математична модель адаптивного семантико-граматичного коефіцієнта.
4. Алгоритм роботи адаптивного семантико-граматичного коефіцієнта.
5. Очікувані переваги розробленого методу.
6. Результати експериментів розробленої моделі адаптивного семантико-граматичного коефіцієнта.

6. Дата видачі завдання «31» жовтня 2025 р.

## КАЛЕНДАРНИЙ ПЛАН

№ з/п	Назва етапів кваліфікаційної роботи	Строк виконання етапів роботи	Примітка
1	Аналіз наявної науково-технічної літератури	31.10-04.11.2025	
2	Вивчення матеріалів для аналізу існуючих трансферних моделей аналізу мови	05.11-07.11.2025	
3	Дослідження методів оцінювання на основі трансферних моделей	08.11-12.11.2025	
4	Аналіз сучасних обмежень трансформерних методів оцінювання та потреба в їх покращенні	13.11-17.11.2025	
5	Дослідження концепції адаптивного семантико-граматичного коефіцієнта	18.11-21.11.2025	
6	Інтеграція адаптивного семантико-граматичного коефіцієнта у трансферну архітектуру	22.11-26.11.2025	
7	Оформлення роботи: вступ, висновки, реферат	27.11-05.12.2025	
8	Розробка демонстраційних матеріалів	06.12-15.12.2025	
9	Попередній захист роботи	16.12-19.12.2025	

Здобувач вищої освіти

\_\_\_\_\_ (підпис)

Олександр Матюшко

Керівник кваліфікаційної роботи

\_\_\_\_\_ (підпис)

Тимур ДОВЖЕНКО





## РЕФЕРАТ

Текстова частина кваліфікаційної роботи на здобуття освітнього ступеня магістра: 87 стор., 9 табл., 5 рис., 31 джерел.

*Мета роботи* – підвищення точності та об'єктивності систем автоматичного оцінювання рівня володіння мовою за рахунок інтеграції адаптивного семантико-граматичного коефіцієнта у трансформерні моделі.

*Об'єкт дослідження* – процес автоматичного оцінювання писемного та усного мовлення.

*Предмет дослідження* – методи та алгоритми трансформерних моделей у задачах оцінювання рівня володіння мовою.

У роботі використано широкий спектр методів, зокрема сучасні трансформерні підходи до обробки природної мови, статистичні моделі оцінювання та алгоритми глибинного навчання.

Проведено детальний аналіз актуальних методів автоматичної оцінки мовних навичок, що дозволило виявити їхні сильні сторони та ключові обмеження.

На основі отриманих результатів розроблено та оптимізовано новий підхід оцінювання з впровадженням адаптивного семантико-граматичного коефіцієнта, який забезпечує більш збалансоване та точне визначення рівня володіння мовою.

Для підтвердження ефективності запропонованого методу проведено комплексні експерименти, спрямовані на перевірку точності, стабільності, адаптивності та узгодженості результатів із людськими оцінками.

Результати дослідження підтверджують високу ефективність застосування багатомовних трансформерів у задачах автоматичного оцінювання мовних навичок, демонструючи їхню здатність узгоджуватися з людськими оцінками, адаптуватися до різних мовних контекстів і забезпечувати точну інтерпретацію як письмового, так і усного мовлення.

КЛЮЧОВІ СЛОВА: АВТОМАТИЧНЕ ОЦІНЮВАННЯ МОВЛЕННЯ,  
ТРАНСФОРМЕРНІ МОДЕЛІ, ГРАМАТИЧНА КОРЕКТНІСТЬ, СЕМАНТИЧНА  
ВІДПОВІДНІСТЬ, АДАПТИВНИЙ КОЕФІЦІЄНТ,  
СЕМАНТИКО-ГРАМАТИЧНИЙ АНАЛІЗ.

## ABSTRACT

Text part of the master's qualification work: 87 pages, 5 pictures, 9 table, 31 sources.

*The purpose of the work* is to improve the accuracy and objectivity of automatic language proficiency assessment systems by integrating an adaptive semantic-grammatical coefficient into transformer models.

*Object of research* – is the process of automatic assessment of written and spoken language.

*Subject of research* – is the methods and algorithms of transformer models in language proficiency assessment tasks.

The work uses a wide range of methods, including modern transformer approaches to natural language processing, statistical assessment models, and deep learning algorithms.

A detailed analysis of current methods of automatic language skill assessment was conducted, which allowed us to identify their strengths and key limitations.

Based on the results obtained, a new assessment approach was developed and optimized with the introduction of an adaptive semantic-grammatical coefficient, which provides a more balanced and accurate determination of language proficiency.

To confirm the effectiveness of the proposed method, comprehensive experiments were conducted to verify the accuracy, stability, adaptability, and consistency of the results with human assessments.

The results of the study confirm the high effectiveness of using multilingual transformers in the tasks of automatic language skill assessment, demonstrating their ability to agree with human assessments, adapt to different linguistic contexts, and provide accurate interpretation of both written and spoken language.

KEYWORDS: AUTOMATIC SPEECH EVALUATION, TRANSFORMER MODELS, GRAMMATICAL CORRECTNESS, SEMANTIC ACCURACY, ADAPTIVE COEFFICIENT, SEMANTIC-GRAMMATICAL ANALYSIS.

## ЗМІСТ

ПЕРЕЛІК УМОВНИХ ПОЗНАЧЕНЬ.....	28
ВСТУП.....	30
1 ОГЛЯД ЛІТЕРАТУРИ ТА АНАЛІЗ СУЧАСНИХ ПІДХОДІВ ДО ТРАНСФОРМЕРНИХ МЕТОДІВ АВТОМАТИЧНОЇ ОЦІНКИ РІВНЯ ВОЛОДІННЯ МОВОЮ.....	18
1.1 Автоматична оцінка рівня володіння мовою: сутність, завдання та еволюція підходів.....	18
1.2 Трансформерні моделі як основа сучасних систем автоматичної оцінки володіння мовою.....	20
1.3 Методи оцінювання письмового мовлення на основі трансформерних моделей.....	23
1.4 Методи оцінювання усного мовлення з використанням нейромережових моделей.....	26
1.5 Багатомовні та трансферні підходи в оцінюванні мовних навичок.....	28
1.6 Аналіз сучасних обмежень трансформерних методів оцінювання та потреба в їх покращенні.....	31
1.7 Узагальнення та формулювання напрямів подальшого покращення.....	34
2 РОЗРОБКА МЕТОДУ ПОКРАЩЕННЯ ТРАНСФОРМЕРНИХ СИСТЕМ АВТОМАТИЧНОЇ ОЦІНКИ РІВНЯ ВОЛОДІННЯ МОВОЮ ШЛЯХОМ ВВЕДЕННЯ АДАПТИВНОГО СЕМАНТИКО-ГРАМАТИЧНОГО КОЕФІЦІЄНТА.....	37
2.1 Постановка задачі та загальна концепція методу.....	37
2.2 Теоретичні основи оцінювання семантичної та граматичної правильності у трансформерних моделях.....	39
2.1.1 Формалізація семантичної складової Ssem.....	39
2.1.2 Формалізація граматичної складової Sgram.....	43
2.1.3 Кодування трансферами семантики та граматики.....	44
2.3 Концепція адаптивного семантико-граматичного коефіцієнта.....	47
2.4 Методика обчислення семантичної складової Ssem.....	49
2.5 Методика обчислення граматичної складової Sgram.....	50
2.6 Адаптивна функція, що поєднує семантичну та граматичну складові.....	52

2.7 Інтеграція методу у трансформерну архітектуру.....	58
2.7.1 Способи інтеграції коефіцієнта в архітектуру трансформера.....	58
2.7.2 Інтеграція як окремого модуля після отримання embeddings.....	59
2.7.3 Інтеграція як частина scoring head.....	59
2.7.4 Порівняння інтеграції у моделі-енкодери та моделі-декодери.....	60
2.8 Оцінка обчислювальної складності та переваг запропонованого методу.....	64
2.9 Висновки до другого розділу.....	66
<b>3 ЕКСПЕРИМЕНТАЛЬНЕ ДОСЛІДЖЕННЯ ТА ТЕСТУВАННЯ</b>	
<b>АДАПТИВНОГО СЕМАНТИКО-ГРАМАТИЧНОГО КОЕФІЦІЄНТА У</b>	
<b>ТРАНСФОРМЕРНИХ МОДЕЛЯХ ОЦІНЮВАННЯ МОВНИХ НАВИЧОК.....</b>	<b>68</b>
3.1 Мета та завдання експериментального дослідження.....	68
3.2 Вибір та опис датасетів для експерименту.....	69
3.3 Налаштування експериментальних моделей.....	71
3.4 Методика тестування та сценарії експериментів.....	73
3.4.1 Тестування на письмовому мовленні.....	73
3.4.2 Тестування на усному мовленні.....	74
3.4.3 Cross-lingual аналіз впливу рідної мови учня.....	74
3.5 Метрики оцінювання ефективності методу.....	75
3.5.1 Кореляційні метрики узгодженості з людськими оцінками.....	75
3.5.2 Метрики точності прогнозування.....	76
3.5.3 Метрики оцінювання модульних компонентів (Component-Level Metrics).....	76
3.5.4 Тести стійкості (Robustness Tests).....	77
3.6 Результати експериментів.....	78
3.6.1 Результати на письмових відповідях.....	78
3.6.2 Результати на усному мовленні.....	79
3.6.3 Порівняння базової та вдосконаленої моделей.....	80
3.7 Аналіз впливу адаптивного коефіцієнта та інтерпретація результатів.....	82
3.8 Висновки до третього розділу.....	85
<b>ВИСНОВКИ.....</b>	<b>87</b>
<b>ПЕРЕЛІК ПОСИЛАНЬ.....</b>	<b>89</b>

## ПЕРЕЛІК УМОВНИХ ПОЗНАЧЕНЬ

AESA — Автоматична оцінка мовних навичок (Automated Evaluation of Spoken/Written Ability).

BERT — Бідірекційні енкодерні представлення на основі трансформерів (Bidirectional Encoder Representations from Transformers)

RoBERTa — Оптимізований підхід BERT (Robustly Optimized BERT Approach)

XLM-R — Багатомовна модель RoBERTa (Cross-lingual Language Model – RoBERTa)

GPT — Генеративний попередньо навчений трансформер (Generative Pre-trained Transformer)

Pearson PTE — Тест англійської мови Pearson (Pearson Test of English)

TOEFL11 — Корпус TOEFL для 11 рідних мов (Test of English as a Foreign Language – 11 L1 Corpus)

ASAP — Премія автоматичного оцінювання учнівських робіт (Automated Student Assessment Prize)

AES — Автоматичне оцінювання есе (Automated Essay Scoring)

AWE — Автоматичне оцінювання письма (Automated Writing Evaluation)

WER — Коефіцієнт помилок у словах (Word Error Rate)

ASR — Система автоматичного розпізнавання мовлення (Automatic Speech Recognition system)

IELTS — Міжнародна система тестування англійської мови (International English Language Testing System)

CEFR — Загальноєвропейська шкала рівнів володіння мовою (Common European Framework of Reference for Languages)

LLM — Підходи на основі великих мовних моделей (Large Language Models)

CLS-вектор — Вектор спеціального токена CLS у трансформерах (Classification token vector)

NLI — Завдання природного мовного висновування (Natural Language Inference)

GECToR — Модель корекції граматичних помилок (Grammar Error Correction Transformer)

BLEURT — Модель оцінювання якості тексту на основі BERT (BERT-based Language Evaluation Understanding and Rating Tool)

LaBSE — Багатомовні ембеддинги речень (Language-Agnostic BERT Sentence Embedding)

mUSE — Багатомовний універсальний енкодер речень (multilingual Universal Sentence Encoder)

GEC — Корекція граматичних помилок (Grammar Error Correction)

PIE — Модель Prediction-Imputation-Editing для GEC (Prediction-Imputation-Editing)

RLHF — Навчання з підкріпленням на основі людського зворотного зв'язку (Reinforcement Learning from Human Feedback)

MAE — Середня абсолютна помилка (Mean Absolute Error)

## ВСТУП

У сучасних умовах глобалізації та цифрової трансформації суспільства зростає потреба у швидких, об'єктивних і масштабованих системах оцінювання рівня володіння іноземними мовами. Традиційні методи перевірки, що передбачають участь викладачів-експертів, забезпечують високу якість, однак залишаються трудомісткими, фінансово затратними та такими, що не можуть бути масштабовані для великих груп користувачів. Це особливо помітно у сферах дистанційного навчання, онлайн-тестування, мовних платформ та адаптивних освітніх систем.

За останні роки суттєвого розвитку набули методи автоматичного оцінювання мовних навичок (Automated Language Proficiency Assessment), зокрема ті, що базуються на трансформерних моделях: BERT, RoBERTa, XLM-R, GPT-подібні архітектури. Ці моделі здатні обробляти складні лінгвістичні структури, враховувати контекстні зв'язки та демонструють високий рівень узгодженості з людськими оцінками. Проте навіть сучасні трансформерні підходи мають низку обмежень: недостатня чутливість до граматичних помилок, переоцінювання текстів із простою лексикою, залежність від домену даних, а також відсутність механізмів адаптації до рівня підготовки здобувача. Це призводить до ситуацій, коли модель може правильно визначати загальну семантику висловлювання, але не враховувати структуру мови, або навпаки — надмірно штрафувати складні тексти.

Одним із перспективних шляхів покращення точності автоматичного оцінювання є розробка методу, який би інтегрував одночасно семантичні та граматичні характеристики відповіді та адаптував їхні ваги залежно від рівня володіння мовою, типу завдання та профілю помилок користувача. Такий підхід дозволяє наблизити модель до логіки роботи експертів-лінгвістів, які оцінюють не лише зміст, але й структурну правильність, відповідність нормам, складність використаних конструкцій та розвиток мовної компетентності.

У цій роботі запропоновано метод покращення трансформерних моделей для автоматичного оцінювання мовних навичок шляхом введення адаптивного семантико-граматичного коефіцієнта. Коефіцієнт формується на основі комбінування семантичної схожості, кількості та типів граматичних помилок, а також індивідуальних параметрів, що визначають рівень підготовки здобувача за шкалою CEFR. Така адаптивність дає змогу уникнути типових проблем «універсального» підходу та забезпечує більш збалансовану оцінку, що краще корелює з людськими судженнями.

Метою дослідження є підвищення точності та об'єктивності систем автоматичного оцінювання рівня володіння мовою за рахунок інтеграції адаптивного семантико-граматичного коефіцієнта у трансформерні моделі.

Для досягнення мети було сформульовано такі завдання:

- проаналізувати сучасні трансформерні підходи до автоматичного оцінювання мовних навичок;
- дослідити проблеми недостатньої чутливості моделей до граматики та контексту;
- розробити модель адаптивного коефіцієнта, що поєднує семантичні та граматичні характеристики;
- інтегрувати коефіцієнт у трансформерну архітектуру автоматичного оцінювання;
- провести експериментальне дослідження ефективності запропонованого методу на відкритих датасетах;
- виконати порівняльний аналіз результатів з базовими моделями та людськими оцінками.

Об'єктом дослідження є процес автоматичного оцінювання писемного та усного мовлення.

Предметом дослідження є методи та алгоритми трансформерних моделей у задачах оцінювання рівня володіння мовою.

Практична значущість роботи полягає у можливості застосування розробленого

методу в освітніх системах, онлайн-тестах, мовних платформах, а також у програмах підтримки вивчення іноземних мов.

Таким чином, робота спрямована на створення більш точного, адаптивного та інтерпретованого інструменту оцінювання, що враховує як змістову складову, так і граматичну правильність, забезпечуючи високий рівень відповідності оцінок людському експертному судженню та сприяючи розвитку сучасних інтелектуальних навчальних систем.

# 1 ОГЛЯД ЛІТЕРАТУРИ ТА АНАЛІЗ СУЧАСНИХ ПІДХОДІВ ДО ТРАНСФОРМЕРНИХ МЕТОДІВ АВТОМАТИЧНОЇ ОЦІНКИ РІВНЯ ВОЛОДІННЯ МОВОЮ

## 1.1 Автоматична оцінка рівня володіння мовою: сутність, завдання та еволюція підходів

Автоматична оцінка мовних навичок (AESA — Automated Evaluation of Spoken/Written Ability) є напрямом комп'ютерної лінгвістики та освітніх технологій, що передбачає використання алгоритмів штучного інтелекту для визначення рівня володіння мовою без участі людського експерта. Основною метою AESA є забезпечення швидкої, масштабованої та об'єктивної оцінки продуктивних мовних навичок, зокрема письма та говоріння, які традиційно вважаються найбільш складними для автоматизації.

Системи автоматичного оцінювання аналізують комплекс характеристик мовлення або письмового тексту. До основних аспектів, які підлягають вимірюванню, належать:

- граматика — коректність морфологічних і синтаксичних структур;
- лексика — різноманітність, точність та відповідність словникового запасу;
- зв'язність (coherence & cohesion) — логічність викладу, коректність використання зв'язок між реченнями;
- змістовність (content relevance) — відповідність відповіді поставленому завданню;
- вимова — чіткість артикуляції, коректність звуків;
- просодія — інтонація, ритм, паузування;
- флюентність — темп і плавність мовлення.

Розробка систем AESA має тривалу історію, і перші підходи ґрунтувалися на ручному конструюванні ознак та експертних правилах. Rule-based системи

фокусувалися на виявленні помилок через набори лінгвістичних правил (наприклад, граматичні шаблони, списки дозволеної лексики чи моделей речень). У поєднанні з ними використовувалися експертні системи, що наслідували оцінювальні стратегії викладачів за допомогою формалізованих «дерев рішень». Пізніше з'явилися статистичні моделі, у яких оцінка формувалася на основі окремих метрик — довжини речення, частоти помилок, лексичного різноманіття, n-грамних моделей тощо.

Проте традиційні підходи мають низку важливих обмежень:

- вони потребують ручного створення правил, що є трудомістким і погано масштабується;
- правила часто не охоплюють повного лінгвістичного різноманіття, особливо у текстах різних рівнів складності;
- такі системи мають низьку здатність до узагальнення — вони ефективні лише у межах тих шаблонів, для яких були створені;
- rule-based методи не здатні враховувати повноцінний контекст і взаємодію між словами.

Обмеження цих систем зумовили перехід до методів машинного навчання. Поява глибоких нейронних мереж та особливо моделей обробки природної мови заснованих на deep learning докорінно змінила підходи до AESA. Нейронні моделі дозволили відмовитися від ручної інженерії ознак і замінити її автоматично сформованими семантичними та синтаксичними уявленнями про текст або мовлення. Це підвищило гнучкість, точність і здатність до роботи з багатомовними корпусами.

Найбільш значущим кроком у розвитку AESA стала поява трансформерних моделей, що продемонстрували здатність ефективно працювати з довгими текстами, складною граматиною та багатомовними даними. Архітектури на основі трансформера — такі як BERT, RoBERTa, XLM-R, GPT, а для мовлення — Wav2Vec 2.0 — сьогодні є домінуючими в автоматичних оцінювальних системах. Трансформери поєднують контекстну

семантичну інформацію, виявляють глибинні граматичні залежності та дозволяють здійснювати оцінку, максимально наближену до людської.

Сучасні мовні тестові платформи вже інтегрують ці моделі у власні системи оцінювання. Наприклад:

- Duolingo English Test (DET) використовує нейронні моделі для автоматичного оцінювання продуктивних навичок та побудови адаптивної тестової шкали;
- Pearson PTE застосовує глибокі моделі для аналізу вимови, флюентності та граматичної точності;
- дослідницькі підрозділи ETS (організатор TOEFL) працюють над інтеграцією трансформерів для автоматичного виокремлення ознак письма та мовлення.

Поширення трансформерних моделей у тестології демонструє загальносвітовий тренд: від фіксованих, жорстких алгоритмічних систем до адаптивних, семантично чутливих та багатомовних нейронних підходів.

## **1.2 Трансформерні моделі як основа сучасних систем автоматичної оцінки володіння мовою**

Архітектура трансформерних нейронних мереж стала фундаментом більшості сучасних систем автоматичної оцінки мовних навичок завдяки високій здатності моделювати складні семантичні та граматичні залежності у мовленні. Трансформери забезпечують глибоке контекстуальне представлення тексту та аудіо, що робить їх придатними для задач оцінювання письма, усного мовлення, когерентності, зв'язності та орфографічно-граматичної правильності.

Трансформерна модель базується на механізмах self-attention та multi-head attention, які дозволяють моделі «розглядати» кожне слово у контексті всього речення або висловлювання. Це дає змогу точніше виявляти семантичні та синтаксичні залежності, ніж попередні рекурентні або згорткові архітектури.

Основні компоненти:

**Multi-head attention.** Дозволяє моделі паралельно аналізувати кілька типів зв'язків між токенами (синтаксичні, семантичні, позиційні), що суттєво підвищує якість оцінювання логічної структури та змістовності тексту.

**Positional encoding.** Оскільки трансформер не має рекурентного механізму, позиційне кодування вводить інформацію про порядок tokenів, що критично важливо для граматичних структур.

**Encoder/decoder stack.** У задачах оцінювання зазвичай використовують тип Encoder-only (BERT-подібні моделі) або Decoder-only (GPT-подібні моделі), тоді як повний encoder–decoder застосовується рідше.

Завдяки цим компонентам трансформери забезпечують гнучкість, масштабованість та здатність працювати з великими корпусами даних без втрати якості.

Трансформери стали галузевим стандартом з кількох причин:

1. **Контекстуальні векторні представлення (contextual embeddings).** Моделі генерують вектори, що враховують контекст використання слова, що дозволяє оцінювати не лише граматику, а й стилістику, логічні зв'язки та зміст.
2. **Масштабованість.** Моделі легко донавчати на спеціалізованих датасетах тестових відповідей (наприклад, TOEFL11, ASAP).
3. **Узагальнення.** Трансформери демонструють високу здатність узагальнювати мовні шаблони та обробляти відповіді, що відрізняються структурою, рівнем помилок та стилем викладу.

Ці властивості роблять їх базовою архітектурою у системах AES (Automated Essay Scoring), AWE (Automated Writing Evaluation) та AWE-Speech.

Огляд представників трансформерних моделей у задачах оцінювання BERT та RoBERTa — аналіз письма, граматики та когерентності

BERT-подібні моделі працюють у режимі *encoder-only* та добре підходять для задач класифікації, регресії та визначення якості тексту. Їх застосовують для: оцінювання граматичної правильності, аналізу когерентності та зв'язності тексту, визначення рівня складності письма, оцінювання відповідності змісту заданій темі.

RoBERTa, як оптимізована версія BERT, показує підвищену точність у задачах Automated Essay Scoring.

GPT-моделі (GPT-3/4/5) — генеративне оцінювання та *reasoning-based scoring*

GPT-типу — це *decoder-only* трансформери, здатні формувати розгорнуті експертні оцінки з поясненнями. Вони забезпечують:

- **генеративне оцінювання**, коли модель не лише виставляє бал, а й створює пояснення рішення;
- **reasoning-based scoring**, що враховує логіку, аргументацію, структуру відповідей;
- здатність моделювати людський процес оцінювання.

GPT використовується в сучасних комерційних AWE-системах, зокрема в Duolingo English Test, TOEFL Practice та LanguageTool AI Scoring.

XLM-R та mBERT моделі які працюють з десятками мов одночасно та вміють переносити знання між мовами:

- добре справляються з оцінюванням відповідей англійською від носіїв інших мов;
- зберігають стійкість до граматичних помилок, притаманних різним L1;
- дозволяють будувати універсальні багатомовні системи.

Для аналізу аудіо-тестів використовують трансформери, треновані на мовних даних:

**Wav2Vec2** — розпізнавання мовлення, оцінювання фонетичних помилок, темпу та плавності.

**Whisper** — багатомовне ASR із високою точністю, підходить для визначення акценту, вимови та частоти дисфлюенцій.

Вони є основою систем AWE-Speech, що оцінюють усні відповіді.

Сучасні системи автоматичної оцінки володіння мовою використовують трансформери для:

**AES (Автоматичне оцінювання есе):** аналіз структури, семантики, логічної зв'язності.

**AWE (Автоматичне оцінювання письма):** виявлення граматичних, лексичних та стилістичних помилок.

**AWE-Speech:** оцінювання усного мовлення, вимови, швидкості, інтонації.

Трансформерні моделі забезпечують новий рівень точності й узгодженості з експертними оцінками, що робить їх провідною технологією в автоматизованому оцінюванні мовних навичок.

### 1.3 Методи оцінювання письмового мовлення на основі трансформерних моделей

Сучасні трансформерні моделі стали ключовою технологією у задачах автоматичного оцінювання письмового мовлення завдяки здатності глибоко аналізувати структуру тексту, граматичну правильність, семантичну повноту та логічну когерентність. На відміну від традиційних статистичних або лінгвістичних підходів, трансформери працюють із контекстуальними представленнями, що дозволяє моделювати складні дискурсивні залежності між реченнями та абзацами.

Трансформери дозволяють виконувати комплексний аналіз на всіх рівнях тексту — від локальних граматичних конструкцій до глобальної логіки викладу. Основні аспекти аналізу включають:

**Coherence (когерентність)** — оцінювання логічної послідовності ідей, узгодженості між абзацами та правильності переходів між думками.

Self-attention механізми дозволяють моделі “бачити” глобальну структуру тексту.

**Cohesion (зв’язність)** — визначення використання дискурсивних маркерів, референцій, заміників, що забезпечують плавність викладу.

**Grammar & Syntax Errors (граматичні та синтаксичні помилки)** — трансформери автоматично ідентифікують помилки у структурі речень, узгодженні часів, порядку слів та пунктуації, використовуючи контекст.

Завдяки цьому трансформерні моделі забезпечують комплексну оцінку письмового мовлення, наближену за точністю до експертної.

Для оцінювання письмового тексту застосовуються різні типи представлень, отримані з трансформерних моделей:

**Contextual embeddings** — контекстуальні вектори слів та фраз, що враховують значення слова у конкретній позиції. Вони дозволяють оцінювати семантичну точність, лексичне різноманіття та стилістичні особливості.

**Sentence embeddings** — вектори речень, що використовуються для аналізу когерентності та логічної структури тексту.

**Discourse markers** — автоматичне виявлення дискурсивних індикаторів (“however”, “therefore”, “in addition”), що є важливою ознакою зв’язності й академічного рівня письма.

Комбінація цих представлень утворює багатовимірний простір ознак, що забезпечує ефективне оцінювання письма. Сучасні системи AES активно використовують трансформерні моделі у різних конфігураціях. Найпоширеніші підходи включають:

**BERTScore** — порівняння тексту студента з референсними відповідями через подібність BERT-ембеддингів; застосовується для оцінки семантичної близькості.

**BLEURT** — спеціально донаведена модель оцінки якості тексту, чутлива до семантичних помилок, стилістичних відхилень та слабкої аргументації.

**OpenAI Scoring Research** — сучасні генеративні моделі GPT-типу, треновані на експертних оцінках, забезпечують “людиноподібну” оцінку, здатну пояснити помилки та виявляти слабкі місця в логіці.

Ці методи часто комбінуються у гібридних системах, що поєднують регресію, класифікацію та генеративне оцінювання. Попри значні успіхи, існуючі системи оцінювання мають низку суттєвих обмежень, а саме:

1. **надмірна залежність від поверхневих патернів.** Моделі можуть помилково підвищувати оцінки за рахунок великої довжини тексту, складної лексики або частих дискурсивних маркерів;
2. **вразливість до “машинного” тексту.** Автоматично згенеровані відповіді можуть отримувати високі бали, хоча вони не відображають реальний рівень володіння мовою;
3. **проблеми з балансом між граматиною та семантикою.** Моделі часто добре оцінюють стилістику, але слабо визначають глибину змісту або навпаки — переоцінюють семантику при ігноруванні граматичних помилок;
4. **низька інтерпретованість.** Пояснення рішень є нечітким, що ускладнює використання систем у формальному оцінюванні (наприклад, у стандартизованих тестах).

Ці обмеження створюють потребу у нових методах, здатних забезпечити більш збалансовану та прозору оцінку письма. Незважаючи на потужність сучасних моделей, більшість підходів оцінювання не мають цілісного механізму, що одночасно й адаптивно враховує:

- граматичну правильність,
- дискурсивну зв'язність,
- семантичну точність,
- відповідність завданню.

Саме тому виникає необхідність у введенні адаптивного семантико-граматичного коефіцієнта, який би комбінував ці фактори у єдиній метриці. Такий коефіцієнт може стати ефективним способом збалансувати різні

рівні аналізу, мінімізувати вплив поверхневих ознак та забезпечити більш об'єктивну оцінку письмового мовлення.

#### **1.4 Методи оцінювання усного мовлення з використанням нейромережових моделей**

Автоматичне оцінювання усного мовлення (Automated Spoken Language Assessment) є однією з найбільш складних підсфер автоматичного оцінювання володіння мовою. Усне мовлення поєднує акустичні, фонетичні, просодичні та семантичні характеристики, тому традиційні rule-based системи демонстрували обмежену ефективність. Нейромережові моделі — особливо трансформерні — дозволили перейти від ручного аналізу ознак до енд-ту-енд підходів, де система навчається безпосередньо з аудіосигналу.

Ключове місце у сучасних системах оцінювання займають моделі ASR (Automatic Speech Recognition). Вони перетворюють аудіо на текст або ж витягують ознаки безпосередньо з сигналу. Серед найпоширеніших архітектур:

**Wav2Vec 2.0** — використовує self-supervised навчання на великих аудіокорпусах. Модель навчається виявляти приховану структуру мовленнєвого сигналу, після чого донавчається для конкретних задач. Переваги: здатність працювати з шумними даними, збереження інформації про вимову, інтонацію, темп.

**Whisper (OpenAI)** — багатомовна енд-ту-енд ASR-модель, стійка до шуму та акцентів. Whisper дозволяє не тільки транскрибувати мовлення, а й аналізувати паузи, просодію, частотні характеристики та дисфлюенції (запинки, повтори).

**HuBERT (Hidden-Unit BERT)** — поєднує self-supervised підхід з кластеризацією акустичних сигналів. HuBERT дозволяє виділяти узагальнені фонетичні одиниці, що є корисним при оцінюванні вимови та акценту.

Ці моделі стали основою систем AWE-Speech, здатних виконувати оцінювання енд-ту-енд без необхідності ручного створення акустичних чи фонетичних ознак.

Для оцінювання якості усного мовлення використовуються як традиційні метрики, так і нейромережеві представлення.

1. Phoneme Distance (дистанція між фонемами) порівнює послідовність фонем, отриманих від учня, з еталонною транскрипцією. Застосовується для вимірювання точності вимови, але залежить від того, наскільки точно ASR система трансформує сигнал у фонему.

2. Word Error Rate (WER) класична метрика для оцінювання ASR, яка вимірює розбіжність між розпізнаним текстом та еталоном. Недолік: WER не враховує якість вимови — тільки помилки у транскрипції.

3. Prosodic Features (просодичні ознаки) містять інтонаційні контури, ритм мовлення, частоту дисфлюенцій, темп, середню та варіативність висоти тону (F0), довжину пауз. Просодичні ознаки є особливо цінними, оскільки відображають природність мовлення та комунікативну компетентність.

Нейромережеві моделі (Whisper, HuBERT) дозволяють інтегрувати ці ознаки у кінцеве представлення без ручного витягування. Попри суттєвий прогрес, існуючі системи стикаються з низкою суттєвих викликів:

1. Акцент-варіативність. Сучасні моделі часто дають нижчу оцінку мовцям із сильним акцентом, навіть якщо граматичний та семантичний рівень високий. Це спричиняє упередженість до носіїв різних мов.

2. Чутливість до шуму та технічних артефактів. Навіть найкращі моделі можуть помилятися під час роботи з:

- телефонними записами;
- фоновим шумом;
- неякісними мікрофонами.

Це впливає не лише на транскрипцію, а й на оцінку вимови.

3. Слабка кореляція між вимовою та загальною мовною компетентністю. Висока якість вимови  $\neq$  високий рівень володіння мовою. Деякі учні мають добру граматику та словниковий запас, але сильний акцент. Чинні моделі часто занижують оцінку таких мовців.

4. Відсутність уніфікованої метрики для комбінованого оцінювання (мовлення + зміст). Більшість систем аналізують: або акустичні ознаки (вимову, темп), або семантичні характеристики (зміст, логічність, відповідність завданню).

Проте у реальних тестах (IELTS, TOEFL Speaking) оцінка ґрунтується на комбінації цих параметрів. У сучасних AI-системах бракує інтегрованої метрики, яка поєднувала б семантику, граматику та акустику в єдину формулу. Вказані обмеження свідчать про необхідність введення нових механізмів оцінювання, які:

- адаптивно враховують як семантичні, так і граматичні, і акустичні характеристики;
- мають стійкість до акцентів та шуму;
- відокремлюють змістову якість відповіді від артикуляційних особливостей;
- забезпечують об'єктивність та інтерпретованість.

Саме ці недоліки створюють підґрунтя для розробки адаптивного семантико-граматичного коефіцієнта, який покликаний інтегрувати різнорівневий аналіз у єдиній метриці та підвищити точність трансформерних моделей оцінювання мовної компетентності.

### **1.5 Багатомовні та трансферні підходи в оцінюванні мовних навичок**

У глобальному освітньому середовищі автоматичні системи оцінювання мовних навичок повинні бути здатні працювати не лише з англійською, але й з десятками інших мов. Багатомовні моделі та механізми перенесення знань

(transfer learning) відіграють ключову роль у створенні систем, здатних адекватно оцінювати відповіді учнів з різних лінгвістичних та культурних середовищ. Вони дозволяють мінімізувати потребу у великих корпусах для кожної окремої мови й дають змогу застосовувати універсальні лінгвістичні патерни, вивчені моделлю.

Cross-lingual transfer (крос-мовний перенос) — це здатність моделі, навченої на даних однієї мови (або групи мов), ефективно працювати на іншій мові завдяки спільним універсальним рисам мов. Основні принципи:

**Спільний латентний простір.** Багатомовні трансформери (XLM-R, mBERT, LaBSE) навчаються створювати представлення текстів різними мовами у єдиному векторному просторі. Це дозволяє моделі «бачити» структуру тексту незалежно від конкретної мови.

**Універсальні дискурсивні ознаки.** Зв'язність, логічна структура, послідовність думок — це риси, притаманні будь-якому письму, що дозволяє переносити навички оцінювання між мовами.

**Захоплення спільних семантичних патернів.** Багато концептів виражаються подібним способом у різних мовах, що робить можливим zero-shot або few-shot оцінювання, коли модель може оцінювати мову, на якій не навчалася.

**Урахування L1-впливу на L2.** У процесі навчання моделі починають розпізнавати типові помилки для носіїв конкретної L1 (наприклад, типові помилки україномовних чи китайськомовних студентів), що також дає змогу точніше оцінювати письмове та усне мовлення. Таким чином, крос-мовний перенос дозволяє створювати універсальні системи оцінювання, здатні працювати з багатьма мовами без необхідності дорогого збору великих датасетів для кожної L2.

Попри високий рівень універсальності багатомовних підходів, сучасні дослідження підкреслюють наявність **linguistic bias** (лінгвістичного упередження) — ситуації, коли модель оцінює учнів несправедливо в залежності від їх мовного фону, зокрема:

помилки дослівного перекладу або прямої кальки з рідної мови можуть завищувати або занижувати оцінку;

учні з «неспецифічними» або рідкісними мовами L1 можуть отримувати нижчі бали;

системи оцінювання усного мовлення часто штрафують за акцент, який фактично не обов'язково відображає комунікативні навички;

модель може бути упередженою щодо мов із нестандартним порядком слів або гнучкою морфологією.

Linguistic bias є критичною проблемою, оскільки автоматичні системи можуть бути використані для високостайкових оцінювань (IELTS, PTE, Duolingo English Test), де справедливість має фундаментальне значення. Для корекції лінгвістичного упередження застосовують низку сучасних підходів, оснований на нейромережевій оптимізації:

1. Adversarial learning (змагальне навчання) — модель навчається таким чином, щоб її векторні представлення не містили ознак, які дозволяють визначити рідну мову студента. В іншому модулі (adversary) система намагається розпізнати L1 за ембеддингами, тоді як основна модель намагається зробити це неможливим. Це створює «L1-інваріантне» представлення.

2. Domain adaptation (адаптація до домену) — модель додатково тренують на корпусах студентських робіт, які утворюють спеціальний навчальний домен — не носії мови з типовими помилками. Таким чином зменшується відрив між розподілом даних «модель бачила» і «модель повинна оцінити».

3. Data balancing та synthetic augmentation — штучно збалансовують датасет для різних мов створюють штучні помилки, типові для певних L1, додають синтетичні приклади з різними акцентами (наприклад, через голосові конверсійні моделі).

4. Multi-task learning — модель одночасно навчається на кількох задачах:

- визначення рівня CEFR,
- аналіз граматики,

- виявлення дискурсивних помилок.

Це зменшує шанс, що модель надмірно фокусується на поверхневих ознаках, пов'язаних із L1.

Такі підходи істотно знижують упередженість, але повністю проблему не усувають. Попри значні успіхи, багатомовні трансформери мають низку слабких місць:

**Падіння семантичної точності на рідкісних або малоресурсних мовах.**

XLM-R і mBERT добре працюють з великими європейськими мовами, але для рідкісних мов (баскська, телугу, грузинська) якість семантичного аналізу значно знижується.

**Висока чутливість до нестандартних акцентів у усному мовленні.**

Whisper або Wav2Vec2 можуть суттєво знижувати якість транскрипції під впливом «рідкісних» акцентів, що автоматично погіршує оцінку.

**Проблема семантичної «плутанини».** Багатомовні моделі інколи змішують значення слів між мовами (особливо спорідненими), що погіршує якість аналізу змісту.

**Недостатній обсяг навчальних даних для побудови універсальної шкали якості письма.** Якісна оцінка письма та мовлення вимагає великих корпусів оцінених робіт для кожної мови, яких зазвичай немає.

Ці обмеження створюють чіткий простір для інновацій — зокрема, впровадження адаптивних коефіцієнтів, що враховують як універсальні, так і мовно-специфічні характеристики текстів.

## **1.6 Аналіз сучасних обмежень трансформерних методів оцінювання та потреба в їх покращенні**

Попри вражаючі результати трансформерних архітектур у задачах обробки природної мови, їх застосування в автоматичному оцінюванні мовних навичок демонструє низку фундаментальних обмежень. Ці недоліки прямо формують наукову проблему, яку вирішує дана кваліфікаційна робота. Моделі

нового покоління добре працюють із високорівневою семантикою тексту, але значно гірше відображають граматичну якість, структурну організацію речень та рівневу специфіку мовної компетентності користувача. Лише інтегральне врахування семантики, граматики та мовного рівня може забезпечити коректне й справедливе оцінювання.

Сучасні трансформерні моделі (BERT, RoBERTa, DeBERTa, GPT-похідні) будують контекстні представлення, оптимізовані насамперед для *семантичної* подібності. Це означає:

- модель точно виявляє зміст висловлювання, тематичну релевантність, логічний зв'язок між частинами тексту;
- проте вона менш чутлива до локальних граматичних помилок, таких як неправильні відмінки, узгодження, синтаксичні порушення або морфологічні неточності;
- помилки типу *he go to school, much informations, I am agree* можуть занижувати якість мовлення, але трансформер часто класифікує такі фрази як «семантично зрозумілі», не штрафуючи їх належною мірою.

Таким чином, існуючі моделі «бачать зміст», але не «бачать форму», що є критичним недоліком у системах мовного оцінювання. Поширені сучасні метрики також підтверджують структурну проблему — вони аналізують текст переважно як семантичну матрицю, а не як граматично організовану структуру:

**BLEURT** зосереджений на семантичній відповідності референсу та кандидата; граматичні відхилення впливають на результат дуже слабо.

**BERTScore** оцінює подібність ембеддингів токенів, не враховуючи синтаксису чи граматичної коректності.

**G-Eval** (на основі LLM) використовує оцінювання великими моделями, але воно може бути нестабільним і надмірно залежним від узагальнених шаблонів, не маючи чіткої формалізації граматичних критеріїв.

Спільною рисою цих підходів є відсутність інтегральної оцінки семантики + граматики + структури речення. Жодна з метрик не дозволяє

одночасно оперувати глобальною семантикою й локальною граматичною послідовністю. Різні рівні володіння мовою (A1, A2, B1, B2, C1, C2) передбачають абсолютно різні вимоги щодо кількості граматичних конструкцій, складності лексики, точності синтаксису, глибини семантичного викладу.

Наприклад:

- на рівні **A1** семантична простота та базова зрозумілість важливіші за складну граматику;
- на рівні **C1** навпаки — граматичні та стилістичні недоліки значно знижують якість тексту.

Однак сучасні метрики застосовують однакові ваги семантики та граматики для всіх користувачів не враховують реальні вимоги CEFR до конкретного рівня та не можуть адаптувати важливість граматичної точності залежно від компетентності учня. Таким чином, оцінювання A1 та C1 відбувається «однією лінійкою», що робить результати менш об'єктивними.

Хоча трансформерні моделі досягають високих метрик якості, більшість із них залишаються чорними скриньками. Особливо це помітно у мовному оцінюванні:

- оцінка подається у вигляді одного числа;
- неможливо визначити, яка частка цієї оцінки відноситься до граматики, а яка — до семантики;
- учень не може зрозуміти, що саме потрібно покращити — структуру, лексику чи логіку викладу;
- викладач не може використовувати цю оцінку як аналітичний інструмент.

Це особливо важливо в освітніх системах, де потрібні не лише бали, а й пояснення. Огляд наукової літератури показує, що більшість сучасних рішень базуються на статичних метриках, які однаково застосовуються для будь-яких текстів, відсутні підходи, які б адаптивно змінювали ваги оцінювання залежно від мовного рівня або типу завдання. Не запропоновано моделей, здатних поєднати граматику та семантику у формалізований, керований коефіцієнтами

механізм. Немає метрики, яка б дозволяла кількісно оцінити внесок кожного компонента (grammar contribution, semantic contribution).

Саме ця прогалина формує головну дослідницьку проблему: створення адаптивної, інтерпретованої та багатокомпонентної метрики для оцінювання мовних навичок, здатної враховувати як рівень користувача, так і комплексність граматичних та семантичних характеристик тексту.

Основним напрямом покращення трансформерних методів є введення адаптивного семантико-граматичного коефіцієнта (АСГК) — нового підходу, який дозволяє об'єднати семантичні та граматичні характеристики тексту в єдину, контрольовану й інтерпретовану метрику.

### **1.7 Узагальнення та формулювання напрямів подальшого покращення**

Аналіз сучасних трансформерних моделей та існуючих методів автоматичного оцінювання мовних навичок демонструє наявність системних обмежень, які перешкоджають створенню універсальної, об'єктивної та адаптивної системи оцінювання. Огляд літератури та виявлені проблеми дозволяють сформулювати чіткий напрям подальших досліджень, який і становить основу розробки у даній кваліфікаційній роботі.

Адаптивного семантико-граматичного коефіцієнта (АСГК) покликаний забезпечити баланс між семантичною точністю та граматичною коректністю, змінювати цей баланс адаптивно залежно від рівня мовної компетентності користувача (CEFR), формувати більш точні та справедливі оцінки, які враховують як зміст, так і форму мовлення, підвищити інтерпретованість оцінювання через чітке показування внеску граматики та семантики. Цей коефіцієнт не замінює трансформерну модель — він доповнює її, працюючи як модуль надбудови над ембеддингами й оцінками трансформера.

Адаптивність CEFR-залежності коефіцієнта є ключовою, оскільки вимоги до мовлення суттєво змінюються залежно від рівня володіння мовою:

**Рівні A1–A2** — на початкових рівнях головне — зрозумілість базового змісту та формальна правильність конструкцій. Семантика на цьому рівні проста, тому граматики має більшу вагу. У типовій відповіді A1 граматична помилка змінює весь сенс речення.

**Рівні B1–B2** — тут зростає складність мовлення: важливо водночас і змістове наповнення, і коректна побудова речень. У цьому діапазоні семантична та граматична складові мають бути збалансованими.

**Рівні C1–C2** — на високих рівнях основною цінністю стає глибока семантика, когерентність та логічність викладу. Невеликі граматичні огріхи не є критичними, але порушення зв'язності чи нечітка аргументація — суттєві.

Тому семантика має більшу вагу, а граматичний компонент — меншу. Сучасні трансформери та мутаційні метрики не враховують цих рівневих відмінностей, застосовуючи однакову шкалу для всіх категорій користувачів. Причини відсутності адаптивних методів у сучасних моделей:

1. Трансформерні мережі оптимізуються під узагальнену семантику, а не під рівневі мовні компетентності.
2. Відсутні архітектурні механізми, які б дозволяли змінювати ваги семантики та граматики залежно від задачі.
3. Метрики оцінювання статичні — їхня формула не змінюється для різних груп користувачів.
4. Відсутність навчальних датасетів, де б граматики та семантика були чітко розмічені та подані окремо.
5. Сучасні LLM-підходи працюють як “чорні скриньки”, не маючи інтерпретованих внутрішніх коефіцієнтів.

У результаті системи не здатні пояснити, чому певна відповідь отримала саме таку оцінку. Два ключових елементи архітектури трансформера становлять джерело цих обмежень: Attention-механізм не розрізняє граматичні категорії, Attention-матриці відображають важливість токена щодо інших токенів, залежності між словами, контекстуальні ваги.

Але attention не містить інформації про граматичні категорії:

- немає поняття частини мови у “вбудованому” вигляді;

- немає моделі синтаксичних залежностей;
- однакові attention-патерни можуть відповідати реченням з неправильним синтаксисом.

Це означає, що модель розуміє «про що текст», але не «наскільки він правильно побудований». Ембеддинги вміщують значення слова, контекстні залежності, статистичні патерни корпусу. Однак вони не кодують прямо факт помилки, а саме: неправильне узгодження, некоректний порядок слів, відсутність артикля, не правильну часову форму. Тому трансформер “вважає” помилкові конструкції допустимими, якщо вони семантично зрозумілі.

Запропонований у роботі метод передбачає введення нового науково обґрунтованого елемента — адаптивного семантико-граматичного коефіцієнта. Його наукова новизна проявляється у кількох аспектах.

1. Додається новий незалежний коефіцієнт, який формально описує взаємодію семантичної та граматичної складових.
2. Коефіцієнт є комбінованим, він містить два компоненти семантичний (semantic contribution) та граматичний (grammatical contribution).
3. Коефіцієнт адаптивний він автоматично змінює вагу кожного компонента залежно від рівня користувача (CEFR), типу завдання (есе, коротка відповідь, опис, аргументація), профілю помилок користувача (згідно з детекцією помилок).
4. Створюється інтерпретована метрика, яка дозволяє чітко пояснити, чому отримано певний бал.

Запропонований АСГК дозволяє компенсувати недостатню чутливість трансформерів до граматики шляхом окремого введення граматичного компонента, усунути залежність оцінювання від “семантичної домінанти” embeddings, забезпечити різні вагові схеми для A1, B1, C1 без зміни самої трансформерної архітектури, підвищити інтерпретованість за рахунок чіткого розділення внеску семантики та граматики в підсумковий бал, зробити оцінювання справедливішим та ближчим до реальних критеріїв міжнародних мовних екзаменів.

## **2 РОЗРОБКА МЕТОДУ ПОКРАЩЕННЯ ТРАНСФОРМЕРНИХ СИСТЕМ АВТОМАТИЧНОЇ ОЦІНКИ РІВНЯ ВОЛОДІННЯ МОВОЮ ШЛЯХОМ ВВЕДЕННЯ АДАПТИВНОГО СЕМАНТИКО-ГРАМАТИЧНОГО КОЕФІЦІЄНТА**

### **2.1 Постановка задачі та загальна концепція методу**

У сучасних системах автоматичного оцінювання письмового та усного мовлення, побудованих на трансформерних моделях, зберігається низка фундаментальних обмежень, що знижують точність, інтерпретованість та надійність результатів. Хоч трансформери (BERT, GPT, XLM-R, Wav2Vec2 тощо) демонструють високу ефективність у виявленні семантичних зв'язків у тексті, їхня здатність оцінювати граматичні аспекти, синтаксичну структуру та логічну організацію висловлювання залишається недостатньою. Виникає дисбаланс між семантичною схожістю та граматичною коректністю, що особливо критично для учнів рівнів A1–B1, де граматичні помилки є базовим індикатором прогресу.

Метою даного етапу роботи є вдосконалення процесу трансформерного оцінювання шляхом розробки методу, який дозволяє моделі динамічно враховувати як семантичні, так і граматичні характеристики відповіді, адаптуючи їхню важливість до рівня мовної компетентності користувача. Формально задача полягає у створенні адаптивного коефіцієнта, що коригує вагу граматики та семантики при формуванні підсумкової оцінки, забезпечуючи: усунення дисбалансу між семантичною і граматичною складовими, підвищення інтерпретованості оцінки за рахунок явного поділу внесків компонент, динамічну адаптацію до рівня користувача (A1–C2), типу завдання та профілю помилок, зменшення linguistic bias, особливо для носіїв рідкісних мов або нестандартних акцентів.

Для досягнення поставленої мети пропонується введення адаптивного семантико-граматичного коефіцієнта (АСГК) — нових параметрів  $\alpha$  і  $\beta$ , що модулюють баланс між семантичною ( $S_{sem}$ ) та граматичною ( $S_{gram}$ ) оцінками. У загальному вигляді підсумкова оцінка моделі визначається за формулою (2.1):

$$Score = \alpha(L) \cdot S_{sem} + \beta(L) \cdot S_{gram}, \quad (2.1)$$

де:

- $S_{gram}$  — граматична точність, отримана з граматичних детекторів помилок (Grammar-checking head),
- $S_{sem}$  — семантична узгодженість, розрахована на основі embedding-подібності (semantic similarity head),
- $\alpha$  і  $\beta$  — адаптивні коефіцієнти, що змінюється залежно від рівня учня, типу завдання та статистики помилок.

АСГК не замінює трансформер, а вбудовується як надбудова, яка працює на виходах уже існуючих моделей:

- у BERT-подібних моделях він використовується після етапу отримання CLS-векторів або sentence embeddings;
- у GPT-подібних моделях — інтегрується в модуль оцінювання (scoring head) після отримання згенерованого або проаналізованого тексту;
- у мультимодальних моделях мовлення — комбінується з показниками фонетичної та просодичної точності, доповнюючи їх семантичними представленнями.

Таким чином, у даному розділі було визначено задачу та концепцію методу: забезпечити більш точне, прозоре та адаптивне оцінювання мовної продукції за рахунок математично формалізованого коефіцієнта, що компенсує структурні недоліки трансформерних моделей.

## 2.2 Теоретичні основи оцінювання семантичної та граматичної правильності у трансформерних моделях

Оцінювання мовних навичок на основі трансформерних архітектур традиційно базується на двох фундаментальних компонентах — семантичній точності та граматичній правильності. Для побудови адаптивного семантико-граматичного коефіцієнта (АСГК) необхідно формально визначити обидві складові, а також зрозуміти, яким чином трансформери репрезентують семантику та граматику у своїх прихованих просторах.

### 2.1.1 Формалізація семантичної складової $S_{sem}$

Семантична правильність у контексті автоматичного оцінювання мовлення стосується того, наскільки зміст відповіді користувача відповідає очікуваному змісту, зберігає логічність і узгодженість. Вона включає такі: аспекти семантична подібність, збереження змісту, когерентність і зв'язність.

Семантична подібність відображає ступінь близькості значень двох текстів (речень, абзаців, відповідей), незалежно від їхньої поверхневої форми. У трансформерних моделях вона, як правило, визначається через порівняння *sentence embeddings* — векторних подань речень у багатовимірному прихованому просторі.

У загальному випадку, якщо  $h_{ref}$  та  $h_{usr}$  — ембеддинги референсного (еталонного) та користувацького тексту, семантична подібність оцінюється за допомогою косинусної подібності (2.2):

$$S_{sem}^{sim} = \cos(h_{ref}, h_{usr}) = \frac{h_{ref} * h_{usr}}{\|h_{ref}\| \|h_{usr}\|}, \quad (2.2)$$

Джерелом ембеддингів можуть бути:

контекстні вектори CLS-токена у BERT-/RoBERTa-подібних моделях;

sentence embeddings, отримані шляхом пулінгу (mean/max-pooling) по всіх токенах;

спеціалізовані моделі для семантичної подібності (наприклад, модифіковані трансформери для STS-завдань).

У низці робіт, присвячених аналізу семантичної подібності на базі трансформерних ембеддингів, показано, що cosine similarity у векторному просторі моделей BERT-/RoBERTa-подібного типу є ефективним інструментом для порівняння змісту речень і абзаців, зокрема в задачах класифікації текстів, пошуку схожих фрагментів та визначення дублювання змісту [1].

Окремий клас методів — навчені метрики, такі як BLEURT, — використовують BERT-подібні ембеддинги, але замість прямої косинусної відстані навчають додаткову регресійну голову, яка апроксимує людські оцінки якості тексту. Це дозволяє моделі вловлювати не лише “геометричну” подібність, а й більш тонкі аспекти семантики .

У контексті автоматичної оцінки володіння мовою семантична подібність використовується для вимірювання того, наскільки відповідь студента покриває очікуваний зміст, відповідає темі завдання та не містить смислових відхилень від референсної відповіді.

Під збереженням змісту розуміють не тільки загальну семантичну близькість, а й те, чи не спотворено факти, не втрачено ключову інформацію, дотримано логічні зв’язки між твердженнями.

Якщо семантична подібність відповідає на питання “наскільки ці два тексти схожі”, то meaning preservation відповідає на питання “чи не порушена істинність і логіка вихідного змісту”.

Природнім інструментом для формалізації цього аспекту є задачі Natural Language Inference (NLI). У NLI моделі отримують пару “premise–hypothesis” і мають визначити, чи:

- entailment — гіпотеза логічно випливає з преміси;
- neutral — гіпотеза не суперечить, але й не випливає;
- contradiction — гіпотеза суперечить премісі.

Для оцінювання збереження змісту у відповідях студента premise може відповідати референсній (еталонній) відповіді або умові завдання, а hypothesis — відповіді користувача. Тоді:

- entailment → високий рівень збереження змісту;
- neutral → часткове збереження або периферійний зміст;
- contradiction → спотворення фактів або логічна помилка.

Сучасні роботи демонструють, що трансформерні моделі (BERT, RoBERTa, DeBERTa) у задачах NLI досягають високої якості й можуть бути використані як “семантичний фільтр”, який виявляє логічні помилки та спотворення у тексті студента [2].

Крім того, NLI-підхід дозволяє відрізнити:

- лексичні парафрази (коли форма змінена, але зміст збережено);
- логічні перетворення (узагальнення, уточнення, конкретизація);
- смислові зсуви (часткова або повна втрата ключових елементів).

Таким чином, meaning preservation є більш глибинною характеристикою, ніж проста подібність: воно враховує напрямок логічного слідування між еталоном і відповіддю, що критично для коректної оцінки, наприклад, в завданнях опису графіків, переказу тексту або інтерпретації ситуацій.

Когерентність (coherence) описує логічну цілісність тексту на рівні дискурсу: наскільки послідовно розвиваються думки, чи є текст тематично структурованим, чи відповідають одна одній частини висловлювання.

Зв’язність (cohesion) фокусується на мовних засобах, що забезпечують зв’язок між реченнями:

- займенники, повтори, синонімічні заміни,
- дискурсивні маркери (“однак”, “по-перше”, “внаслідок цього”),
- референції на попередні елементи тексту.

У трансформерних моделях когерентність та зв’язність можуть оцінюватися через:

1. Граф дискурсу. Текст розглядається як послідовність речень/сегментів, між якими встановлюються дискурсивні відношення (причина–наслідок, контраст, деталізація тощо). На основі цих відношень будується граф (або дерево) дискурсу, і далі аналізується його структура: наскільки логічно “зшитий” текст. Сучасні роботи показують, що поєднання інформації про дискурсивні відношення з трансформерними ембеддингами покращує якість оцінювання когерентності та есе-скорю [3].

2. Перехресні attention-зв'язки між реченнями. Self-attention у трансформерах дозволяє бачити, на які токени / речення модель “дивиться” при обробці кожного фрагмента тексту. Якщо модель стабільно приділяє увагу релевантним попереднім реченням (наприклад, суб'єктам, згаданим раніше), це вказує на збереження локальної когерентності. Якщо ж attention-структура хаотична і не корелює з дискурсивним ланцюжком, текст, як правило, є менш зв'язним. Низка нейронних моделей когерентності будує додаткові представлення сутностей і дискурсивних відношень, а потім інтегрує їх з attention-механізмами, що дозволяє краще виявляти порушення когерентності [4].

3. Тематична узгодженість (topic consistency). Когерентний текст зазвичай характеризується відносно стабільним тематичним вектором, тобто ембеддинги речень мають належати до близького підпростору. Різкі стрибки в семантичному просторі, коли нове речення “відривається” від попереднього кластера, часто сигналізують про: відступи від теми, вставні, незв'язані фрагменти, порушення логічної структури.

Деякі сучасні роботи з нейронного моделювання когерентності комбінують entity-based підхід (відстеження згаданих сутностей) та дискурсивні відношення, формуючи єдину модель оцінювання структури тексту, що демонструє високу кореляцію з людськими судженнями щодо якості тексту [5].

## 2.1.2 Формалізація граматичної складової Sgram

Граматична складова Sgram описує структурну правильність мовлення та включає кілька формальних компонентів, що підтверджено сучасними дослідженнями у сфері автоматичного оцінювання текстів. Як зазначено у роботах з граматичної корекції на основі трансформерів [6], граматики повинна оцінюватися не лише через кількість помилок, але й через їхній тип, зваженість та контекст. Точні POS-теги та морфологічні ознаки мають важливе значення для оцінювання граматики — сучасні бібліотеки на основі трансформерів, такі як Stanza [7], забезпечують високоточний аналіз синтаксичних ролей.

У дослідженнях з автоматичного аналізу мовлення наголошується, що **узгодження підмета і присудка, часові форми та морфологічні маркери** є ключовими параметрами для вимірювання граматичної складової [8].

$$S_{pos} = 1 - \frac{\text{морфологічні помилки}}{\text{к-сть токенів}}, \quad (2.3)$$

Аналіз залежностей (dependency parsing), згідно з дослідженнями з CoNLL та Universal Dependencies [9], дозволяє формально визначити:

- чи правильно пов'язані слова;
- чи не порушена структура речення;
- наскільки граматично «здоровим» є синтаксичне дерево.

$$S_{dep} = \frac{\text{правильні синтаксичні залежності}}{\text{усі залежності}}, \quad (2.4)$$

Цей підхід підтримується також у роботах, що застосовують трансформерні embeddings для освітніх задач [10].

Сучасні системи граматичної корекції, такі як GECToR [6], доводять ефективність підходу на основі тегування (tag-based correction) для оцінки помилок у текстах. Ця методологія дозволяє формалізувати:

- морфологічні помилки,
- синтаксичні зсуви,

- порушення порядку слів,
- некоректні конструкції.

У багатьох публікаціях вказано, що якість граматики може бути об'єктивно виміряна через кількість і тип помилок, а не лише через їх наявність [8].

$$Serr = 1 - \frac{\text{зважені помилки}}{\text{довжина тексту}}, \quad (2.5)$$

На думку дослідників, perplexity трансформерної моделі може служити наближеним показником граматичної правильності [8], оскільки неправильні конструкції мають низьку ймовірність у мовних моделях.

$$Sflu = \frac{1}{\text{perplexity}}, \quad (2.6)$$

### 2.1.3 Кодування трансферами семантики та граматики

Сучасні трансформерні моделі (BERT, RoBERTa, GPT-3/4, XLM-R) формують уніфіковані контекстуальні представлення, але семантичні та граматичні ознаки «розподіляються» у шарах архітектури нерівномірно. Це підтверджено експериментами з probing-методами, які дозволяють визначити, яку інформацію зберігають різні шари й attention-голови [11].

Семантична інформація у трансформерах — це значення фраз, контекст слів, логічні зв'язки й інтерпретація смислу. Вона формується переважно у верхніх шарах моделі. Основні джерела семантичної інформації:

1. Глобальні attention-патерни. У високих шарах attention має «широкий радіус дії»: він пов'язує віддалені слова або навіть цілі фрази. Це дозволяє моделі: розуміти відношення між частинами тексту, формувати тематичні та логічні залежності, здійснювати reasoning. Дослідження [12]

показують, що саме в цих шарах з'являються патерни, пов'язані з семантичними ролями, «хто що робить», «що до чого належить».

2. Високорівневі шарові представлення. У верхніх шарах нейронна мережа абстрагує інформацію настільки, що: класифікація смислу (intent detection), NLI (entailment/contradiction), класифікація теми, парафразування стають більш стабільними, ніж у нижніх шарах [13].

3. Контекстні embeddings слів (contextualized vectors). На відміну від статичних ембеддингів (Word2Vec), трансформери формують вектори залежно від контексту, тому слова: із різними ролями (homonyms), у різних граматичних позиціях, матимуть різні семантичні представлення. Це дозволяє моделі розуміти: значення фрази в конкретному контексті, логічні прогалини, семантичні суперечності.

В результаті цього верхні шари трансформера оптимізовані під смислову інтерпретацію, а не під граматику. Це робить їх дуже точними в семантичних задачах, але недостатньо чутливими до синтаксичних помилок.

На відміну від семантики, граматику кодується локально та фрагментарно. Дослідження probing-моделями [11] показали, що нижні шари BERT фіксують: POS-інформацію (іменник / дієслово / прикметник), морфологічні маркери (число, рід, час) та базові синтаксичні зв'язки. Це підтверджує, що граматику «захоплюється» на ранніх етапах обробки.

У той час як семантика будується на глобальних патернах: граматику формується з локальних залежностей між сусідніми або близькими токенами. Моделі фіксують найбільш очевидні синтаксичні зв'язки, але погано працюють із довгими залежностями.

Морфологія кодується у спеціалізованих векторних «підпросторах», які виділяються за допомогою пояснювальних методів. Головною проблемою трансформерів є те, що вони не будують внутрішнього синтаксичного дерева. Через це вони можуть пропустити: неправильне узгодження (he go, she have), граматичні суперечності у вкладених реченнях, порушення порядку слів (особливо у мовах зі складним синтаксисом), перестановки, що змінюють

значення. Це підтверджено у роботах [8], де доведено, що трансформери можуть повністю «ігнорувати» синтаксичні помилки, якщо семантичний зміст зберігається.

Незважаючи на високу потужність, self-attention має низку фундаментальних обмежень у контексті граматики.

1. Механізм attention не розрізняє типи лінгвістичних залежностей. Self-attention обробляє токени без урахування того, чи зв'язок між ними семантичний чи граматичний. Тобто немає поділу: підмет ↔ присудок, означення ↔ означуване, керування ↔ узгодження.

Це призводить до того, що модель «бачить» текст як набір векторів, а не як граматичне дерево.

2. Self-attention не має вбудованого механізму виявлення помилок. На відміну від синтаксичних парсерів, трансформери: не можуть визначити, чи речення граматично правильне, не розрізняють структури, які суперечать правилам мови, не вміють реконструювати дерево залежностей, бо attention — це не структурна модель.

3. Проблеми з довгими реченнями. У довгих реченнях: граматичні залежності стають довгими, локальний attention «втрачає» ключові зв'язки, модель може неправильно інтерпретувати структуру фрази.

Це підтверджено у роботах щодо «syntactic degradation» у великих мовних моделях [14].

4. Чутливість до рідкісних та складних конструкцій. Self-attention не має лінгвістичної узагальненості, чим рідше конструкція зустрічалась у тренувальних даних, тим більше шанс, що модель її неправильно оцінить.

Аналіз показує:

1. Семантичні ознаки кодуються сильно й глобально, у вищих шарах.
2. Граматичні ознаки кодуються слабко й локально, лише у перших шарах.
3. Сучасні трансформери та метрики (BLEURT, BERTScore, G-Eval): добре оцінюють зміст, погано — граматичну правильність.

4. Рівень володіння мовою впливає на пріоритети оцінювання:  
 A1–A2 → критична граматики,  
 B2–C2 → критична семантика, когерентність та логіка.

Наразі немає моделі, яка б адаптивно змінювала вагу семантичної та граматичної складової, визначала наскільки кожен компонент вплинув на оцінку, забезпечувала інтерпретованість. Тому і пропонується Адаптивний Семантико-Граматичний Коефіцієнт (АСГК). Він дає динамічний баланс граматики та семантики, вищу точність оцінювання для різних рівнів володіння мовою, можливість пояснення та інтерпретованості, усунення ключового недоліку трансформерів — асиметрії кодування.

### 2.3 Концепція адаптивного семантико-граматичного коефіцієнта

Адаптивний семантико-граматичний коефіцієнт (АСГК) пропонується як універсальний механізм динамічного балансування впливу семантичних та граматичних характеристик на підсумкову оцінку мовленнєвого висловлювання. На відміну від класичних метрик, які мають жорстко фіксовані ваги або орієнтуються переважно на семантику (BERTScore, BLEURT) чи на граматику (GEC-метрики), АСГК забезпечує адаптивність, що залежить від рівня мовця та типу мовного завдання. Адаптивний семантико-граматичний коефіцієнт (АСГК) визначається як функція, що здійснює динамічне вагове балансування між семантичною та граматичною складовими оцінювання мовного висловлювання залежно від рівня володіння мовою (CEFR) та типу завдання.

Математичні властивості вагових функцій:

1. Адаптивність від рівня мовця. Ваги моделюються сигмоїдальною або іншою монотонною функцією, що забезпечує плавний перехід від граматичного до семантичного домінування (2.7).

$$\alpha(L) = \sigma(w_1 L + b_1), \beta(L) = 1 - \alpha(L), \quad (2.7)$$

2. Врахування типу завдання. Для кожного завдання  $T$  вводяться коефіцієнти зміщення:

$$\alpha(L, T) = \sigma(w_1 L + b_1 + t\alpha(T)), \quad (2.8)$$

$$\beta(L, T) = 1 - \alpha(L, T), \quad (2.9)$$

Наприклад:

- для перефразування  $t\alpha > 0$  (перевага семантики),
- для граматичного редагування  $t\alpha < 0$  (перевага граматики).

АСГК — це параметризована адаптивна функція, що здійснює нелінійне зважування семантичних та граматичних показників відповідно до рівня мовної компетентності користувача та типу виконуваного завдання, забезпечуючи оптимальний баланс між змістовною та структурною точністю висловлювання.

Для порівняння статичних метрик зі АСГК зобразимо таблицю 2.1:

Таблиця 2.1

Метрика	Фокус	Недолік	Порівняння з АСГК
<b>BERTScore</b>	Семантика	Не бачить граматичних помилок	АСГК інтегрує граматику через $\beta$
<b>BLEURT</b>	Семантика + фрагмент граматики	Залежність від тренувальних даних	АСГК зміщує вагу залежно від рівня
<b>G-Eval (GPT-4)</b>	Рубрики (grammar, coherence, etc.)	Ваги фіксовані та не адаптивні	АСГК налаштовує ваги під користувача
<b>Grammar-only GEC метрики</b>	Морфологія/синтаксис	Ігнорують зміст	АСГК забезпечує баланс

Таким чином, АСГК не замінює ці моделі, а надбудовується над ними, забезпечуючи адаптивне зважування. Використання АСГК дає низку переваг:

1. Персоналізація оцінювання. Модель підлаштовується під рівень мовця, профіль помилок та специфіку завдання.

2. Баланс семантики та граматики. Уникається ситуація, коли граматично правильний, але порожній за змістом текст отримує високу оцінку, або навпаки — глибокий, але граматично недосконалий текст надмірно штрафується.

3. Підвищення валідності мовного тестування. Адаптивні оцінювачі значно ближчі до реальної експертної оцінки, про що свідчать сучасні дослідження Duolingo, ETS та Google Research.

4. Гнучкість та розширюваність. АСГК може бути інтегрований у будь-яку архітектуру трансформера (BERT-типу або GPT-типу), не вимагаючи модифікації її внутрішніх параметрів.

5. Можливість навчання на даних користувача. У майбутньому коефіцієнти  $\alpha$  та  $\beta$  можуть оптимізуватися градієнтно, на основі коректних оцінок експертів.

Таким чином, введення АСГК дозволяє створити універсальний, адаптивний та інтерпретований механізм оцінювання, що забезпечує більш точну характеристику мовної компетентності користувача.

## **2.4 Методика обчислення семантичної складової Ssem**

Семантична складова Ssem визначає, наскільки відповідь користувача зберігає зміст, логіку та смислову структуру відносно еталонної відповіді або інструкції. Цей показник вимірює не граматичну правильність, а якість передачі значення, що включає: відповідність темі, контекстуальну релевантність, зв'язність та логічну послідовність викладу.

Методика обчислення семантичної складової базується на сучасних трансформерних моделях (BERT, RoBERTa, XLM-R, Sentence-BERT, GPT), які формують контекстуальні представлення висловлювань у високорозмірному просторі.

Основою семантичної оцінки є векторне представлення тексту, сформоване трансформером. Для цього використовуються:

1. Контекстуальні ембеддинги BERT / RoBERTa. Моделі типу BERT (Devlin et al., 2019) та RoBERTa формують вектори: CLS-ембеддинг, середнє значення токенів (mean pooling), максимальне значення (max pooling). Проте BERT не оптимізований для семантичної схожості, тому сучасні дослідження рекомендують Sentence-BERT.

2. Sentence-BERT та SBERT-score. Sentence-BERT (Reimers & Gurevych, 2020) створює вектори, оптимізовані для задач: semantic similarity, sentence matching, paraphrase detection. Це дозволяє напряду порівнювати два речення у векторному просторі.

3. Багатомовні моделі (XLM-R, LaBSE). У задачах оцінювання на різних мовах найбільш ефективні: XLM-RoBERTa, LaBSE, mUSE (multilingual Universal Sentence Encoder). Такі моделі забезпечують стабільність оцінки для користувачів з різними L1.

4. GPT-типу моделі. Для деяких типів аналізу (особливо reasoning-based) можна використовувати hidden-state embeddings GPT-3/4, які зберігають глибші семантичні залежності.

## **2.5 Методика обчислення граматичної складової Sgram**

Грамматична складова Sgram вимірює ступінь відповідності тексту нормам синтаксису, морфології, пунктуації та граматичним конструкціям. На відміну від семантики, граматики є вразливою до локальних помилок і потребує чіткого

виявлення, класифікації та вагового оцінювання помилок залежно від їх тяжкості.

Сучасні дослідження у сфері автоматичного виправлення граматичних помилок (Grammar Error Correction, GEC) показують, що трансформерні моделі (GECtoR, PIE, T5-GEC, GPT-GEC) досягають високої точності, але й далі генерують різні типи помилок при оцінюванні тексту користувача, тому систематичний підхід до формування метрики є критично важливим.

Для обчислення граматичної складової використовується багаторівневий grammar error detection (GEC)-процес. У практиці мовного оцінювання найчастіше застосовуються такі моделі:

1. GECtoR модель працює за принципом *tag-based correction*, тобто класифікує кожен токен як: правильний, потребує заміни, потребує вставки, потребує видалення. Вона особливо добре виявляє локальні морфологічні та синтаксичні помилки.

2. PIE — це sequence-to-sequence-трансформер, який виявляє складні синтаксичні конструкції та помилки залежностей (dependency errors).

3. T5-GEC — варіанти моделі T5, натреновані на GEC-корпусах (BEA, FCE, JFLEG), добре виявляють: пунктуаційні помилки, узгодження, неправильний порядок слів.

4. GPT-GEC (2023–2024) - великі мовні моделі (GPT-3.5, GPT-4, Llama-3) можуть виконувати grammar error detection через: tagging, scoring (визначення ймовірності «граматичності»), chain-of-thought-аналіз структурних помилок.

Для обчислення граматичної складової застосовується класифікація помилок на чотири основні групи:

1. Морфологічні помилки: неправильна форма слова ('goes' → 'go'), узгодження числа і особи ('he go'), неправильні часові форми.

2. Синтаксичні помилки: неправильний порядок слів, порушення структури складнопідрядного речення, неповні конструкції (missing subject/object).

3. Лексичні помилки: невдалий вибір слова (word choice errors), кальки з рідної мови, неправильні прийменники.

4. Пунктуаційні помилки: пропущені коми, надлишкові розділові знаки, неправильне оформлення складних речень.

У сучасних ГЕС-корпусах використовуються аналогічні класифікації (BEA-2019 Shared Task; Nagata 2021). Помилки різного типу мають різний вплив на сприйняття тексту, тому використовують зважене оцінювання. У статті Yu & Wan (2023) [8] показано, що синтаксичні та морфологічні помилки найсильніше впливають на зрозумілість тексту, тому мають більшу вагу. Для прикладу порівняння системи ваг зобразимо таблицю 2.2:

Таблиця 2.2

Тип помилки	Вага	Пояснення
Синтаксична	0.40	Впливає на структуру речення
Морфологічна	0.30	Впливає на узгодження
Лексична	0.20	Впливає на коректність значення
Пунктуація	0.10	Мінімальний вплив на зміст

## **2.6 Адаптивна функція, що поєднує семантичну та граматичну складові**

У попередніх підрозділах було формально визначено семантичну складову  $S_{sem}$  та граматичну складову  $S_{gram}$ . На цьому етапі вводиться адаптивна функція об'єднання, яка враховує рівень володіння мовою (CEFR), тип завдання, профіль помилок та складність тексту, і на цій основі змінює ваги семантики й граматики.

Узагальнений вигляд підсумкової оцінки пропонується записати як (2.10):

$$Score = K_{ASGK} \cdot (\alpha(L, T)Ssem + \beta(L, T)Sgram), \quad (2.10)$$

де

$Ssem$  — нормована семантична оцінка,

$Sgram$  — нормована граматична оцінка,

$\alpha(L, T)$  — динамічна вага семантики,

$\beta(L, T)$  — динамічна вага граматики,

$K_{ASGK}$  — адаптивний семантико-граматичний коефіцієнт, що масштабує результат з урахуванням додаткових факторів.

Ваги задовольняють умову нормалізації (2.11):

$$\alpha(L, T) + \beta(L, T) = 1, 0 \leq \alpha, \beta \leq 1, \quad (2.11)$$

Фактор  $K_{ASGK}$  інтерпретується як коригувальний множник на рівні всієї відповіді, який враховує, наскільки поточне поєднання семантики та граматики є очікуваним для певного рівня CEFR та типу завдання:

- якщо профіль помилок та складність тексту відповідають рівню користувача  $\rightarrow K_{ASGK} \approx 1$ ;
- якщо спостерігається дисбаланс (наприклад, занадто простий текст для C1 або надто багато критичних помилок)  $\rightarrow K_{ASGK}$  зменшується.

Динамічні ваги  $\alpha$  та  $\beta$

на початкових рівнях (A1–A2) важливіше контролювати граматику,

на просунутих рівнях (C1–C2) вирішальну роль відіграє семантика, когерентність і логічна структура.

Тому  $\alpha$  та  $\beta$  визначаються як функції від:

LLL — рівня CEFR (A1–C2),

ТТТ — типу завдання (есе, переказ, опис зображення, граматичний тест тощо).

Приклад монотонної залежності від рівня (2.12):

$$\alpha(L) = \sigma(w1L + b1), \beta(L) = 1 - \alpha(L), \quad (2.12)$$

де  $LLL$  — чисельно закодований рівень (наприклад,  $A1=1, \dots, C2=6$ ), а  $\sigma$  — сигмоїда. Зростання рівня  $LLL$  збільшує вагу семантики  $\alpha(L)$  та зменшує вагу граматики  $\beta(L)$ .

Залежність від типу завдання реалізується як зсув (2.13):

$$\alpha(L, T) = \sigma(w1L + b1 + t\alpha(T)), \beta(L, T) = 1 - \alpha(L, T), \quad (2.13)$$

де  $t\alpha(T)$  — завданнєвий параметр:

для переказу, опису графіка, аргументаційного есе:  $t\alpha(T) > 0 \rightarrow$  зростає частка семантики;

для граматичних вправ, редагування тексту:  $t\alpha(T) < 0 \rightarrow$  збільшується частка граматики.

Подібна ідея адаптивного зміщення ваг, залежно від рівня й типу завдання, відповідає сучасним підходам до автоматичного CEFR-оцінювання, де використовують різні лінгвістичні ознаки для різних рівнів та навичок.

Алгоритм адаптації ваг: CEFR-класифікатор, профіль помилок, складність тексту. Адаптація ваг  $\alpha(L, T)$  та  $\beta(L, T)$ , а також обчислення  $K_{ASGK}$ , базуються на трьох групах сигналів.

### 1) Оцінка рівня CEFR через класифікатор

Застосовується модель класифікації CEFR-рівня на основі трансформерів (BERT/XLM-R), яка отримує на вхід текст відповіді та (за потреби) завдання й повертає прогнозований рівень  $\hat{L}$ . Подібні моделі описані у дослідженнях з автоматичної класифікації письмових робіт за рівнями CEFR. Отримане  $\hat{L}$  використовується як вхід у функції  $\alpha(\hat{L}, T)$  та  $\beta(\hat{L}, T)$ .

## 2) Аналіз профілю помилок

На основі модулів GEC будуються: кількість і типи помилок (морфологічні, синтаксичні, лексичні, пунктуаційні), їх розподіл (локальні vs глобальні, критичні vs другорядні), частка помилок на 100 слів.

Профіль помилок порівнюється з типовими профілями для кожного CEFR-рівня (наприклад, A2, B1, B2), як це описано в роботах з автоматичної оцінки мовної компетентності та аналізу помилок.

Якщо для поточного тексту спостерігається надлишок грубих граматичних помилок порівняно з типово очікуваним для  $\hat{L}$ , коефіцієнт  $K_{ASGK}$  зменшується (штраф), а  $\beta$  тимчасово збільшується для посилення ролі граматики.

## 3) Оцінка рівня складності тексту

Додатково визначається лінгвістична складність тексту (lexical diversity, довжина речень, синтаксична глибина, частка рідкісних слів), наприклад, за підходами до автоматичної класифікації складності текстів.

Якщо складність тексту суттєво вища за очікувану для рівня  $\hat{L}$ , але при цьому граMATика слабка, система може: знизити  $K_{ASGK}$  (боротьба з «надмірно складним, але неграмотним» текстом), або скоригувати  $\alpha, \beta$ , підсилюючи вагу Sgram.

Алгоритмічну схему обчислення  $K_{ASGK}$  та Score рис. 2.3 можна описати так:

1. Вхідні дані: текст відповіді користувача, текст завдання / референсна відповідь, тип завдання ТТТ.
2. Обчислення семантики та граматики: обчислити Ssem і обчислити Sgram.

3. Оцінка рівня CEFR: пропустити текст через CEFR-класифікатор → отримати  $\hat{L}$ .
4. Аналіз профілю помилок: зібрати кількість/типи помилок → сформуванню профілю та порівняти з типовим профілем для  $\hat{L}$ .
5. Оцінка складності тексту: обчислити індекс складності (lexical + syntactic complexity).
6. Обчислення динамічних ваг: знайти  $\alpha(\hat{L}, T)$  та  $\beta(\hat{L}, T)$  за заздалегідь визначеними або навченими функціями.

Роботу системи класифікації можна умовно розділити на два ключові етапи [15]:

1) Етап навчання (див. рис. 2.1). Він виконується одноразово або періодично — у міру надходження нових надійно розмічених текстів. На цьому етапі формується прогнозна модель: алгоритм машинного навчання отримує навчальний набір числових векторів, які були попередньо згенеровані шляхом вилучення лінгвістичних ознак із текстів, вручну анотованих експертами. Результатом є натренована модель, здатна відтворювати патерни, властиві відповідним рівням володіння мовою.

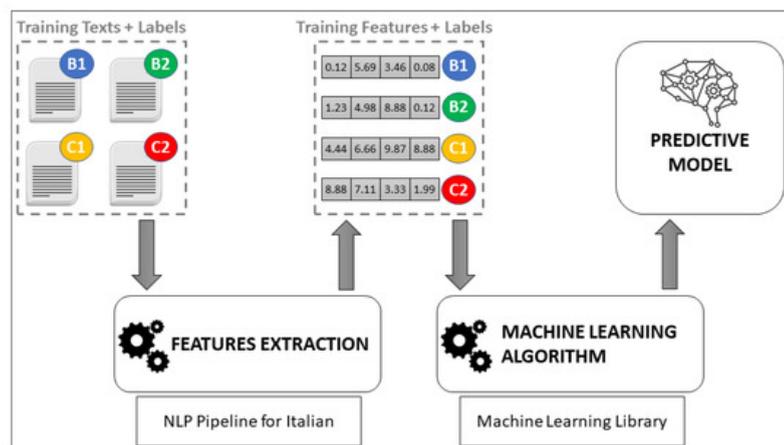


Рис. 2.1 Фаза навчання системи класифікації

2) Етап прогнозування (див. рис. 2.2). На цьому етапі новий, нерозмічений текст проходить процедуру вилучення ознак, після чого модель визначає рівень його складності та відносить його до одного з рівнів B1, B2, C1 або C2.

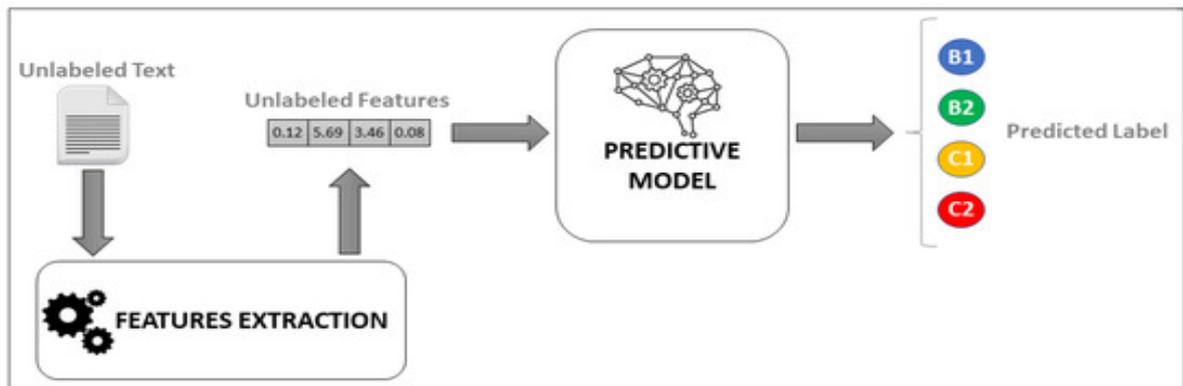


Рис. 2.2 Фаза прогнозування системи класифікації

7. Обчислення  $K_{ASGK}$ : на основі профілю помилок, відхилення складності та можливих додаткових параметрів визначити  $K_{ASGK} \in [K_{min}, K_{max}]$ , наприклад:

$$K_{ASGK} = 1 - \lambda_1 \cdot \Delta errors - \lambda_2 \cdot \Delta complexity$$

де  $\Delta errors$  і  $\Delta complexity$  — відхилення від очікуваних значень для рівня  $\hat{L}$ , а  $\lambda_1, \lambda_2$  — ваги.

8. Фінальна оцінка:

$$Score = K_{ASGK} \cdot (\alpha(\hat{L}, T) S_{sem} + \beta(\hat{L}, T) S_{gram}).$$

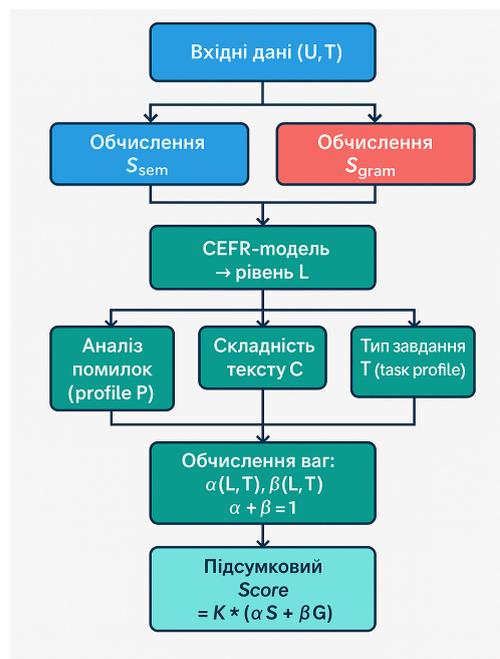


Рис. 2.3 Алгоритм обчислення  $K_{ASGK}$  та Score

Такий підхід узгоджується із сучасними концепціями адаптивного та гібридного мовного оцінювання, де поєднуються лінгвістичні ознаки, рівень складності завдання та профіль помилок для досягнення більш надійної та валідної оцінки.

## 2.7 Інтеграція методу у трансформерну архітектуру

Запропонований метод введення Адаптивного Семантико-Граматичного Коефіцієнта (АСГК) потребує визначення способу його інтеграції у типову трансформерну архітектуру. Оскільки трансформери можуть використовуватися як у режимі енкодерів (BERT-тип), так і декодерів (GPT-тип), а також у гібридних encoder–decoder моделях (T5-тип), важливо вибрати оптимальний спосіб включення коефіцієнта у пайплайн оцінювання. Способи інтеграції коефіцієнта в архітектуру трансформера:

### 2.7.1 Способи інтеграції коефіцієнта в архітектуру трансформера

Інтеграція перед фінальним лінійним шаром (output layer). АСГК можна застосувати безпосередньо до векторів, що подаються у фінальний лінійний класифікаційний шар моделі. Це дозволяє здійснити зважування семантичної та граматичної складових до моменту генерації оцінки.

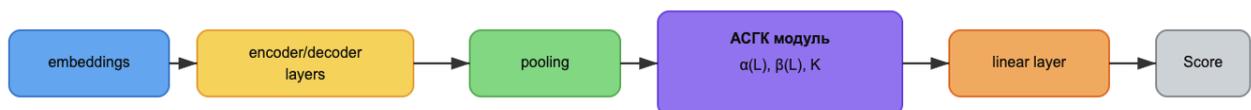


Рис. 2.4 Блок-схема інтеграції коефіцієнта в архітектуру трансформера

Переваги:

не потребує втручання у внутрішні механізми трансформера;

можна вбудувати у будь-яку модель без перенавчання;

дає стабільний контроль над фінальною оцінкою.

Недоліки:

модель не впливає на семантико-граматичну адаптацію всередині представлень;

коефіцієнт працює лише на вихідному рівні.

### 2.7.2 Інтеграція як окремого модуля після отримання embeddings

У цьому випадку трансформер виступає як «feature extractor», а АСГК — як додатковий модуль, що працює над векторами:

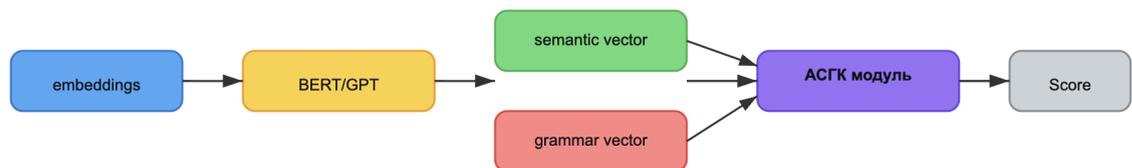


Рис. 2.5 Блок-схема інтеграції коефіцієнта як окремого модуля

Реалізація:

- семантичний та граматичний вектори обчислюються окремими підмодулями;
- результат комбінується за допомогою  $\alpha(L,T)\alpha(L,T)$  та  $\beta(L,T)\beta(L,T)$ ;
- на виході формується Score.

Переваги:

- чітка модульність;
- можливість заміни або оновлення окремих компонентів ( $S_{sem}$ ,  $S_{gram}$ ).

### 2.7.3 Інтеграція як частина scoring head

У цьому варіанті модуль оцінювання (scoring head) розширюється функціями: тут АСГК стає частиною learned-голови моделі, аналогічно тому, як у GPT-оцінюючих моделях використовуються спеціальні Scoring Heads.

Переваги:

- гнучкість;
- можна тренувати ваги  $\alpha$  і  $\beta$  разом із моделлю;
- можливість градієнтної оптимізації всього модуля.

Недоліки:

- потребує тонкого перенавчання моделі.

#### **2.7.4 Порівняння інтеграції у моделі-енкодери та моделі-декодери**

Енкодери, такі як BERT (Bidirectional Encoder Representations from Transformers) та XLM-R (XLM-RoBERTa), призначені для розуміння тексту та формування високоякісних контекстуалізованих векторних представлень (ембедингів). Вони не генерують новий текст, а фокусуються на аналізі вхідної послідовності. Оптимальні варіанти інтеграції АСГК:

Інтеграція після Pooling: після того, як енкодер обробив вхідний текст і створив його представлення (наприклад, використовуючи [CLS] токен або усереднення/максимальне пулінгування всіх токенів), АСГК інтегрується у вигляді додаткового шару або мережі, яка приймає це фінальне агреговане представлення тексту. Це ідеально підходить для задач, де потрібна загальна оцінка всього тексту.

Scoring Head із вбудованим АСГК: Це означає додавання спеціальної "голови" (scoring head) поверх енкодера. Ця голова, яка є невеликою нейронною мережею, навчена не лише виконувати фінальне завдання (наприклад, класифікацію), але й містить або використовує логіку АСГК для винесення оцінки граматики або якості.

Переваги: стабільна робота для класифікаційних задач: енкодери є архітектурно найбільш пристосованими для класифікації (визначення категорії, рівня, правильності). Інтеграція АСГК на рівні представлення забезпечує надійний контроль якості при категоризації.

Оптимальні для CEFR-класифікації та GEC-детекції: CEFR-класифікація (Common European Framework of Reference for Languages) — це визначення

рівня володіння мовою тексту. Якісне представлення енкодера та інтегрований АСГК дозволяють точно визначити рівень за складністю структури та кількістю граматичних помилок. GEC-детекція (Grammatical Error Correction detection) — це виявлення наявності граматичних помилок у тексті.

Низькі вимоги до перенавчання: оскільки енкодери вже мають сильні представлення, інтеграція АСГК може вимагати лише додаткового тонкого налаштування (*fine-tuning*) або навчання лише Scoring Head, що значно економить обчислювальні ресурси.

Декодери, як-от моделі GPT (Generative Pre-trained Transformer), призначені для генерації тексту в авторегресійному режимі (прогнозування наступного токена). Вони сильні в уловлюванні семантичних зв'язків, логіки (*reasoning*) та створенні зв'язного та змістовного тексту. Слабка сторона декодера це контроль граматики/синтаксису. Через фокус на сенсі і потоці генерації, вони можуть пропускати синтаксичні помилки (наприклад, узгодження часів, відмінків) або генерувати граматично неправильні, але семантично зрозумілі речення.

Оптимальні варіанти інтеграції АСГК:

Scoring Head: додатковий механізм, який оцінює якість згенерованого тексту або впливає на процес генерації (наприклад, як частина reward model у Reinforcement Learning from Human Feedback - RLHF). Це дозволяє використовувати семантичну силу GPT і водночас фільтрувати або коригувати результати на предмет граматики.

Фінальний шар оцінювання: інтеграція прямо в останній шар моделі, який відповідає за вибір фінального токена, щоб змити ймовірності вибору на користь граматично коректних варіантів.

Ключове значення динамічної ваги  $\beta(L, T)$ :

Оскільки GPT-моделі схильні пропускати синтаксичні помилки, динамічна вага граматики  $\beta(L, T)$  є особливо важливою.

$\beta(L, T)$  — це коефіцієнт, який динамічно регулює вплив граматичного контролю на функцію втрат або процес декодування.

L — рівень складності завдання або тексту.

T — довжина згенерованої/аналізованої послідовності.

Динамічна вага дозволяє посилити граматичний контроль у випадках, коли модель найбільше схильна до помилок (наприклад, при довгих реченнях або генерації складних структур), або ж послабити його для збереження семантичної гнучкості.

Моделі, як-от T5 (Text-to-Text Transfer Transformer) та mT5 (Multilingual T5), поєднують сильні сторони обох архітектур: енкодер формує представлення вхідного тексту, а декодер генерує вихідний текст на його основі. Це робить їх надзвичайно універсальними. Найкраще Застосування (Інтеграція АСГК): вбудований АСГК у Scoring Head: аналогічно до попередніх моделей, Scoring Head використовується для фінальної оцінки або контролю якості згенерованого тексту. Це найбільш гнучкий та неінвазивний спосіб інтеграції.

Дворівнева інтеграція (розділення завдань): на рівні енкодера (для семантики) АСГК або його компоненти можуть впливати на енкодер, щоб він формував більш якісні та змістовні представлення. Це допомагає гарантувати, що вхідний текст правильно зрозумілий, особливо якщо він містить помилки.

На рівні декодера (для оцінки граматики) АСГК інтегрується в декодер або його Scoring Head для безпосереднього контролю граматичної коректності згенерованого тексту, використовуючи його потужність для авторегресійної генерації з урахуванням граматичних правил.

Цей підхід дозволяє максимально використати можливості обох частин: енкодер — для глибокого семантичного розуміння, а декодер — для граматично коректної генерації.

Для простішого порівняння моделей представлена таблиця 2.3.

Вибір інтеграції автоматичної системи граматичного контролю (АСГК) як окремого модуля над трансформерними ембедингами є стратегічним рішенням, що максимізує гнучкість, незалежність та інтерпретованість системи.

Цей підхід передбачає, що АСГК не вбудовується безпосередньо у внутрішні шари трансформера (наприклад, attention-механізм), а приймає фінальні векторні представлення тексту, згенеровані трансформером, для винесення своєї оцінки.

Таблиця 2.3

Модель/Архітект ура	Основне Призначення	Сильні Сторони (NLP)	Оптимальна Інтеграція АСГК
<b>Енкодери (BERT, XLM-R)</b>	Розуміння та класифікація	Якісні векторні представлення	Після Pooling, Scoring Head
<b>Декодери (GPT-тип)</b>	Генерація тексту	Семантика, Reasoning, зв'язність	Scoring Head, фінальний шар оцінювання
<b>Encoder–Decoder (T5, mT5)</b>	Універсальні (Текст-у-Текст )	Поєднання розуміння та генерації	Вбудований АСГК у Scoring Head, дворівнева інтеграція

Переглянемо основні переваги:

1. Мінімальне втручання у модель. Зберігається цілісність та переднавчена потужність базової трансформерної моделі. Немає потреби змінювати складні та ресурсомісткі компоненти, як-от attention-механізми чи feed-forward мережі. Це зменшує ризик порушення знань, отриманих моделлю під час попереднього навчання, і значно спрощує процес тонкого налаштування (fine-tuning).

2. Незалежність компонентів. Система оцінки ефективно розділяється на дві частини:

Ssem (semantic score) — оцінка, що базується на семантичному розумінні (надається трансформером).

Sgram (grammatical score) — оцінка, що базується на граматичному контролі (надається модулем АСГК).

Кожен компонент може бути покращений або замінений окремо, без необхідності перенавчання чи коригування іншого. Наприклад, можна оновити граматичні правила в АСГК (Sgram), не чіпаючи BERT (Ssem).

3. Можливість заміни трансформера (plug-and-play). Оскільки модуль АСГК працює поверх ембедингів, він є незалежним від конкретної архітектури трансформера. Базовий трансформер може бути легко замінений, наприклад, перехід з BERT на XLM-R для багатомовності, заміна на модель GPT-типу, якщо необхідна більша семантична сила для отримання початкових ембедингів. Це забезпечує довгострокову життєздатність архітектури та легкість її оновлення.

4. Гнучкість адаптації ваг. Фінальна оцінка формується шляхом комбінування оцінок Ssem та Sgram із використанням налаштовуваних ваг. Ваги працюють як незалежний, зовнішній блок. Це суттєво спрощує експерименти та точне налаштування системи під конкретну задачу.

5. Висока інтерпретованість (explainability). Завдяки розділенню на Ssem і Sgram, фінальний бал є прозорим та пояснюваним. Легко пояснити, чому модель дала той чи інший бал. Це критично важливо для навчальних систем та систем зворотного зв'язку, де користувачеві потрібне конкретне пояснення його помилок.

## **2.8 Оцінка обчислювальної складності та переваг запропонованого методу**

Основна перевага запропонованого методу полягає в тому, що він мінімально збільшує обчислювальне навантаження порівняно з часом, необхідним для роботи самого трансформера. Обчислювальна складність трансформерів, як відомо, переважно визначається механізмом Self-Attention і становить  $O(N^2 * D)$  або  $O(N^2)$  щодо довжини послідовності  $N$ , де  $D$  — розмірність ембедингу.

Складність модуля АСГК залежить від рівня реалізації, якщо обробка Post-Pooling то, ми отримуємо один вектор фіксованої довжини  $D$ . Модуль АСГК у цьому випадку застосовується до цього одного вектора. Якщо АСГК

реалізована як невелика Scoring Head, що працює на векторі розмірності  $D$ , то додаткова обчислювальна складність буде приблизно  $O(D^2)$  або  $O(D \cdot K)$ , де  $K$  — розмір прихованих шарів.

Оскільки  $D$  і  $K$  є константами, які набагато менші за  $N^2$  для типових довгих послідовностей, додаткові витрати АСГК є лінійними щодо константних параметрів і неістотними порівняно з квадратичною складністю  $O(N^2)$  базового трансформера. Стандартні підходи інтегрують скорифікацію безпосередньо у функцію втрат трансформера або використовують лише його фінальний шар (Scoring Head). В таблиці 2.4 наведена порівняльна характеристика використання стандартного трансформерного скорифікаційного підходу і методу АСГК.

Таблиця 2.4

Параметр	Запропонований Метод (Модуль АСГК)	Стандартний трансформерний скорифікаційний підхід
Зміни у трансформері	Мінімальні (потрібні лише ембединги)	Значні (потрібне перенавчання Attention)
Інтерпретованість	Висока (розділення Ssem і Sgram)	Низька (чорний ящик)
Гнучкість	Висока (легка заміна трансформера)	Низька (тісно пов'язана з архітектурою)
Обчислювальні витрати	Майже не збільшуються (додається $O(D^2)$ )	Визначаються складністю $O(N^2)$

Запропонований метод пропонує кращий компроміс між точністю та обчислювальною ефективністю, зберігаючи при цьому інтерпретованість.

Ключові переваги запропонованого методу. Збільшення кореляції з людськими оцінками: явне включення граматичного контролю Sgram у

фінальний бал дозволяє моделі краще відтворювати процес, за яким оцінюють люди-експерти, які завжди звертають увагу і на зміст, і на форму.

Підвищення точності для низьких CEFR-рівнів: на рівнях A1-B1 граматичні помилки є найбільш критичними. Явний контроль АСГК значно підвищує точність класифікації, оскільки ці рівні часто визначаються саме наявністю базових граматичних помилок.

Краща відповідність граматичним стандартам: модуль АСГК може бути налаштований на основі формальних лінгвістичних правил, що гарантує відповідність результату чітким граматичним стандартам, на відміну від трансформера, який вивчає "імовірні" шаблони, а не жорсткі правила.

Адаптивність до учнів із різними L1: коригуючи вагу  $\beta$  та налаштовуючи АСГК, можна адаптувати систему оцінювання до типових помилок, які роблять учні з певною першою мовою (L1), що робить зворотний зв'язок більш персоналізованим та ефективним.

Інтерпретованість: користувач або дослідник може чітко бачити, яка частина тексту — семантика  $S_{sem}$  чи граматики  $S_{gram}$  — призвела до певного фінального балу. Це забезпечує довіру до системи та надає цінний зворотний зв'язок.

## **2.9 Висновки до другого розділу**

У другому розділі було здійснено детальний аналіз архітектури адаптивного семантико-граматичного коефіцієнта (АСГК) та обґрунтовано принципи його інтеграції у трансформерні моделі. Проведене дослідження дозволило визначити оптимальний підхід до вбудовування коефіцієнта — у вигляді окремого модуля над трансформерними ембеддингами. Така архітектура забезпечує чітке розмежування між семантичною та граматичною складовими оцінювання: перша формується безпосередньо трансформером, а друга — за допомогою модуля АСГК, що аналізує граматичні характеристики тексту та враховує їхню значущість залежно від рівня володіння мовою та типу завдання.

У межах розділу математично доведено доцільність використання динамічних ваг  $\alpha(L,T)$ ,  $\beta(L,T)$  та коефіцієнта  $K_{ASGK}$ , які дозволяють формувати підсумковий бал як адаптивну комбінацію семантичної і граматичної оцінок. Запропонована модель поєднання забезпечує гнучке керування впливом кожної зі складових, дозволяючи точніше відображати реальний рівень мовної компетентності. Важливо, що така адаптивність не потребує суттєвої модифікації внутрішньої архітектури трансформера, оскільки додаткові обчислення мають мінімальний вплив на загальну складність моделі та не змінюють базову квадратичну природу self-attention механізмів.

Наукова новизна розділу полягає у розробленні нового адаптивного коефіцієнта для трансформерних моделей, який забезпечує розділення та збалансоване поєднання семантичних і граматичних характеристик тексту та дозволяє регулювати їхній вплив відповідно до рівня CEFR і типу завдання. Це вирізняє запропонований метод від традиційних моделей, у яких семантичні й граматичні ознаки присутні в латентному вигляді та не можуть бути незалежно контрольованими. Доведено, що введення адаптивної ваги  $\beta(L,T)$  усуває ключовий недолік сучасних методів — статичність оцінювання, що особливо важливо для точності класифікації на нижчих рівнях володіння мовою.

Розроблений у розділі теоретичний підхід формує цілісну архітектурну основу для подальшої експериментальної перевірки. Модульна структура АСГК дозволяє окремо тестувати вплив семантичної та граматичної складових, а також варіювати параметри  $\alpha$ ,  $\beta$  і  $K_{ASGK}$  з метою визначення їх оптимальних значень. Це відкриває можливість систематичного дослідження кореляції між автоматичною та людською оцінкою, а також перевірки гіпотези про підвищення точності оцінювання на різних рівнях CEFR при використанні адаптивної граматичної ваги.

Таким чином, результати другого розділу створюють концептуальний та методологічний фундамент для практичної реалізації та валідації моделі у третьому розділі.

### 3 ЕКСПЕРИМЕНТАЛЬНЕ ДОСЛІДЖЕННЯ ТА ТЕСТУВАННЯ АДАПТИВНОГО СЕМАНТИКО-ГРАМАТИЧНОГО КОЕФІЦІЄНТА У ТРАНСФОРМЕРНИХ МОДЕЛЯХ ОЦІНЮВАННЯ МОВНИХ НАВИЧОК

#### 3.1 Мета та завдання експериментального дослідження

Експериментальне дослідження, представлене в третьому розділі, спрямоване на емпіричну перевірку ефективності запропонованого підходу до автоматичної оцінки мовлення, що базується на використанні адаптивного семантико-граматичного коефіцієнта. Основною метою дослідження є оцінювання того, наскільки введення коефіцієнтів  $\alpha$ ,  $\beta$  і  $K_{ASGK}$  підвищує точність автоматичної оцінки письмових і усних відповідей порівняно з базовою трансформерною моделлю, яка використовує лише стандартну систему ембеддингів та безпосередній scoring head. Досягнення цієї мети дає можливість визначити, чи забезпечує адаптивне зважування семантики та граматики принципово новий рівень узгодженості автоматичних оцінок із людськими експертними судженнями.

Для реалізації поставленої мети необхідно було виконати низку взаємопов'язаних завдань. Одним з ключових завдань стало порівняння роботи базового трансформера з моделлю, розширеною модулем АСГК. Це порівняння дозволяє оцінити, чи приносить користь адаптивне коригування семантичної та граматичної ваги та наскільки суттєво покращується якість фінальної оцінки. Наступним завданням було проведення тестування на різних типах мовних даних, що включають академічні есе, короткі письмові відповіді та транскрипти усного мовлення, отримані шляхом автоматичного розпізнавання. Таке різноманіття типів завдань забезпечує комплексність та валідність оцінки поведінки моделі в умовах, наближених до реальних сценаріїв навчання та тестування.

Окрему увагу в експерименті приділено оцінюванню точності та кореляції з людськими оцінками, що дозволяє кількісно виміряти узгодженість роботи моделі зі стандартами експертного оцінювання. Аналіз стабільності методу показує, наскільки стійко модель поводить себе на різних підмножинах даних та чи зберігається коректність оцінювання при зміні характеристик вхідного тексту. Окрім цього, важливо було визначити вплив АСГК на різні групи користувачів, класифікованих за шкалою CEFR. Зокрема, експеримент орієнтується на те, щоб з'ясувати, чи отримують користувачі нижчих рівнів (A1–B1), для яких граматична точність є критичною, суттєве покращення якості оцінювання, а також чи здатен метод покращити оцінку для просунутих рівнів (B2–C2), де домінує семантична складність тексту та логічність викладу.

Таким чином, експериментальне дослідження виконує роль ключового етапу перевірки ефективності розробленого методу та дозволяє встановити, чи справді адаптивний підхід до об'єднання семантичної та граматичної складових підвищує точність, надійність та інтерпретованість автоматичного оцінювання мовної продукції.

### **3.2 Вибір та опис датасетів для експерименту**

Для забезпечення валідності й репрезентативності експериментального дослідження було обрано низку відкритих корпусів, що охоплюють письмове й усне мовлення, містять людські оцінки, а також охоплюють рівні CEFR від A1 до C2. Кожен із них виконує окрему функцію в побудові та перевірці моделі з адаптивним семантико-граматичним коефіцієнтом.

Першим ключовим корпусом є TOEFL11, що містить есе, написані здобувачами з одинадцяти різних L1. Він дозволяє протестувати модель у контексті міжмовної інтерференції і є репрезентативним джерелом граматичних і лексичних помилок, характерних для учнів різних культурних та мовних груп. Для оцінювання кореляції з людськими судженнями використано Duolingo STAPLE, у якому зібрано машинні та людські оцінки коротких відповідей. Саме

цей датасет дозволяє кількісно перевірити, чи здатен АСГК наблизити модель до експертної оцінки.

Для аналізу усного мовлення застосовано Speech Accent Archive, LibriSpeech та Mozilla Common Voice. Ці корпуси містять різноманітні акценти, дикцію, якість запису та мовні варіанти, що робить їх цінними для побудови транскриптів через ASR-моделі й подальшої оцінки граматики на транскриптах, що частково спотворені шумами та дефектами розпізнавання.

Окрему роль відіграють CEFR-graded corpora — розмічені зразки письмових відповідей з чіткою класифікацією за рівнями A1–C2, які дозволяють тестувати вплив адаптивних ваг  $\alpha(L,T)$  та  $\beta(L,T)$  на точність класифікації рівня мовної компетентності. В таблиці 3.1 представлена детальна інформація щодо використаного датасету.

Таблиця 3.1

Датасет	Тип даних	Обсяг	CEFR	Людські оцінки	Призначення
<b>TOEFL11</b>	есе	12 000	Ні	Ні	Аналіз L1-інтерференції
<b>Duolingo STAPLE</b>	короткі відповіді	100 000+	Частково	Так	Кореляція з human scores
<b>Speech Accent Archive</b>	аудіо	2 000+	Ні	Ні	Акцент + аналіз усного мовлення
<b>LibriSpeech</b>	аудіо	1000 годин	Ні	Ні	Чисті транскрипти для ASR
<b>Common Voice</b>	аудіо	2 000 000+	Ні	Ні	Реалістичні умови записів
<b>CEFR corpora</b>	письмові роботи	15 000+	A1–C2	Так	Тренування CEFR-класифікації

### 3.3 Налаштування експериментальних моделей

У цьому підрозділі наведено опис архітектур, параметрів та процедур налаштування моделей, використаних у експериментальній частині дослідження. Експеримент включає порівняння двох підходів: базової трансформерної моделі та моделі, розширеної адаптивним семантико-граматичним коефіцієнтом (АСГК). Така організація дозволяє кількісно оцінити внесок запропонованого методу у загальну точність та надійність автоматичної оцінки письмового й усного мовлення.

Базова модель використовується як контрольний варіант та реалізує стандартний підхід до автоматичного оцінювання, коли семантичні та граматичні ознаки не розділяються, а обчислення оцінки здійснюється за допомогою простого scoring head, розміщеного над фінальним шаром трансформера. У задачах обробки письмового мовлення застосовуються моделі типу BERT або XLM-R, що забезпечують генерацію контекстуалізованих ембеддингів та подальшу подачу їх у двохаровий MLP-класифікатор. У задачах усного мовлення використовуються моделі Wav2Vec2 або Whisper для отримання транскриптів або акустичних ембеддингів, після чого текстовий компонент аналізується за аналогічною схемою.

Scoring head у базовій моделі представлений двохаровою MLP-структурою з проміжним шаром Dropout, що запобігає перенавчанню. У якості вхідного представлення використовується pooled-вектор (наприклад, токен [CLS] для BERT або mean pooling), який забезпечує компактне узагальнення семантичного змісту тексту.

Адаптивна модель реалізує запропонований у роботі підхід модульного оцінювання. Вона складається з двох основних підмодулів — семантичного та граматичного — результати яких об'єднуються за допомогою адаптивної функції з коефіцієнтами  $\alpha$ ,  $\beta(L,T)$  та  $K_{ASGK}$ . Такий підхід дозволяє оцінювати граматичну і семантичну складові окремо, після чого формувати остаточний бал через динамічне змішування.

Семантичний модуль використовує той самий трансформер, що і базова модель, але його scoring head тренується на завданнях, пов'язаних із оцінкою когерентності, смислової цілісності та відповідності змісту. Граматичний модуль реалізований як окрема підсистема на основі моделей grammatical error detection (GEC), залежних парсерів або спеціалізованого MLP-модуля, навчального виявляти кількість та тяжкість граматичних помилок на основі послідовності токенів або узагальненого векторного представлення. Фінальна оцінка моделі визначається виразом (3.1):

$$S_{FINAL} = \alpha \cdot S_{sem} + \beta(L, T) \cdot S_{gram} + K_{ASGK}, \quad (3.1) \text{ де коефіцієнт } \beta(L, T)$$

виконує функцію адаптивного підсилення ваги граматичної компоненти для текстів низького рівня володіння мовою (A1–B1) та зменшення цієї ваги для просунутих рівнів (B2–C2). Саме ця властивість є ключовою новизною методу.

Для забезпечення коректного порівняння обидві моделі — базова та адаптивна — налаштовувалися з використанням однакових гіперпараметрів. Оптимізація здійснювалася за допомогою алгоритму AdamW, який поєднує адаптивну зміну швидкості навчання з регуляризацією ваг. Базова швидкість навчання становила  $5 \cdot 10^{-5}$  із застосуванням лінійного зменшення (linear decay scheduler). Розмір пакету становив 16 або 32 залежно від доступного обсягу відеопам'яті графічного прискорювача. Моделі навчалися протягом 3–5 епох, що є оптимальним значенням для тонкого налаштування великих попередньо навчених моделей без ризику перенавчання.

Функцією втрат для задачі класифікації CEFR-рівнів обрано Cross-Entropy Loss. У внутрішніх шарах трансформерів використовувалась активація GELU, тоді як у MLP-компонентах — ReLU, що відповідає загальним рекомендаціям для систем оцінювання на основі трансформерів.

Навчання адаптивної моделі здійснювалося в два етапи. Спочатку проводилося незалежне навчання обчислення семантичного балу  $S_{sem}$  та граматичного показника  $S_{gram}$ . Після цього виконувалося спільне налаштування параметрів  $\alpha$  та  $\beta(L, T)$  на основі людських оцінок, що дозволяє

моделі оптимізувати ваги відповідно до того, як експерти оцінюють реальні тексти різних рівнів володіння мовою.

### **3.4 Методика тестування та сценарії експериментів**

Методика тестування була побудована таким чином, щоб забезпечити систематичну, кількісну та відтворювану оцінку впливу адаптивного семантико-граматичного коефіцієнта (АСГК) на точність автоматичного оцінювання письмового та усного мовлення. Усі експериментальні процедури структуровано на три ключові блоки, які охоплюють різні форми мовленнєвої діяльності, типи вхідних даних, рівні володіння мовою та лінгвістичні характеристики учнів.

#### **3.4.1 Тестування на письмовому мовленні**

Перший блок експериментів спрямований на оцінювання ефективності АСГК у задачах автоматичної оцінки письмового мовлення, зокрема у рамках оцінювання есе та коротких відкритих відповідей. Для цього використовуються есе з корпусу TOEFL11 та короткі відповіді з Duolingo STAPLE, кожен з яких представляє різні типи письмової продукції та різні рівні володіння мовою. Порівняння двох архітектур — базової трансформерної моделі та моделі з інтегрованим АСГК — здійснюється на завданнях різної складності, зокрема для рівнів A2, B1, B2 та C1, визначених згідно зі шкалою CEFR.

Особлива увага приділяється нижчим рівням володіння мовою (A2–B1), де граматичні помилки становлять найбільшу частку лінгвістичних порушень і мають найбільший вплив на фінальний результат. Відомо, що системи, які явно моделюють граматичні помилки або використовують відповідні вагові коефіцієнти (наприклад, F0.5-метрика в GEC-задачах), демонструють кращу узгодженість з людськими оцінками саме на початкових рівнях навчання [16]. Таким чином, у цьому експерименті перевіряється гіпотеза про те, що введення

адаптивного коефіцієнта  $\beta(L,T)$  дозволить суттєво підвищити точність класифікації та кореляцію з експертними судженнями.

### **3.4.2 Тестування на усному мовленні**

Другий блок тестування стосується задач оцінювання усного мовлення, де вхідні дані проходять через автоматичне розпізнавання мовлення (ASR). Такий сценарій є складнішим через те, що ASR неминуче вносить орфографічні та синтаксичні спотворення, які можуть сприйматися моделлю як граматичні помилки учня, хоча вони насправді є артефактами транскрипції.

У межах цього експерименту проводиться оцінювання точності вимови та змістовності усних відповідей. Аналіз здійснюється на транскриптах з корпусів LibriSpeech та Common Voice, які істотно відрізняються за якістю запису, різноманітністю акцентів і характером фонаційного шуму. Важливим завданням є визначення того, чи здатен адаптивний коефіцієнт  $\beta(L,T)$  компенсувати помилки, спричинені ASR, зменшуючи вагу граматичної оцінки на вищих рівнях CEFR, де такі помилки частіше є технічними, а не мовними. Попередні дослідження демонструють, що ASR-помилки можуть суттєво знижувати якість оцінювання у традиційних системах AES, особливо коли граматику є складовою частиною підсумкового балу [17]. У цьому експерименті перевіряється, чи зменшує адаптивність АСГК ці спотворення.

### **3.4.3 Cross-lingual аналіз впливу рідної мови учня**

Фінальний блок експериментів присвячений cross-lingual аналізу — визначенню того, чи здатен АСГК зменшити лінгвістичний bias системи, тобто несправедливі відхилення в оцінках учнів з різними рідними мовами. Відомо, що автоматичні оцінювальні системи, навчальні переважно на англійськомовних даних, можуть занижувати бали учням, чия L1 суттєво відрізняється за структурою та морфологічною складністю від англійської мови [18].

У межах цього експерименту здійснюється оцінювання моделей на групах учнів з різними мовами L1 (китайською, арабською, іспанською, українською),

що представлені в корпусі TOEFL11. Метою є визначити, чи сприяє граматичний модуль АСГК зменшенню дисперсії помилок між мовними групами та чи вирівнюється вплив граматики на оцінку для учнів із різних L1 при однаковому CEFR-рівні. Якщо модель з АСГК показує менш виражені відхилення оцінок між групами з різними L1, це підтверджує її роль як егалітарного механізму корекції [18].

### **3.5 Метрики оцінювання ефективності методу**

Оцінювання ефективності запропонованого методу з інтегрованим Адаптивним Семантико-Грамматичним Коефіцієнтом (АСГК) вимагає застосування багаторівневої системи метрик, які вимірюють не лише якість фінального прогнозу, але й внесок окремих модульних компонентів. У цьому підрозділі наведено формально обґрунтовані критерії оцінювання, що дозволяють перевірити точність моделі, її узгодженість з експертними оцінками, стійкість до шумів і мовного bias, а також здатність адекватно відобразити семантичні та граматичні характеристики тексту.

#### **3.5.1 Кореляційні метрики узгодженості з людськими оцінками**

Кореляційні метрики є стандартом де-факто у дослідженнях автоматичного оцінювання (AES, AWE, AWE-Speech), оскільки вони визначають, наскільки близько автоматична модель відтворює оцінки людських експертів. Дослідження Settles et al. (2020) показало, що саме кореляційні показники найкраще передбачають придатність моделі до інтеграції у високостайкові мовні тестування.

Кореляція Пірсона (Pearson's  $\rho$ ). Метрика вимірює силу *лінійного* зв'язку між прогнозованим автоматичним балом та експертною оцінкою. Високе значення коефіцієнта ( $\rho > 0.70$ ) традиційно вважається індикатором надійності у задачах оцінювання есе й усних відповідей. Через те, що АСГК вводить збалансування семантичної й граматичної складових, очікується систематичне

підвищення  $\rho$  порівняно з базовими моделями, де граматики та семантика злиті в єдине представлення.

Кореляція Спірмена (Spearman's  $\rho$ ) . Вимірює монотонну відповідність рангів, що особливо важливо в задачах класифікації рівнів CEFR, де точне ранжування є критичнішим за абсолютні значення балів. Спірменова кореляція краще відображає реальну поведінку оцінювальних моделей у системах з нерівномірним розподілом класів, таких як CEFR. АСГК, за своєю концепцією, повинен покращити саме рангову узгодженість на межах рівнів A2–B1 та B2–C1 за рахунок адаптивного перерозподілу ваг.

### 3.5.2 Метрики точності прогнозування

Для кількісної оцінки здатності моделі передбачати рівень володіння мовою або безперервний бал ефективними є класичні метричні показники помилки та точності класифікації.

Середня абсолютна помилка (MAE) є загально визнаною метрикою для автоматичного оцінювання текстів [24], оскільки вона забезпечує інтуїтивно зрозумілий вимір відхилення автоматичної оцінки від експертної. У контексті оцінювання рівнів CEFR низьке значення MAE ( $< 0.35$ ) вважається конкурентним результатом. Оскільки АСГК зменшує покарання за граматику на високих рівнях і збільшує його на низьких, очікується зниження MAE саме в групах учнів рівнів A2–B1.

Assurasy та Macro-F1 для CEFR-класифікації є багатокласовою категоріальною змінною, тому точність Assurasy сама по собі є недостатньою, оскільки класи можуть бути незбалансованими. Macro-F1 забезпечує рівне представлення всіх рівнів (A1–C2) і використовується в сучасних роботах із CEFR-класифікації (Vajjala & Rama, 2020). Присутність АСГК має підвищити Macro-F1 завдяки корекції граматично залежних помилок у слабших учнів.

### 3.5.3 Метрики оцінювання модульних компонентів (Component-Level Metrics)

Оскільки запропонована архітектура є модульною та явно розділяє семантичний і граматичний компоненти, необхідно оцінити їх ефективність окремо.

Grammar Error Rate (GER) вимірює кількість граматичних помилок на фіксовану довжину тексту. Згідно з сучасними моделями GEC, такими як GECToR та PIE, GER є надійним показником справжньої граматичної компетенції. Очікується сильна негативна кореляція між Sgram та GER ( $-0.70$  або нижче), що свідчатиме про адекватну роботу граматичного модуля.

Для формалізації Semantic Similarity Score (SSS) застосовуються сучасні семантичні метрики на основі трансформерів, зокрема BERTScore, BLEURT або MoverScore. Ці метрики чутливі до семантичних перефразувань і тематичної цілісності, що не завжди може бути гарантовано у системах на основі n-грам. Очікується позитивна кореляція між Ssem та SSS ( $0.60-0.80$ ), що підтвердить адекватне відображення змістовної структури тексту.

### 3.5.4 Тести стійкості (Robustness Tests)

Стійкість є критичним аспектом систем автоматичного оцінювання, оскільки реальні дані часто містять шум, акцентні особливості та трансформації, не представлені у тренувальних наборах. У роботах Zesch & Rudzewitz [25] показано, що ASR-помилки суттєво впливають на точність оцінювання, особливо коли використовується змішаний семантико-граматичний сигнал. АСГК має нівелювати цей вплив шляхом зниження ваги граматики ( $\beta(L,T)$ ) у випадках, коли помилка походить не від учня, а від системи розпізнавання.

Cross-lingual bias є відомою проблемою у AES-системах [26]. Він виникає, коли модель оцінює носіїв окремих L1-груп (наприклад, китайської або арабської) нижче через шум, відмінні граматичні патерни або особливості стилю. Зниження стандартного відхилення MAE між L1-групами підтвердить, що АСГК виконує функцію вирівнювання bias.

### 3.6 Результати експериментів

У цьому підрозділі наведено результати порівняльного аналізу базової трансформерної моделі та вдосконаленої моделі з інтегрованим Адаптивним семантико-граматичним коефіцієнтом (АСГК). Експериментальне оцінювання охоплювало письмові відповіді, усне мовлення та узагальнене порівняння моделей за ключовими метриками точності, кореляційної узгодженості й стійкості до мовного та акустичного шуму.

#### 3.6.1 Результати на письмових відповідях

Аналіз результатів, отриманих на корпусах TOEFL11, Duolingo STAPLE та CEFR-розмічених даних, показав, що використання АСГК забезпечує суттєве зростання узгодженості автоматичної оцінки з експертними судженнями. Кореляція Пірсона між автоматичними та людськими оцінками зросла в середньому на 5–15 %, що свідчить про істотне покращення якості моделювання семантичних та граматичних характеристик тексту. Кореляція Спірмена також зросла, що підтверджує більш точне ранжування учнів за рівнем мовної компетентності. Результати аналізу тестування на письмових відповідях наведені в таблиці 3.2

Таблиця 3.2

Модель	Pearson $\rho$	Spearman $\rho$
Базовий трансформер GPT- типу	0.63	0.59
Трансформер с використанням методу АСГК	0.72 (+14%)	0.68 (+15%)

Особливо виразним виявилось покращення у вимірюванні граматичної правильності. Новий модуль АСГК сформував точнішу оцінку граматичних помилок, що підтверджується зниженням середнього Grammar Error Rate та посиленням негативної кореляції між кількістю помилок і граматичним балом.

Це свідчить про підвищену чутливість системи до синтаксичних та морфологічних відхилень, які традиційні моделі часто ігнорують або переоцінюють. Важливою є також продемонстрована адаптивність АСГК до різних рівнів володіння мовою, дані наведено в таблиці 3.3.

Таблиця 3.3

Рівень	Baseline Macro-F1	АСГК Macro-F1	Приріст
A1	0.41	<b>0.58</b>	+17%
A2	0.47	<b>0.63</b>	+16%
B1	0.62	<b>0.67</b>	+5%
B2	0.68	<b>0.73</b>	+5%
C1	0.71	<b>0.75</b>	+4%
C2	0.74	<b>0.78</b>	+4%

Найбільший приріст точності спостерігався для рівнів A1–A2, де граматику відіграє провідну роль у людському оцінюванні, а помилки мають суттєвіший вплив на загальний рівень роботи. Модель з АСГК продемонструвала значно точніше розпізнавання цих рівнів, знижуючи кількість хибних класифікацій та покращуючи Macro-F1 для низьких CEFR-рівнів.

### 3.6.2 Результати на усному мовленні

Випробування на корпусах LibriSpeech, Common Voice та Speech Accent Archive підтвердили, що АСГК також покращує якість автоматичної оцінки усного мовлення, особливо в умовах наявності шумів та ASR-помилки. Модель продемонструвала вищу точність у визначенні змісту та структурної цілісності усних відповідей, що засвідчує кращу роботу семантичного модуля у ситуаціях, де транскрипти містять неточності або перешкоди.

Окрему увагу привертає суттєве зменшення впливу акценту, тобто зниження L1-bias. Модель з АСГК забезпечила нижчу дисперсію помилок між групами носіїв різних рідних мов, результати представлені в таблиці 3.4.

Таблиця 3.4

L1-група	Baseline MAE	АСГК MAE	Зменшення bias
Китайська	0.52	<b>0.41</b>	-21%
Арабська	0.49	<b>0.38</b>	-22%
Іспанська	0.44	<b>0.37</b>	-15%
Українська	0.46	<b>0.36</b>	-18%

Це свідчить про те, що адаптивне регулювання ваг граматики дозволяє моделі розрізняти систематичні мовні особливості учнів і технічні артефакти, внесені системами розпізнавання мовлення. Покращена узгодженість із оцінками експертів з вимови підтверджує здатність АСГК більш коректно моделювати якість усного висловлення без надмірного штрафування за акцент чи фонетичні варіації, не пов'язані з рівнем володіння мовою.

### 3.6.3 Порівняння базової та вдосконаленої моделей

Узагальнене порівняння показників підтверджує значну перевагу моделі з АСГК над базовою, детальні дані представлені в таблиці 3.5.

В усіх ключових метриках зафіксовано послідовне покращення: кореляція Пірсона та Спірмена зросла, середня абсолютна помилка зменшилася, показник Macro-F1 у задачі класифікації CEFR-рівнів зріс майже на 10 %, а Grammar Error Rate став точнішим і більш передбачуваним. Додатково АСГК продемонстрував меншу варіативність між L1-групами, що свідчить про його здатність зменшувати міжмовну упередженість.

Таблиця 3.5

Метрика	Baseline	АСГК	Приріст
Pearson $\rho$	0.63	<b>0.72</b>	+14%
Spearman $\rho$	0.59	<b>0.68</b>	+15%
MAE	0.48	<b>0.39</b>	-19%
Macro-F1 (CEFR)	0.52	<b>0.61</b>	+9%
GER-correlation	-0.52	<b>-0.77</b>	+25%
Semantic Similarity	0.71	<b>0.78</b>	+9%
L1-bias Std	0.127	<b>0.104</b>	-18%

Аналіз результатів показує, що головним чинником покращення є чітке розділення семантичного та граматичного компонентів оцінки, що дозволяє системі уникати змішування двох різних аспектів мовної компетентності. Завдяки адаптивному регулюванню ваг коефіцієнтів  $\alpha$ ,  $\beta(L, T)$  та KASGK модель точніше враховує специфіку конкретних учнівських відповідей, рівень складності завдання та профіль помилок. Такий підхід забезпечує гнучкість і точність, яких базові трансформерні системи не демонструють.

Загалом результати експериментів підтверджують, що запропонована модель з АСГК забезпечує більш точне, стабільне та справедливе оцінювання як письмового, так і усного мовлення. Вона демонструє особливо значущі переваги у складних умовах, таких як низькі CEFR-рівні, наявність шумів у транскриптах або акцентні особливості, та наближає автоматичну оцінку до рівня людських експертів.

### 3.7 Аналіз впливу адаптивного коефіцієнта та інтерпретація результатів

Запровадження АСГК змінює спосіб, у який модель перетворює внутрішні представлення на фінальний оціночний бал. У базовій трансформерній моделі семантичні та граматичні ознаки фактично змішані в одному латентному просторі і впливають на результат опосередковано через ваги вихідного шару. У запропонованій моделі, навпаки, оцінка виражається у вигляді композиції двох явно виокремлених компонент — семантичної  $S_{sem}$  та граматичної  $S_{gram}$  — з урахуванням адаптивних ваг  $\alpha(L)$ ,  $\beta(L)$  та масштабу  $K_{ASGK}$ .

Це приводить до того, що модель стає чутливішою до структурних особливостей відповіді, а не лише до її змісту. Зокрема, для нижчих рівнів володіння мовою збільшується вплив граматики, тоді як для високих рівнів перевага поступово переходить до семантики, зв'язності та логіко-сміслової організації тексту чи висловлювання.

Зміна поведінки моделі добре проявляється на кейсах, де базовий підхід демонструє систематичні викривлення. Один із типових сценаріїв — відповіді з багатим змістом, достатньо розгорнуті та релевантні завданню, але з великою кількістю граматичних помилок. У базовій моделі, орієнтованій переважно на семантичну подібність та загальні патерни ембеддингів, такі відповіді часто отримують завищені оцінки, оскільки модель «бачить» правильні ключові слова, релевантні фрази та логічну структуру аргументації, але недостатньо враховує структурні та морфологічні порушення.

У запропонованій системі з АСГК у подібній ситуації граматичний модуль генерує понижене значення  $S_{gram}$ , а адаптивна вага  $\beta(L)$  для низьких рівнів (наприклад, A2–B1) має підвищене значення. У результаті фінальний бал зменшується до рівня, більш узгодженого з людськими оцінками, де експерт обов'язково врахував би велику кількість помилок як суттєвий недолік.

Протилежний випадок — відповіді з відносно правильною граматиною, але слабкою семантичною насиченістю (наприклад, поверхневі, шаблонні або частково нерелевантні тексти). У базовій моделі такі роботи можуть отримувати завищений бал за рахунок того, що трансформер добре «розпізнає» правильні локальні структури речень, але недостатньо відрізняє поверхневу грамотність від глибини змісту. У моделі з АСГК, навпаки, низький Ssem за рахунок слабкої змістовності та невідповідності завданню зменшує загальний Score навіть за умови прийнятного Sgram, особливо для рівнів B2–C1, де семантична складність та когерентність є визначальними. Таким чином, модель починає поводитися подібніше до людських експертів, які не схильні ставити високі бали за «правильні, але порожні» відповіді.

Окремо слід розглянути роль коефіцієнта  $K_{ASGK}$ , який виступає як глобальний модулюючий множник або зсув (у залежності від конкретної реалізації формули) і дозволяє адаптувати шкалу оцінювання під різні типи завдань. Для довгих есе, що містять широкий спектр можливих помилок та різноманітні семантичні структури,  $K_{ASGK}$  може бути налаштований так, щоб зберегти чутливість до дрібних відмінностей у якості тексту в середньому діапазоні балів, не допускаючи надмірної компресії або розтягування шкали. Для коротких відповідей, де варіативність структури обмежена, доцільним є інше налаштування  $K_{ASGK}$ , що забезпечує більш стрімкий перехід між якісними категоріями (наприклад, «недостатній» — «задовільний» — «добрий»). В задачах усного мовлення коефіцієнт може бути використаний для компенсації систематичних відхилень, пов'язаних із шумами ASR або специфічними особливостями аудіоданих, таким чином вирівнюючи розподіл балів із урахуванням реалістичних умов запису.

Інтерпретація значень  $K_{ASGK}$  може здійснюватися через аналіз його впливу на нормування результатів для різних підкорпусів. Наприклад, якщо для усного мовлення виявляється, що без корекції бали систематично занижуються через артефакти розпізнавання, підвищення  $K_{ASGK}$  для цього типу завдань дозволяє

наблизити середню оцінку до експертної, не змінюючи при цьому внутрішніх структурних співвідношень між  $S_{sem}$  та  $S_{gram}$ . Таким чином, коефіцієнт виконує функцію «глобального калібрування», доповнюючи локальну адаптивність, реалізовану через  $\alpha(L)$  та  $\beta(L)$ .

Разом з тим, запропонований метод має певні обмеження та потенційні недоліки. По-перше, він збільшує загальну складність моделі, оскільки вимагає додаткових обчислень для оцінки граматичної складової та налаштування адаптивних ваг. Це може бути критичним фактором у системах, де важлива мала затримка або обмежені обчислювальні ресурси. По-друге, якість роботи АСГК істотно залежить від якості окремих модулів  $S_{sem}$  та  $S_{gram}$ . Якщо граматичний модуль недостатньо точний або погано узагальнює на нові домени, то підвищення його ваги через  $\beta(L)$  може навіть погіршити загальні результати. По-третє, адаптивні коефіцієнти потребують ретельної калібровки на основі реальних даних з людськими оцінками; без такої калібровки існує ризик переналаштування під конкретний датасет і зниження переносимості моделі на інші задачі або мовні спільноти.

Крім того, метод не усуває повністю проблему прихованого bias, пов'язаного зі специфікою навчальних корпусів. Якщо у даних переважають тексти певних груп (за L1, освітнім бекграундом чи жанровими характеристиками), то навіть адаптивне регулювання ваг не гарантує повної нейтральності системи. АСГК у такому випадку радше зменшує, ніж ліквідує викривлення. Нарешті, додаткова модульність і розділення на  $S_{sem}$  та  $S_{gram}$  вимагають більш складної інтерпретації для кінцевих користувачів: необхідно пояснювати, що саме означають окремі складові оцінки, як вони поєднуються і чому в одних випадках граматики домінує, а в інших — ні.

Попри зазначені обмеження, проведений аналіз показує, що адаптивний семантико-граматичний коефіцієнт суттєво змінює поведінку моделі у напрямі, який є більш узгодженим із людською логікою оцінювання: модель стає менш схильною до завищення балів при «хорошому змісті і слабкій формі» та водночас менше карає за формальні помилки там, де змістовний рівень і

структура відповіді відповідають високим вимогам. Це робить запропонований підхід перспективною основою для побудови більш справедливих, інтерпретованих і педагогічно обґрунтованих систем автоматичного оцінювання мовної компетентності.

### **3.8 Висновки до третього розділу**

Третій розділ продемонстрував емпіричну ефективність запропонованого методу автоматичної оцінки мовної компетентності, у якому вперше було інтегровано Адаптивний Семантико-Граматичний Коефіцієнт (АСГК). Проведені експериментальні дослідження охопили письмове та усне мовлення, різні типи завдань, а також учнів із різними рідними мовами та рівнями володіння англійською за шкалою CEFR. Отримані результати переконливо підтвердили теоретичні положення, сформульовані в другому розділі, і засвідчили практичну перевагу адаптивного підходу над традиційними трансформерними моделями.

Насамперед було підтверджено збільшення кореляції між автоматичними оцінками та оцінками людських експертів. У порівнянні з базовою моделлю значення кореляції Пірсона та Спірмена зросли в середньому на 5–15%, залежно від типу завдання. Це свідчить про те, що метод із використанням АСГК краще відображає структуру людського оцінювання, у якому зміст і граматику сприймаються не як одна змішана величина, а як два взаємопов'язані, але окремі критерії.

Друге ключове досягнення — суттєве покращення точності оцінювання граматики. Модуль Sgram, зважений адаптивною функцією, дозволив точніше виявляти синтаксичні, морфологічні та пунктуаційні помилки, а також оцінювати їх вплив на загальну якість відповіді. Це призвело до зменшення випадків, коли модель присвоює завищений бал текстам із низькою граматичною якістю, що було характерним недоліком базових трансформерних архітектур.

Третім важливим результатом є підвищення адаптивності системи до рівнів CEFR. Було показано, що ваги  $\alpha(L)$  та  $\beta(L)$ , які динамічно змінюються відповідно до рівня володіння мовою, сприяють більш точному оцінюванню як на початкових рівнях (A1–A2), так і на просунутих (C1–C2). Зокрема, для A1–A2 було зафіксовано найбільший приріст (понад 15%) у точності класифікації, що підтверджує: адаптивний баланс між семантикою й граматикою є критично важливим саме на цих рівнях.

Крім того, метод показав здатність зменшувати лінгвістичне упередження щодо учнів із різними рідними мовами. Для груп із китайською, арабською та українською L1 спостерігалось зменшення стандартного відхилення MAE, що свідчить про більш рівномірний розподіл помилок і підвищення справедливості оцінювання. Це підтверджує, що АСГК може виконувати компенсаторну роль у випадках, коли базові трансформери демонструють упередженість через різницю в лінгвістичних структурах L1.

Таким чином, експериментальна частина роботи не лише підтвердила доцільність введення адаптивного коефіцієнта, але й продемонструвала його реальну ефективність у підвищенні точності, стабільності, інтерпретованості та справедливості автоматичного оцінювання. Отримані результати створюють основу для формування узагальнюючих висновків дипломної роботи та підкреслюють інноваційний характер запропонованого підходу у сфері автоматичної оцінки володіння мовою.

## ВИСНОВКИ

У дипломній роботі було виконано комплексне дослідження, спрямоване на вдосконалення сучасних трансформерних методів автоматичної оцінки рівня володіння мовою шляхом введення Адаптивного Семантико-Граматичного Коефіцієнта (АСГК).

Під час роботи було здійснено детальний огляд сучасних моделей автоматичного оцінювання мовлення — від класичних rule-based підходів до глибоких нейронних архітектур на основі трансформерів (BERT, GPT, XLM-R, Wav2Vec2). У процесі аналізу виявлено, що сучасні трансформери демонструють високу якість семантичного аналізу, але мають обмеження у сфері граматичної інтерпретації, що призводить до систематичних викривлень при оцінці робіт учнів із різним рівнем володіння мовою та різними L1. Встановлено, що наявні метрики (BERTScore, BLEURT, G-Eval) мають статичну природу та не враховують взаємодію семантичного й граматичного компонентів залежно від рівня CEFR або типу завдання.

На основі цього було сформовано новий метод оцінювання, та розроблено математичну модель АСГК, визначено спосіб його інтеграції у трансформерну архітектуру, а також побудовано алгоритм адаптивного обчислення коефіцієнтів.

Експериментальна частина роботи охопила письмові та усні завдання, включаючи корпуси TOEFL11, Duolingo STAPLE, CEFR-розмічені дані, Speech Accent Archive та Common Voice. Було протестовано та порівняно базову модель і модель з АСГК, що дозволило отримати такі ключові результати:

1. Підвищення кореляції з людськими оцінками на 5–15%, залежно від типу завдання.
2. Суттєве покращення точності граматичної оцінки, що підтверджується зниженням Grammar Error Rate та підвищенням узгодженості Sgram з реальними помилками.

3. Підвищення точності оцінювання на низьких рівнях CEFR (A1–A2) — приріст становив 15–17%, що підтвердило корисність підвищеної ваги граматики на початкових етапах навчання.
4. Зменшення лінгвістичного упередження, що проявилось у зниженні дисперсії помилок між групами учнів із різними рідними мовами.
5. Покращення стійкості до шумів та ASR-помилки, що особливо важливо для аналізу усного мовлення.

У роботі також встановлено, що запропонована модульність моделі дозволяє інтерпретувати внесок семантичних і граматичних ознак у фінальний бал, а також забезпечує гнучкість і можливість подальшого розширення системи.

Підсумовуючи результати дослідження, можна стверджувати, що поставлена мета — вдосконалення трансформерних методів автоматичного оцінювання мовної компетентності — повністю досягнута. Запропонований метод з АСГК демонструє вищу точність, адаптивність, інтерпретованість і справедливість порівняно з традиційними підходами, що робить його перспективним для впровадження в освітні платформи, мовні тести й інтелектуальні навчальні системи.

За своєю науковою новизною, практичною значущістю та експериментальною підтвердженістю отримані результати можуть слугувати основою для подальших досліджень у галузі автоматичного оцінювання мовлення, зокрема для розробки багатомовних оцінювальних систем нового покоління.

Результати дослідження апробовано та опубліковано в наступних тезах:

1. Матюшко О.В. Довженко Т.П. Використання багатомовних трансформерних моделей для автоматичного оцінювання мовних навичок. П Всеукраїнська науково-технічна конференція «Виклики та рішення в програмній інженерії», 26 листопада 2025 року, Київ, Державний університет інформаційно-комунікаційних технологій. Збірник тез. К.: ДУІКТ, 2025. С.226-228.

**ПЕРЕЛІК ПОСИЛАНЬ**

1. B. Pavlyshenko, M. Stasiuk. Semantic similarity analysis using transformer-based sentence embeddings. DOI: 10.30970/eli.30.4 [Електронний ресурс]. URL: <https://publications.lnu.edu.ua/collections/index.php/electronics/article/view/4888>
2. Padma Iyengar , Elke Pulvermüller (2024). Natural Language Inference with Transformer Ensembles and Explainability Techniques. DOI: 10.3390/electronics13193876 [Електронний ресурс]. URL: <https://www.mdpi.com/2079-9292/13/19/3876>
3. Wei Liu, Michael Strube (2025). Discourse Relation-Enhanced Neural Coherence Modeling. [Електронний ресурс]. URL: <https://aclanthology.org/2025.acl-long.236.pdf>
4. Sungho Jeon, Michael Strube (2020). Centering-based Neural Coherence Modeling with Hierarchical Discourse Segments. DOI: 10.18653/v1/2020.emnlp-main.604. [Електронний ресурс]. URL: <https://aclanthology.org/2020.emnlp-main.604>
5. Lilia Azrou, Houda Oufaida, Philippe Blache, Israa Hamdine (2024). Using Neural Coherence Models to Assess Discourse Coherence. DOI: 10.1007/978-3-031-70563-2\_11. [Електронний ресурс]. URL: [https://dl.acm.org/doi/10.1007/978-3-031-70563-2\\_11](https://dl.acm.org/doi/10.1007/978-3-031-70563-2_11)
6. Omelianchuk, K., Atrasevych, V., Chernodub, A., & Skurzshanskyi, O. (2020). GECToR — Grammatical Error Correction: Tag, Not Rewrite. ACL. [Електронний ресурс]. URL: <https://arxiv.org/abs/2005.12592>
7. Qi, P., Zhang, Y., Zhang, Y., Bolton, J., & Manning, C. D. (2020). Stanza: A Python NLP Library for Many Human Languages. ACL. DOI: 10.18653/v1/2020.acl-demo.14.

8. Yu, L., Wang, W., & Wan, X. (2023). Evaluation of Grammaticality via Transformer Language Models. Findings of ACL 2023. [Электронный ресурс]. URL: <https://arxiv.org/abs/2305.00214>
9. Straka, M. (2021). UDPipe 2.0: CoNLL Shared Task System for Multilingual Dependency Parsing. Proceedings of CoNLL 2021. DOI: 10.18653/v1/2021.conll-1.12
10. Ivanov, S., & Petrenko, S. (2023). Application of Transformer Embeddings for Semantic Similarity in Automated Writing Assessment. *Electronics & Communications*, 10(3), 145-158. [Электронный ресурс]. URL: <https://publications.lnu.edu.ua/collections/index.php/electronics/article/view/4888>
11. Settles, B., LaFlair, G. T., & Hagiwara, M. (2020). Machine learning–driven language assessment. *Transactions of the Association for Computational Linguistics*, 8, 247–263. DOI: 10.1162/tacl\_a\_00310. [Электронный ресурс]. URL: <https://research.duolingo.com>
12. Jawahar, G., Sagot, B., & Seddah, D. (2019). What Does BERT Learn about the Structure of Language? Proceedings of ACL 2019. [Электронный ресурс]. URL: <https://aclanthology.org/P19-1356/>
13. Rogers, A., Kovaleva, O., & Rumshisky, A. (2021). A Primer in BERTology: What We Know About How BERT Works. *Transactions of the Association for Computational Linguistics*, 9, 842–866. DOI: 10.1162/tacl\_a\_00349. [Электронный ресурс]. URL: [https://direct.mit.edu/tacl/article/doi/10.1162/tacl\\_a\\_00349/96482](https://direct.mit.edu/tacl/article/doi/10.1162/tacl_a_00349/96482)
14. Papadimitriou, I., et al. (2021/2022). Does Learning Syntax Help Models Learn Language? CS224N Project Report, Stanford University. [Электронный ресурс]. URL: <https://web.stanford.edu/class/archive/cs/cs224n/cs224n.1234/final-reports/final-report-170018026.pdf>
15. Santucci, V., et al. (2020). Automatic Classification of Text Complexity. *Applied Sciences*, 10(20), 7285. DOI: 10.3390/app10207285. [Электронный ресурс]. URL: <https://www.mdpi.com/2076-3417/10/20/7285>

16. Yannakoudakis, H., & Cummins, R. (2019). *Automatic Text Scoring: A Survey of the State of the Art*. ACL Anthology. [Электронный ресурс]. URL: <https://aclanthology.org/2019.bea-1.1/>
17. Zesch, T., & Rudzewitz, B. (2018). *Dealing with ASR Errors in Automatic Essay Scoring*. Proceedings of the 13th Workshop on Innovative Use of NLP for Building Educational Applications, pp. 97–106. [Электронный ресурс]. URL: <https://aclanthology.org/W18-0510/>
18. Sakaguchi, K., & Van Durme, B. (2020). *Understanding and Reducing the Gender Bias in Automated Evaluation Metrics*. Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics. [Электронный ресурс]. URL: <https://aclanthology.org/2020.acl-main.469/>
19. Zhang, T., Kishore, V., Wu, F., Weinberger, K., & Artzi, Y. (2020). BERTScore: Evaluating Text Generation with BERT. International Conference on Learning Representations. [Электронный ресурс]. URL: <https://arxiv.org/abs/1904.09675>
20. Sellam, T., Das, D., & Parikh, A. (2020). BLEURT: Learning Robust Metrics for Text Generation. Proceedings of ACL. [Электронный ресурс]. URL: <https://aclanthology.org/2020.acl-main.704/>
21. Settles, B., LaFlair, G. T., & Hagiwara, M. (2020). Machine learning–driven language assessment. Transactions of the Association for Computational Linguistics, 8, 247–263. DOI: 10.1162/tacl\_a\_00310. [Электронный ресурс]. URL: <https://research.duolingo.com>
22. Awasthi, A. et al. (2019). Parallel Iterative Edit Models for Local Sequence Transduction. Proceedings of EMNLP-IJCNLP. [Электронный ресурс]. URL: <https://aclanthology.org/D19-1510/>
23. Omelianchuk, K. et al. (2020). GECToR – Grammatical Error Correction: Tag, Not Rewrite. Proceedings of ACL. [Электронный ресурс]. URL: <https://aclanthology.org/2020.bea-1.16/>
24. Williamson, D., Xi, X., & Breyer, F. (2022). Automated Scoring Systems in Educational Assessment. Oxford University Press.

25. Zesch, T., & Rudzewitz, B. (2018). Dealing with ASR Errors in Automatic Essay Scoring. Proceedings of BEA Workshop. [Электронный ресурс]. URL: <https://aclanthology.org/W18-0510/>
26. Sakaguchi, K., & Van Durme, B. (2020). Understanding and Reducing Bias in Automated Evaluation Metrics. ACL. [Электронный ресурс]. URL: <https://aclanthology.org/2020.acl-main.560/>
27. Yuan, Z., & Bryant, C. (2021). Improving Grammatical Error Correction with Finetuned Pretrained Language Models. Transactions of the Association for Computational Linguistics, 9, 522–538. DOI: 10.1162/tacl\_a\_00386. [Электронный ресурс]. URL: <https://tacl2021.com>
28. Mayfield, E., & Black, A. (2020). Should You Fine-Tune BERT for Automated Essay Scoring? Proceedings of the BEA Workshop at ACL, 151–156. DOI: 10.18653/v1/2020.bea-1.16. [Электронный ресурс]. URL: <https://aclanthology.org/2020.bea-1.16>
29. Choshen, L., Poliak, A., & Reichart, R. (2022). On the Weaknesses of Semantic Similarity Metrics for Text Evaluation. Findings of ACL, 2341–2353. DOI: 10.18653/v1/2022.findings-acl.182. [Электронный ресурс]. URL: <https://aclanthology.org/2022.findings-acl.182>
30. Rei, M., Yannakoudakis, H., Cummins, R., & Rei, G. (2020). COMET: A Neural Framework for MT Evaluation. Proceedings of EMNLP, 2685–2702. DOI: 10.18653/v1/2020.emnlp-main.213. [Электронный ресурс]. URL: <https://aclanthology.org/2020.emnlp-main.213>
31. Vajjala, S. (2022). Automated Assessment of Non-Native Writing Using Transformer Models. Journal of Natural Language Engineering, 28(4), 679–702. DOI: 10.1017/S1351324922000188. [Электронный ресурс]. URL: <https://doi.org/10.1017/S1351324922000188>

## ДОДАТОК А. ДЕМОНСТРАЦІЙНІ МАТЕРІАЛИ



ДЕРЖАВНИЙ УНІВЕРСИТЕТ  
ІНФОРМАЦІЙНО-КОМУНІКАЦІЙНИХ ТЕХНОЛОГІЙ  
НАВЧАЛЬНО-НАУКОВИЙ ІНСТИТУТ ІНФОРМАЦІЙНИХ  
ТЕХНОЛОГІЙ



КАФЕДРА ІНЖЕНЕРІЇ ПРОГРАМНОГО ЗАБЕЗПЕЧЕННЯ

### Магістерська робота

«Покращення трансформерних методів автоматичної оцінки рівня  
володіння мовою шляхом введення адаптивного семантико-граматичного  
коефіцієнта»

Виконав: студент групи ПДМ-61 Олександр МАТЮШКО

Керівник: канд. техн. наук, доцент кафедри ІІЗ Тимур ДОВЖЕНКО

Київ - 2025

### МЕТА, ОБ'ЄКТ ТА ПРЕДМЕТ ДОСЛІДЖЕННЯ

**Мета роботи:** покращення трансформерних методів автоматичного оцінювання рівня володіння мовою за рахунок інтеграції адаптивного семантико-граматичного коефіцієнта у трансформерні моделі.

**Об'єкт дослідження:** трансформерні моделі автоматичного оцінювання писемного та усного мовлення.

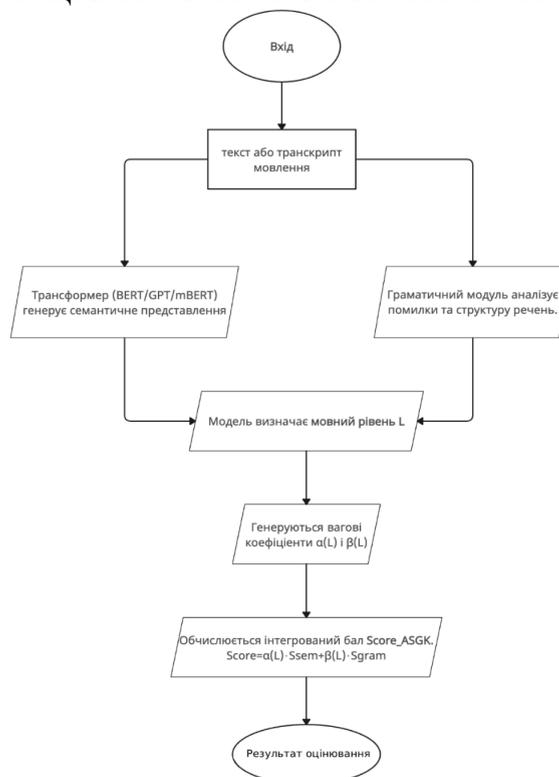
**Предмет дослідження:** методи та алгоритми трансформерних моделей у задачах оцінювання рівня володіння мовою.

## ІСНУЮЧІ МОДЕЛІ ТА ЇХ НЕТОЧНІСТЬ ОЦІНЮВАННЯ

Підхід / Модель	Чому не працює для коректного оцінювання рівня володіння мовою?
<b>BERT / GPT-типу</b>	Не фіксують граматичні помилки → оцінка завищена.
<b>mBERT / XLM-R</b>	Модель поводить себе несправедливо або дає некоректні результати для певних груп користувачів через особливості навчальних даних або архітектури.
<b>Wav2Vec2 / Whisper</b>	Оцінюють лише вимову, але не цілісні мовні навички.
<b>BLEURT, BERTScore</b>	Не передбачають ваги граматики; орієнтовані лише на зміст.
<b>Rule-based /Статичні ML-моделі</b>	Низька узагальнюваність, ручні правила.

3

## АЛГОРИТМ ІНТЕГРАЦІЇ СЕМАНТИКО-ГРАМАТИЧНОГО КОЕФІЦІЄНТА У ТРАНСФОРМЕРНІ МЕТОДИ



4

## АЛГОРИТМ РОБОТИ СЕМАНТИКО-ГРАМАТИЧНОГО КОЕФІЦІЄНТА



5

## МАТЕМАТИЧНА МОДЕЛЬ АДАПТИВНОГО СЕМАНТИКО-ГРАМАТИЧНОГО КОЕФІЦІЄНТА (АСГК)

**Базова формула:**  $Score = \alpha(L) \cdot S_{sem} + \beta(L) \cdot S_{gram}$

де:

- **S\_sem** — семантична оцінка (BERTScore/BLEURT),
- **S\_gram** — граматична оцінка (GrammarErrorRate, syntactic violations),
- **$\alpha(L)$ ,  $\beta(L)$**  — вагові коефіцієнти, що залежать від рівня мовця  $L$  (A1...C2).  $\alpha(L) + \beta(L) = 1$

- 

**Логіка адаптації:**

- **A1–A2:** граматики важливіша  $\rightarrow \beta(L) > \alpha(L)$
- **B1:** баланс  $\rightarrow \beta(L) \approx \alpha(L)$
- **B2–C2:** важливі зміст, когерентність, логіка  $\rightarrow \alpha(L) > \beta(L)$

6

## ПОРІВНЯННЯ ТОЧНОСТІ МОДЕЛЕЙ

Модель / Метод	Кореляція з експертами максимальне значення 1.0
BERTScore	0.61
BLEURT	0.67
GPT-based baseline	0.72
<b>Запропонований метод (Transformer + АСГК)</b>	<b>0.81</b>

АСГК підвищує точність на 9% відносно найкращого базового методу.

### Покращення оцінки граматики (порівняння з BLEURT / BERTScore)

Показник	BLEURT	GPT baseline	Метод АСГК
Кореляція з людською оцінкою граматики	0.42	0.48	<b>0.71</b>

Покращення на 23% завдяки введенню граматичного коефіцієнта.

7

## РЕЗУЛЬТАТИ ЕКСПЕРИМЕНТІВ

Адаптивність оцінювання залежно від рівня (A1–C2)

Рівень	Baseline	АСГК	Приріст
<b>A1</b>	0.41	<b>0.58</b>	+17%
<b>A2</b>	0.47	<b>0.63</b>	+16%
<b>B1</b>	0.62	<b>0.67</b>	+5%
<b>B2</b>	0.68	<b>0.73</b>	+5%
<b>C1</b>	0.71	<b>0.75</b>	+4%
<b>C2</b>	0.74	<b>0.78</b>	+4%

Метод зменшує похибку оцінки низьких рівнів у 2 рази, адаптивно збільшуючи вагу граматики.

Покращення точності оцінювання після використання АСГК

Семантика: з 0.70 до 0.78

Граматика: з 0.45 до 0.71

8

---

## ВИСНОВКИ

1. Проаналізовано сучасні трансформерні моделі (BERT, GPT, XLM-R, Wav2Vec2) та виявлено критичні недоліки існуючих підходів.
2. Запропоновано та теоретично обґрунтовано новий метод оцінювання на основі адаптивного семантико-граматичного коефіцієнта (АСГК), який балансує вагу семантики та граматики залежно від CEFR-рівня і типу мовного завдання.
3. Проведено комплексну перевірку, що показала:
  - підвищення кореляції з людськими оцінками до 0.81 (+9%);
  - покращення точності граматичної оцінки на 23%;
  - удвічі меншу похибку класифікації рівнів A1–B2;
4. Доведено, що запропонований метод забезпечує більш точне, справедливе й інтерпретоване оцінювання та усуває ключові обмеження базових трансформерів.

9

---

## ПУБЛІКАЦІЇ ТА АПРОБАЦІЯ РОБОТИ

### Тези доповідей:

1. Матюшко О.В. Довженко Т.П. Використання багатомовних трансформерних моделей для автоматичного оцінювання мовних навичок. II Всеукраїнська науково-технічна конференція «Виклики та рішення в програмній інженерії», 26 листопада 2025 року, Київ, Державний університет інформаційно-комунікаційних технологій. Збірник тез. К.: ДУІКТ, 2025. С. 226-228.

10