



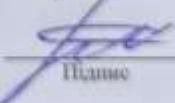
КВАЛІФІКАЦІЙНА РОБОТА БАКАЛАВРА

на тему Метод автоматизованого визначення семантичної подібності
цифрових текстів для виявлення запозичень

Галузь знань 12 – Інформаційні технології
Шифр і назва галузі знань
Спеціальність 122 – Комп'ютерні науки
Шифр і назва спеціальності
Освітня програма Комп'ютерні науки
Назва освітньої програми

Виконав: студентка 4 курсу, група КН-18-1  A.V. Raeva
Курс, група виконавця Шифр ініціали, прізвище

Керівник: к.т.н., доцент кафедри КН  O.V. Mazurec
Науковий ступінь, посада Шифр ініціали, прізвище

Нормоконтроль: к.т.н., доцент кафедри КН  P.O. Bagrii
Науковий ступінь, посада Шифр ініціали, прізвище

До захисту допускаю:
Зав. кафедри КН, д.т.н., професор  O.V. Barma
Шифр ініціали, прізвище

13 червня 2022 р.

ХМЕЛЬНИЦЬКИЙ НАЦІОНАЛЬНИЙ УНІВЕРСИТЕТ

Факультет інформаційних технологій

Кафедра комп'ютерних наук


Освітній ступінь бакалавр

Галузь знань 12 – Інформаційні технології

Спеціальність 122 – Комп'ютерні науки

ЗАТВЕРДЖУЮ

Завідувач кафедри комп'ютерних наук


(підпис)
д.т.н., професор О.В. Бармак
« 25 » березня 2022 року

**ЗАВДАННЯ
НА КВАЛІФІКАЦІЙНУ РОБОТУ БАКАЛАВРА**

1. Тема кваліфікаційної роботи бакалавра: «Метод автоматизованого визначення семантичної подібності цифрових текстів для виявлення запозичень»

2. Завдання видано студентці Раєвій Анні Вячеславівні
(прізвище, ім'я, по батькові)

3. Керівник роботи доцент кафедри КН Мазурець Олександр Вікторович
(посада, прізвище, ім'я, по батькові)

4. Затверджено наказом університету від «01» березня 2022 р. № 18

5. Зміст пояснювальної записки (перелік задач) та вихідні дані:

Мета роботи – розробка та апробація методу автоматизованого визначення семантичної подібності цифрових текстів для виявлення запозичень. Слід забезпечити виконання функцій аналізу тексту, пошуку ключових слів, обрахунку оцінок рівня запозичень тестового цифрового тексту щодо кожного цифрового тексту базової вибірки, сортування результуючого переліку за оцінкою рівня запозичень та формування висновку щодо визначення семантичної подібності цифрових текстів.

Виконавець: студентка 4 курсу, група КН-18-1  А.Р. Расва
Курс, група виконавця (підпис) Ініціали, прізвище

Керівник: к.т.н., доцент кафедри КН  О.В. Мазурець
Науковий ступінь, посада (підпис) Ініціали, прізвище

Анотація

Тема кваліфікаційної роботи бакалавра: «Метод автоматизованого визначення семантичної подібності цифрових текстів для виявлення запозичень»

Виконавець кваліфікаційної роботи бакалавра: студентка групи КН-18-1 Раєва Анна Вячеславівна

Керівник кваліфікаційної роботи бакалавра: к.т.н., доцент кафедри КН Мазурець Олександр Вікторович

Кваліфікаційна робота бакалавра містить:

Пояснювальна записка				Кількість додатків
Сторінок	Рисунків	Таблиць	Джерел інформації	
73	29	8	50	3

Метою кваліфікаційної роботи бакалавра є розробка методу автоматизованого визначення семантичної подібності цифрових текстів для виявлення запозичень, а також його програмна реалізація. Для розробки інформаційної системи було використано платформу .NET, мову програмування C#, систему керування базами даних MS SQL Server.

Напрямами практичного використання розробленої інформаційної системи автоматизованого визначення семантичної подібності цифрових текстів для виявлення запозичень визначено порівняння цифрових текстів та визначення умовного показника рівня запозичень.

Ключові слова: ключові слова, семантичні одиниці, плагіат.

Виконавець: студентка 4 курсу, група КН-18-1
Курс, група виконавця


Підпис

А.Р. Раєва
Ініціали, прізвище

Зміст

Перелік скорочень	2
Вступ.....	3
Розділ 1 Характеристика предметної області: аналіз моделей, методів та реалізацій.....	5
1.1 Аналіз інформаційних моделей.....	5
1.2 Огляд теоретичних підходів до розв’язку подібних задач	9
1.3 Аналіз існуючих програмних рішень.....	14
1.4 Аналіз сучасних засобів створення програмного забезпечення	21
1.5 Мета, задачі та вимоги до реалізації інформаційної системи	27
Розділ 2 Проектування інформаційної системи	30
2.1 Метод автоматизованого визначення семантичної подібності цифрових текстів для виявлення запозичень	30
2.2 Інформаційна структура системи	34
2.2.1 Проектна архітектура системи та взаємозв’язок компонентів.....	34
2.2.2 Інформаційна модель.....	36
2.3 Вибір засобів розробки інформаційної системи	40
2.3.1 Вибір мови програмування	41
2.3.2 Вибір фреймворку.....	42
2.3.3 Вибір СКБД	43
Розділ 3 Програмна реалізація інформаційної системи	44
3.1 Структура та функціональне призначення програмних складових системи.....	44
3.2 Особливості реалізації програмних складових системи	45
3.3 Тестування інформаційної системи	49
3.4 Інструкція користувача.....	57
Висновки	67
Перелік посилань.....	69
Додатки	

Перелік скорочень

Скорочення, термін, позначення	Пояснення
БД	База даних
ІС	Інформаційна система
ІТ	Інформаційні технології
КРБ	Кваліфікаційна робота бакалавра
КН	Комп'ютерні науки
ПП	Програмний продукт
СКБД	Система керування базами даних
ХНУ	Хмельницький національний університет.
CLR	Common Language Runtime
JVM	Java Virtual Machine
CLR	Common Language Runtime
BCL	Base Class Library
DE	Disperce Evaluation
LSA	Latent Semantic Analysis
MS	Microsoft
MFC	Microsoft Foundation Class
TF	Term Frequency
TF-IDF	Term Frequency – Inverse Document Frequency

Вступ

Метою кваліфікаційної роботи бакалавра є розробка методу автоматизованого визначення семантичної подібності цифрових текстів для виявлення запозичень, а також його програмна реалізація.

Об'єкт дослідження – процес виявлення запозичень у цифрових текстах.

Предмет дослідження – інформаційні технології, моделі, методи та засоби для автоматизованого визначення семантичної подібності та рівня запозичень у цифрових текстах.

Основою людського існування є потреба творити. Це може бути різного роду творчість, проте уся вона має право бути захищеною від відтворювання сторонніми особами. Часто така творчість має на меті здобуття матеріальної винагороди. Чим більше людство існує, тим більше результатів людської творчості існує і дуже часто автори своїх творів стикаються з тим, що їхні здобутки використовуються повторно.

З метою захисту авторських творів створено авторське право. Воно дозволяє авторам захистити свої твори від присвоєння іншими особами, або використання в інших цілях.

Присвоєння усього твору або його частини іншим автором є шахрайством та має назву плагіат. Часто з плагіатом стикаються у навчальних закладах, тому виокремлено поняття академічного плагіату. Академічний плагіат зустрічається у студентських роботах, монографіях, дисертаціях, наукових статтях. Плагіат можна поділити на дві категорії: умисний та неумисний. Неумисний плагіат часто виникає по причині невміння цитувати та посилатися на першоджерела.

Для боротьби з такого роду шахрайством було створено спеціальні системи для боротьби з плагіатом. Суть таких систем полягає в тому, щоб

виявити запозичення текстів у творах інших авторів. Для створення бази творів використовуються репозитарії.

Такі системи розроблюються на основі семантичного аналізу текстів, що дозволяє виявити навіть прихований плагіат. Однією із задач семантичного аналізу є пошук ключових слів та словосполучень у тексті, що дозволяють описати головну суть вхідного тексту. Для цього використовуються різні методи, наприклад оцінка частоти слів, оцінка TF-IDF, дисперсійна оцінка.

Отже, можна зробити висновок про те, що задача автоматизованого визначення семантичної подібності цифрових текстів для виявлення запозичень є актуальною у всьому світі. Хоч і існує багато методів для розв'язку цієї задачі, проте майже всі вони потребують вдосконалення, адже шахраї вдаються все до нових методів обходу систем антиплагіату.

Розділ 1 Характеристика предметної області: аналіз моделей, методів та реалізацій

1.1 Аналіз інформаційних моделей

В епоху стрімкого розвитку цифрової інформації виникає проблема авторства та плагіату результатів інтелектуальної діяльності людини.

Результат інтелектуальної діяльності людини (РІДЛ) – це термін, що означає продукт діяльності людини, що є нематеріальним та комерційним, а також підлягає використанню або самим автором, або іншими людьми, також він є об'єктом інтелектуальної власності (ОІВ). РІДЛ надається правова охорона, що регулюється чинним законодавством. Особа, яка володіє РІДЛ, має право встановити обмеження щодо використання чи повністю заборонити використання ОІВ сторонніми особами. [1] З метою отримання соціальних благ автори РІДЛ можуть отримати авторське право, яке закріплює право власності.

Авторські права – це права, які належать автору у зв'язку зі створенням РІДЛ. Вони регулюють відносини, що виникають в наслідок створення та використання РІДЛ. Перш за все авторське право створює сприятливі умови для творчої діяльності, для поширення цієї діяльності у суспільстві [2].

Автор – особа, яка шляхом інтелектуальної діяльності створює твір. При чому тут варто також зазначити, що твір може мати декілька авторів, які різною мірою створювали твір – такі особи називаються співавторами [3]. Проте, не дивлячись на захист авторських прав, поряд існує поняття плагіату. Плагіат не є новим явищем, його поширеність зросла в масштабах до такої міри, що набуває майже епідемічних розмірів [4].

Плагіат – це публікація частини або усього твору особою, яка не являється автором або співавтором цього твору. Плагіат – це порушення авторських прав. Особа, що порушує авторські права може бути притягнена до відповідальності згідно законодавства (ст. 50 Закону України «Про авторське право і суміжні права») [3].

Доволі розповсюдженим є академічний плагіат – це явище використання науковими співробітниками (викладачами, докторантами, студентами) фрагментів чужих творів. Частіше всього плагіат зустрічається у студентських роботах, таких як курсові проекти, дипломні проекти. При чому плагіат у даних видів робіт може бути як навмисним, так і не навмисним із-за невміння цитування джерел. Тому плагіат також став загальноприйнятим терміном для позначення неналежного цитування джерел, копіювання чи перефразування без належної пунктуації чи атрибуції або надання чужої роботи як власної. Більшість вищих навчальних закладів мають офіційну політику щодо плагіату [5].

Часто зустрічається явище фальсифікації, що має на увазі спотворення фактів, результатів проведених досліджень або вигадкування показників з метою обходу плагіату [3]. Проблема плагіату – це неправильне використання матеріалу, воно може бути навмисним, або ж ні. Якщо бути точним, то більш широко поняття плагіату можна пояснити використанням матеріалу без посилання на слова, ідеї та факти, що взяті з певного твору. Таким чином, можна сказати, що особа, яка вчиняє плагіат претендує на право власності та на інтелектуальну власність іншого автора [6].

Плагіат можна поділити на три типи [7]:

– повне копіювання чужої роботи без вказання авторства;

- створення твору шляхом копіювання декількох фрагментів з чужих робіт без вказання авторства на ці фрагменти;

- перефразування (або рерайтинг) чужої роботи без оформлення посилань на цю роботу.

У власних творах допустиме цитування невеликих частин творів інших авторів.

Цитата – це невеликий уривок твору іншого автора, який відтворюється з метою підтвердити свої думки або висвітлити погляди на конкретні речі інших авторів [3].

У мережі Інтернет поширені способи обходу плагіату, серед них [8]:

- заміна слів на синоніми – для уникнення плагіату часто застосовують цей метод. Заміна слів синонімами дає можливість зробити текст оригінальнішим при чому не вимагає великих зусиль;

- використання перекладів з іншомовних джерел – надає можливість також уніфікувати текст твору, адже при пошуку плагіату здебільшого системи антиплагіату перевіряють плагіат серед творів створених однією мовою;

- переказ своїми словами – трудомісткий спосіб, який здебільшого зменшує об'єм вихідного твору, або може давати неточні визначення;

- перефразування тексту, тобто зміна порядку слів – зберігає сенс тексту, при перестановці слів. Цей спосіб здебільшого дієвий для мов, в яких може бути довільний порядок слів зі збереженим сенсом вихідного тексту;

- рерайтинг, тобто застосування заміни синонімами, переказу, перефразування одночасно – дуже дієвий та складний процес, адже поєднує в собі декілька методів обходу плагіату, часто використовується для уніфікації тексту професійними рерайтерами.

– застосування спеціального програмного забезпечення – таке ПЗ змінює кодування тексту, замінює слова синонімами, здійснює рерайтинг тощо, по суті заміняє людську роботу, робить це швидко і доволі вдало.

З метою боротьби з плагіатом в деяких країнах передбачено створення національних консультативних служб з плагіату, а також створення спеціального програмного забезпечення, що виявляє запозичення у роботах у вигляді плагіату.

Основою таких застосунків для перевірки на плагіат є великого розміру репозиторії та технології і методи, що дають можливість перевірити твір на плагіат.

Репозиторій – це ресурс у мережі, що зберігає у собі тексти творів. При чому доступ до репозиторіїв є відкритим, будь-хто може переглянути твір з репозиторію. Зазвичай у кожного університету є свій репозиторій, що містить у собі наукові дисертації, наукові видання, навчальні видання, роботи студентів [9].

В Україні також у 2018 році створено Національний репозиторій академічних текстів, в якому зібрані різноманітні наукові роботи, як от монографії, статті, підручники, кваліфікаційні роботи тощо. Власником даного репозиторію є Міністерство освіти та науки [10]

Часто, щоб уникнути плагіату вдаються до наступних методів [11]:

- замінюють деякі літери українського алфавіту на англійські;
- максимально використовують синоніми;
- використовують спеціальні програми, що роблять текст унікальним;
- вставляють текст зображеннями;
- вставляють переноси.

Проте сучасні засоби для виявлення плагіату можуть виявляти і ці хитрощі. На справді для того, щоб плагіату не було у творі, необхідно притримуватися правил цитування та правильно оформлювати посилання на джерела.

Отже, на сьогоднішній день захист авторських прав стає доволі серйозною проблемою, адже іноді не усвідомлюючи того, у при створенні нових творів автори створюють плагіат на уже опублікований раніше твір іншим автором. Ця проблема особливо чітко простежується у академічному плагіаті – у студентських кваліфікаційних роботах, наукових роботах викладачів та професорів. Основним рушієм створення плагіату є невміння посилатися на твори інших авторів. Визначати плагіат допомагає спеціалізоване програмне забезпечення, проте для його виявлення необхідною умовою є наявність великого репозиторію.

1.2 Огляд теоретичних підходів до розв'язку подібних задач

Для того, щоб якимось чином обробити текст, який користувач подає на перевірку в системи антиплагіату, використовується низка методів. Більше того, такою обробкою текстів займається ціла наука.

Комп'ютерна лінгвістика – це інженерна дисципліна, що займається розумінням мови (як усної, так і письмової) з точки зору. По суті мова є дзеркалом розуму, обчислювальне розуміння мови також забезпечує розуміння мислення та інтелекту. Також мова є природним засобом спілкування, тому комп'ютери з лінгвістичною компетентністю могли би забезпечити спілкування з людиною для ширшого розуміння та виконання її потреб [12]. Методи, якими

маніпулює у дослідженнях комп'ютерна лінгвістика, часто спираються на теорії та висновки теоретичної лінгвістики, філософської логіки, когнітивної науки та інформатики [12]. Одним із основних етапів для автоматичного розуміння тексту є семантичний аналіз.

Семантичний аналіз – це один із головних етапів автоматичного розуміння текстів, що полягає у виділенні семантичних відносин, тобто значення конструкції мови, на противагу її синтаксису, формуванні семантичного представлення текстів [13].

Текст – це структура, яка представляє собою витвір мовлення людини та складається з речень. Існує ряд наук, які займаються вивченням тексту [14]:

- синтактика – зв'язність тексту;
- лінгвістика – закономірність побудови тексту;
- стилістика – читацьке сприйняття тексту, текст як форма комунікації, мовленнєві стилі;
- прагматика – мовленнєва тактика
- семантика – структурування сенсів у тексті, тощо.

Семантичний аналіз тексту починається з лексичної семантики, яка вивчає значення окремих слів (тобто словникові визначення). Потім семантичний аналіз досліджує зв'язки між окремими словами та аналізує значення слів, які об'єднуються в речення. Цей аналіз дає чітке розуміння слів у контексті [15]. Таким чином утворюється семантична мережа, яка представляє значення та відносини слів у тексті. Приклад семантичної мережі зображено на рисунку 1.1.

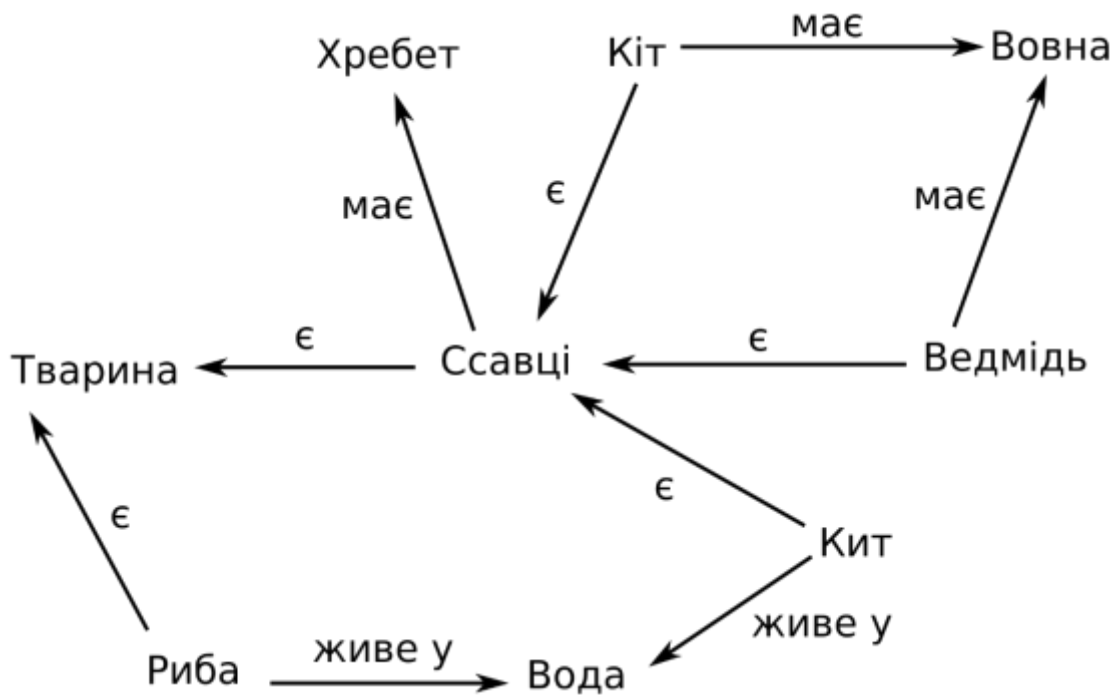


Рисунок 1.1 – Приклад семантичної мережі [16]

У семантичному аналізі є два типи методів, що допомагають дослідити зв'язки слів у реченні. Це семантичні класифікатори та семантичні екстрактори.

Семантичні класифікатори – це моделі класифікації тексту, які призначають будь-які попередньо визначені категорії даному тексту. При чому класифікація проводиться за декількома типами: класифікація тем, настроїв, намірів [17].

Семантичні екстрактори – це методи виділення відповідних слів і виразів або сутностей у будь-якому тексті, щоб з'ясувати детальну інформацію про них. Наприклад, для виділення відповідних слів і виразів семантичний екстрактор в основному використовується разом з різними моделями класифікації. Він використовується для аналізу ключових слів у корпусі тексту та визначення того, які слова мають негативний сенс, а які позитивний. Слова, які найчастіше згадуються у тексті можуть дати зрозуміти про що цей текст, ці слова називаються ключовими [17].

Ключове слово – слово що використовується для вираження деяких моментів у тексті, це слово, що має смислове навантаження [18]. Таких слів у тексті може бути виявлено багато. Ключові слова також використовуються для пошуку в інтернеті, або для створення бібліографі в бібліотеках – це спрощує пошук необхідних книг, журналів, статей тощо. Для визначення ключових слів розроблено безліч інструментів, які допомагають це зробити без втручання людини.

Ключовими можуть бути не тільки слова, а й словосполучення. Іноді для розкриття змісту тексту самих ключових слів недостатньо. В такому випадку ключові словосполучення більш широко розкривають зміст та семантичний зв'язок.

Для визначення ключових слів та словосполучень існує декілька методів, наприклад:

- Term Frequency (TF);
- Term Frequency-Inverse Document Frequency (TF-IDF);
- Disperse Evaluation (DE);
- Text Range;
- тощо.

Далі розглянуто деякі з цих методів.

Term Frequency – даний метод обраховує частоту входження слова у тексті відносно загальної кількості слів в ньому. Якщо аналізується корпус текстів, що частота кожного слова обраховується конкретно для кожного тексту у вхідному корпусі. Таким чином створюється частотні словники для кожного тексту, які використовуються для передавання сенсу мови, створення нових словників та для досліджень в лінгвістиці [19]

Term Frequency-Inverse Document Frequency – статистичний показник, який використовується для визначення важливості кожного слова у вхідному тексті. На відміну, від Term Frequency, в цьому методі має значення частота конкретного слова у інших текстах, що належать до корпусу. Вага слова пропорційна частоті даного слова у тексті та обернено пропорційна частоті цього слова у інших текстах корпусу. Цей метод є доволі розповсюдженим у задачах аналізу текстів, наприклад, визначення релевантності тексту щодо пошукового запиту користувача. Чим частіше в тексті зустрічається слово та найменш часто зустрічається в інших текстах корпусу, тим вища оцінка TF-IDF [20].

Disperse Evaluation – дисперсійна оцінка дискримінантної сили слів у тексті шляхом визначення відстані між словами у тексті. Слід зазначити, що метод визначення ключових слів у тексті шляхом обрахунку дисперсійної оцінки підходить для тих текстів, у яких слова розміщуються відносно нерівномірно. Якщо слова у тексті розміщуються чітко рівномірно, тоді дисперсна оцінка для слів буде дорівнювати 0 [21].

Отже, обробка природної мови людини є доволі популярним та складним завданням. Вивченням обробки мови займається комп'ютерна лінгвістика. Одна із важливих задач, якою займаються науковці – це семантичний аналіз текстів з метою визначення сенсу та взаємозв'язку слів у реченнях. Умовно методи, які допомагають дослідити зв'язки слів у реченнях можна поділити на дві групи: семантичні класифікатори та семантичні екстрактори. При чому для різного типу задач вони можуть поєднуватися. Також у підрозділі розглянуто деякі методи визначення ключових слів у тексті. Деякі з них допомагають визначати ключові слова незалежно від інших текстів у корпусі, а для інших необхідно враховувати усі тексти наявні в корпусі.

1.3 Аналіз існуючих програмних рішень

На теперішній час існує багато сервісів для перевірки текстів на плагіат. Здебільшого вони розміщені як веб-ресурси, щоб кожен бажаючий міг перевірити свою роботу на плагіат. Дані системи антиплагіату використовуються як для перевірки на академічний плагіат, такі і для перевірки інтернет-контенту SEO-спеціалістами перед розміщенням їх на свої сайти. Серед них є як безкоштовні, так і платні ресурси. Далі буде розглянуто найпопулярніші з них.

Одним із популярних сервісів у світі є Unicheck. Його використовують близько 1100 університетів у світі. Цей сервіс працює як програмне забезпечення що допомагає виявити плагіат у тексті, а також щоб порівняти його з іншими роботами, доступними в Інтернеті. Unicheck є платною програмою з можливістю використання безкоштовно деяких її функцій. Розробники системи гарантують безперебійну роботу та доступність, також надають цілодобову підтримку для користувачів. Кожен користувач має свій індивідуальний профіль з типом доступу: індивідуальний, індивідуальний+, школи та університети [22].

В загальному Unicheck надає лише один простий інструмент у використанні, а саме перевірку на плагіат, але гарантує, що він достатньо ефективний, щоб вирішувати навіть найхитріші маніпуляції з текстом. Також дана система антиплагіату надає корисні дані під час перевірки, наприклад, скільки тексту було запозичено, наскільки він схожий на матеріали, які подані в джерелах [23]. Інтерфейс системи Unicheck зображено на рисунку 1.2.

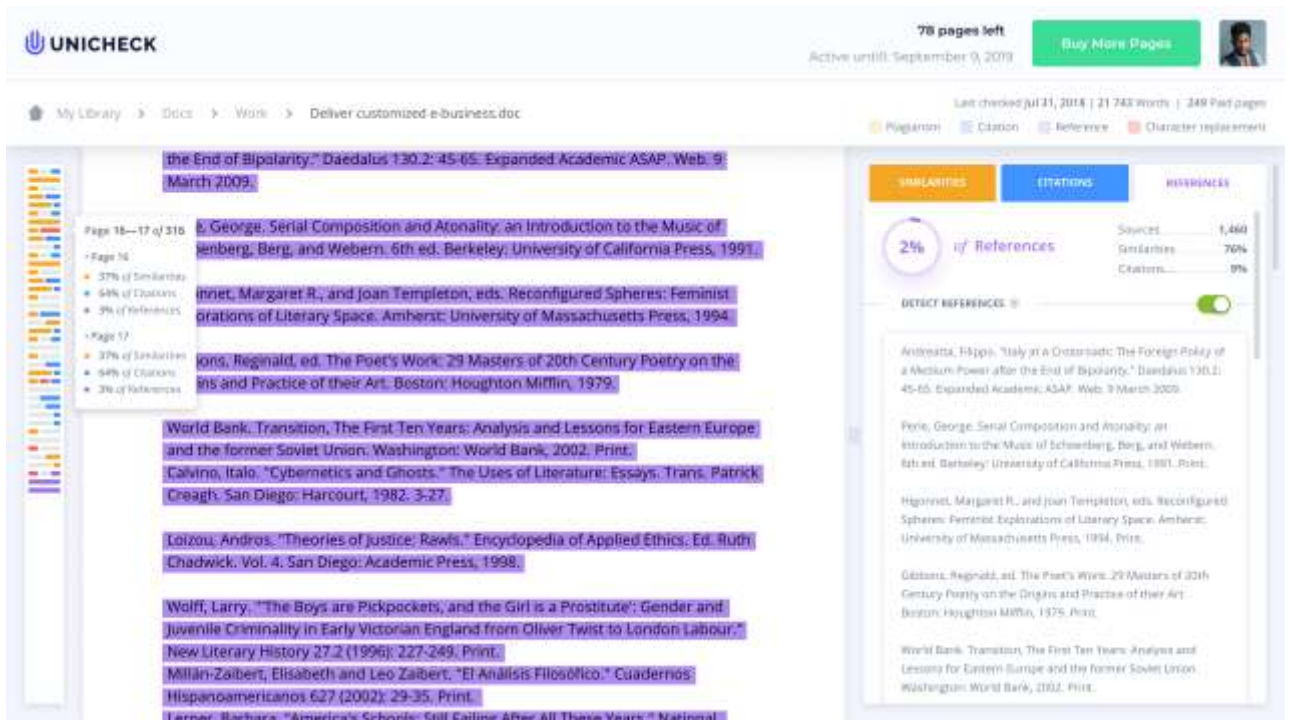


Рисунок 1.2 – Інтерфейс системи Unicheck [22]

Turnitin – система світового рівня для перевірки студентських робіт на академічну доброчесність. Skorистatись нею можуть і студенти, і викладачі. Система перевіряю на плагіат за допомогою своєї бази даних робіт. Якщо завантаженої на перевірку роботи не знайдено у базі, то вона автоматично додається у неї [24]. До того ж, база даних містить копію загальнодоступних ресурсів мережі Інтернету, для цього використовується сканер, який постійно сканує сайти. База даних також містить комерційні, захищені авторським правом книги, журнали та газети.

Дана система є дещо суперечливою, адже проводились дослідження, які показали, що вона не розпізнала плагіат, який був створений іншою системою шляхом перестановки слів. Також деякі студенти та викладачі вважають, що завантажуючи роботу для перевірки, створюється презумпція вини [25]. Інтерфейс описаної системи перевірки на плагіат зображено на рисунку 1.3.

feedback studio Joe Bloggs The Goliath of the Sea 20 /100 1 of 30

The majestic blue whale, the goliath of the sea, certainly stands alone within the animal kingdom for its adaptations beyond its massive size. At 30 metres (98 ft) in length and 190 tonnes (210 short tons) or more in weight, it is the largest existing animal and the heaviest that has ever existed. Despite their incomparable mass, aggressive hunting in the 1900s by whalers seeking whale oil drove them to the brink of extinction. But there are other reasons for why they are now so endangered.

The blue whale's common name derives from bluish-hue that covers the upper side of its body, while its Latin designation is *Balaenoptera musculus*. The blue whale belongs to the Mysticeti suborder of cetaceans, also known as baleen whales, which means they have fringed plates of fingernail-like material, called baleen, attached to their upper jaws. Blue whales feed almost exclusively on krill, though they also take small numbers of copepods. An adult blue whale can eat up to 40 million krill in a day.

These gargantuan beasts used to dominate all the oceans of the Earth up until the late nineteenth century, when the technology was developed to effectively hunt and harvest them. In 1864, the Norwegian Svend Foyn equipped a steamboat with harpoons specifically designed for catching large

Match Overview

45%

1	animals.nationalgeogr...	14%
2	aganews.com	12%
3	pro-solutions.techhelp...	12%
4	animals.partner.blogsp...	7%

Рисунок 1.3 – Інтерфейс системи Turnitin [25]

Сорускаре – система перевірки на плагіат текстів в мережі Інтернет. Це онлайн-сервіс, який надає постійний доступ до своїх функцій. Дана система виконує пошук за допомогою API від Google. Це потужний інструмент, який стає в нагоді для виявлення випадків копіювання вже опублікованого контенту іншими сайтами, а також для виявлення випадків, коли даний вміст перефразовується та продається під виглядом оригінального контенту. По суті, цей сервіс повертає список сторінок, що підозрюються у використанні неоригінального контенту. Є можливість розмістити банер з логотипом Сорускаре на своєму сайті для застереження шахраїв про те, що вміст сайту захищений [26]. Інтерфейс Сорускаре зображено на рисунку 1.4.

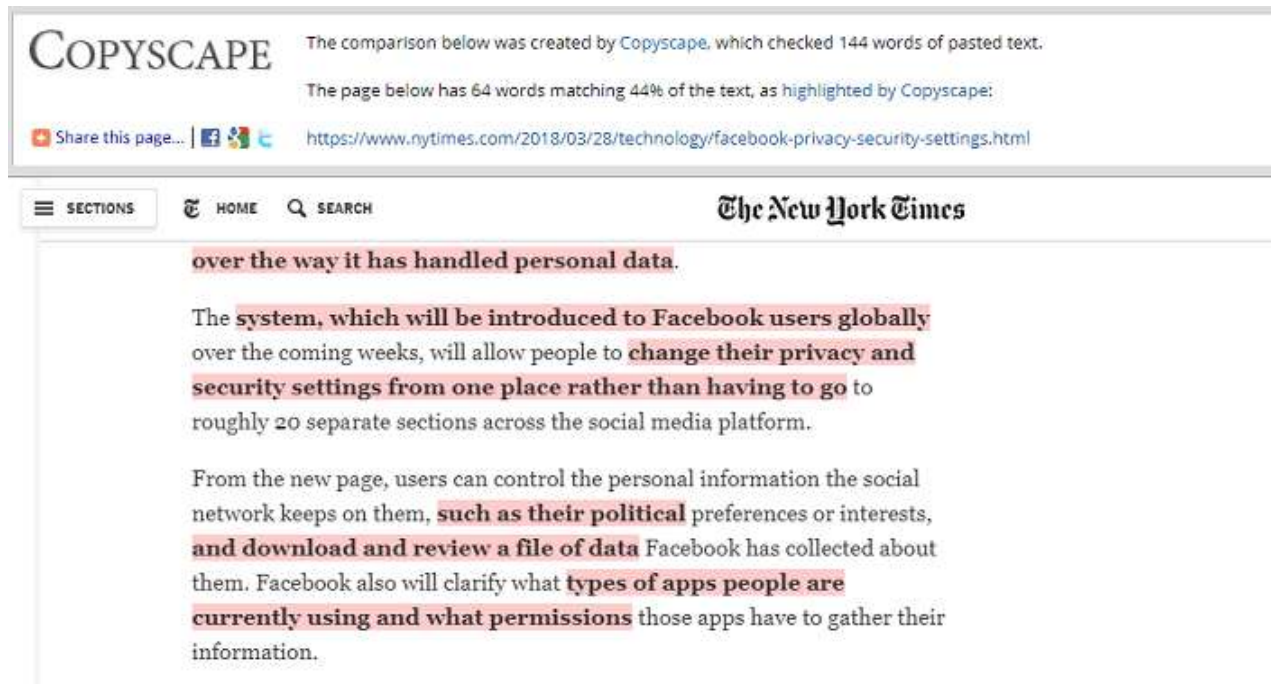


Рисунок 1.4 – Інтерфейс системи Copyscape [27]

Plagiarism Checker by EduBirdie дозволяє перевіряти на плагіат різного типу роботи: резюме, контент для сайт, письмові роботи тощо. Можна перевіряти текст як частинами, так і завантажувати у файлі. Допомогає виявити прямий плагіат, текст, що містить велику частину чужого тексті, самоплагіат, виявляє неіснуючі у посиланнях джерелах. Система є мультимовна. Підтримуються багато форматів, на відміну від більшості шашок, зокрема PDF, .txt, .doc, .docx, .rtf та інші. Процес перевірки відбувається швидко, займає від 10 до 20 секунд. Для використання цієї системи не потрібно реєструватися чи оформлювати підписку. Вона є абсолютно безкоштовною [28]. Інтерфейс Plagiarism Checker by EduBirdie зображено на рисунку 1.5.

⚠ Для того щоб отримати точні результати, будь ласка, заповніть всі поля

Виберіть тип контенту

Письмова робота
 Контент для сайту
 Резюме
 Інше

Вставте заголовок

plagiarism destroys your career life long time ago students cannot seem to know if the paper they submitted contained plagiarized content all they knew was that they were already facing the crime of piracy sadly they do not have access to free tools like plagiarisma today scholars are able to correct and edit similar idioms after scanning their documents with software teachers also will find the tool lessening the burden of returning research works and filing cases

Довжина тексту: 1524 (Без пробілів: 1277) [↻ Перевірити ще один текст](#)

0.0% Унікальність тексту

Текст копіює

Ресурси:	Відсоток співпадіння:	Знайти в тексті:
https://desktop-plagiarism-checker.en.softonic.com/comments	100.0	Показати

[Показати всі співпадіння](#)

Рисунок 1.5 – Інтерфейс системи Plagiarism Checker by EduBirdie [28]

Plagium – потужна система для глибокого пошуку плагіату документів у загальнодоступній всесвітній мережі та приватних репозиторіях. Для користування системою необхідно зареєструватись. Також розробники пропонують обрати один із індивідуальних планів користування. Безкоштовно можна перевірити лише текст довжиною до 1000 символів. Розробниками випущена також мобільна версія, що значно розширює аудиторію користувачів. В залежності від обраного індивідуального плану користувачам можна здійснювати пошук на плагіат обмежену кількість разів, обравши один із просунутих планів стає доступний наступний функціонал: забезпечення

додаткового рівня пошуку плагіату, можливість завантажувати .doc та .pdf документи, створення звітів з роботи, а також порівняння конкретних документів та робота з Google Документами [29]. Інтерфейс Plagium зображено на рисунку 1.6.

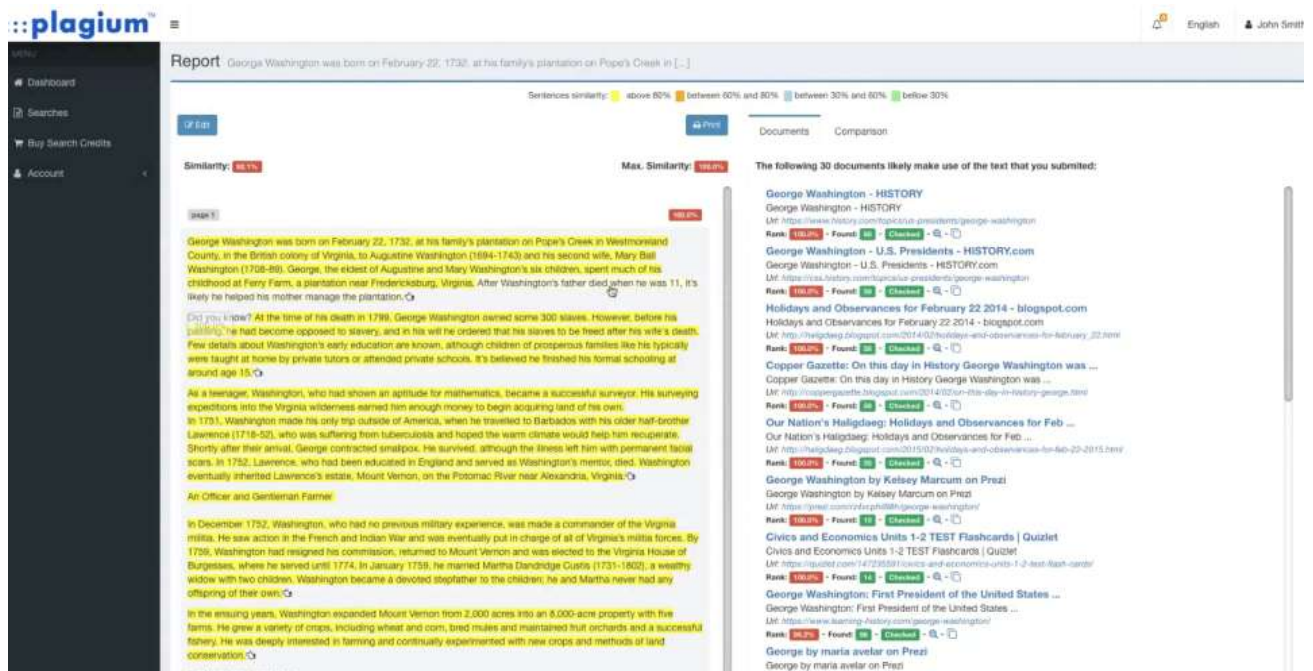


Рисунок 1.6 – Інтерфейс системи Plagium [30]

Питанням виявлення плагіату займаються і науковці в усьому світі. Відомі підходи постійно удосконалюються та пропонуються покращені методи для виявлення запозичень у текстах.

У статті [31] авторами проведено дослідження та створено експериментальний застосунок для перевірки підозрілого документу із корпусом документів та обчислює загальний коефіцієнт плагіату цього документа з використанням методів TFIDF, кластеризації k-середніх та косинусної подібності. Застосунок створює звіт у вигляді PDF-файлу із текстом підозрілого документа з виділенням кольорами речень, що помічені як плагіат відповідно до

документів з корпусу. У звіті також показується загальне співвідношення плагіату підозрілого документа, часткове співвідношення з кожного вихідного документа, кількість речень, слів та іншу інформацію про підозрілий документ. Цей коефіцієнт допомагає визначити, чи відповідає підозрілий документ вимогам, і чи дійсно він містить плагіат. Єдиним недоліком запропонованого підходу є те, що він використовує лише корпус із 300 наукових статей. Підхід можна розширити, додавши багато документів статей до корпусу, а також додавши можливість шукати в Інтернеті пов'язані документи та використовувати інші методи кодування тексту та вимірювання подібності.

Автор статті [32] використав набори даних PAN-PC-11 і обрахував оцінку TF-IDF текстових документів і косинусоподібні показники подібності між вихідними та підозрілими документами, щоб знайти перетин текстів. Модель успішно створює вектори та вимірює показники подібності. Проте, алгоритм не був розширений далі для автоматичного отримання пов'язаних документів для подальшого виконання (перетворення текстів у n-грами для детального аналізу та виявлення найкращого збігу як джерела плагіату та оцінки точності моделі). Модель створює матрицю косинусної подібності для всіх документів, які використовувались, щоб вручну отримати документи та перевірити на перетин за допомогою онлайн-інструментів. Розширення алгоритму на основі запропонованого методу дозволить точно оцінювати зразки документів.

Автори [33] говорять про те, що сучасні засоби виявлення плагіату здебільшого обмежуються порівнянням підозрілих текстів із плагіатом і потенційних оригінальних текстів на рівні рядка. У своєму дослідженні автори пропонують підвищити точність виявлення плагіату шляхом включення методів обробки природної мови (NLP) в існуючі підходи. Автори пропонують

структуру для виявлення зовнішнього плагіату, в якій використовується декілька методів NLP (Trigram Containment Measure: Baseline, Ferret, LM – Bigram Perplexity, LM – Trigram Perplexity, Longest Common Subsequence, Dependency Relations) для обробки набору підозрілих і оригінальних документів, не тільки для аналізу рядків, але й структури тексту, використовуючи ресурси для обліку текстових відносин. Отримані результати, показали, що методи NLP підвищують точність існуючих підходів.

Отже, можна зробити висновок про те, що задача пошуку плагіату є доволі популярною, і на сьогоднішній день існує безліч програмних продуктів, що вирішують цю задачу. Здебільшого всі ці програмні продукти представлені у вигляді веб-застосунків, проте деякі розробники також реалізують свої ПП у вигляді десктопних та мобільних застосунків. В залежності від того який функціонал необхідний користувачу та скільки готовий заплатити за нього, пропонуються також додаткові функції – збереження звітів по пошуку, робота з файлами, порівняння конкретних документів тощо.

Також для задачі виявлення плагіату запропоновано покращені методи на основі базових, часто з використанням методів обробки природньої мови.

1.4 Аналіз сучасних засобів створення програмного забезпечення

На сьогоднішній день існує декілька способів розробки застосунків, а саме: у вигляді вебзастосунків, мобільних додатків та віконного застосунка для операційних систем Windows, Linux, MacOS та інших менш популярних систем.

Вебзастосунок – це програма, що представляє собою клієнт-сервер, у якій клієнт працює у веббраузері. Користувачі отримують доступ застосунку з будь-

якого комп'ютера, мобільного телефона або планшету та навіть телевізора, підключеного до Інтернету, за допомогою стандартного браузера. Це дає змогу користувача не встановлювати додаткові застосунки на свій пристрій. На базовому веб-сайті сторінки статичні, проте вони постійно взаємодіють з користувачами, надаючи необхідну їм інформацію [34].

Серед переваг вебзастосунків варто виділити:

- вартість – розробка вебзастосунків порівняно з іншими типами невисокою;
- оновлення – можна бути впевненим в тому, що користувач працює з останнім оновленим програмним забезпеченням, адже при розробці застосунку, розробники стараються використовувати найновіші технології та інструменти;
- налаштування – вебзастосунки є доволі легкими у налаштуванні, та зазвичай вважаються безпечними у використанні;
- вимога у завантаженні – як така відсутня, так як для використання сайту необхідно мати лише браузер на пристрої.
- підтримка платформ – вебзастосунки є на стільки універсальними, що будь то Windows, MacOS, Android чи iOS – користувач гарантовано отримає доступ до вебзастосунку [35].

Недоліками вебзастосунків слід вважати:

- продуктивність – завжди пов'язана із продуктивністю браузера на пристрої, тому якщо останній має нестабільну роботу, то і робота вебзастосунку буде нестабільною;
- безпека – хоч і вебзастосунки вважаються одними із безпечніших в плані неможливості встановлення вірусного ПЗ, проте у вебзастосунках відсутня фільтрація якості контенту;

– залежність від Інтернету – доступ до веб-сайтів завжди залежить від наявності та якості зв'язку Інтернет [35].

Таким чином можна стверджувати, що вебзастосунки мають багато переваг та невелику кількість недоліків, що робить їх дуже популярними та широко застосовуваними.

Мобільний застосунок створений для роботи на мобільному пристрої, яким може бути планшет або смартфон. Більше того, мобільні застосунки часто мають обмежену функціональність, і не можуть запропонувати користувачам повний функціонал, як наприклад віконний застосунок.

Здебільшого розробка мобільних застосунків здійснюється додатково до віконного застосунку або вебзастосунку для збільшення аудиторії та розширення доступності. Мобільні застосунки допомагають компаніям таким чином збільшувати прибуток, а також краще обслуговувати своїх клієнтів – вони є прямим маркетинговим каналом для роботи з клієнтами [36].

Серед позитивних сторін мобільних застосунків варто відмітити:

– швидкість – мобільні застосунки зазвичай в 1,5 рази швидше опрацьовують запит користувача, ніж вебзастосунок;

– персоналізація та адаптивність – користувачі люблять підлаштовувати інтерфейс під себе, починаючи від мови застосунка і до колірної схеми застосунка, це вказує на те, що мобільні застосунки є доволі адаптивними;

– безперебійна робота – більшість мобільних застосунків мають можливість роботи як «онлайн» так і «офлайн» або з часто збереженим функціоналом при «офлайн» роботі;

- миттєві сповіщення – мобільні застосунки мають можливість миттєво сповіщати користувачів про певні події пов'язані з роботою застосунку та швидкий доступ до цих сповіщень;

- доступ до внутрішніх функцій телефону – є великою перевагою, коли застосунок може використовувати технології NFS, GPS, сканування QR-кодів [37]

Недоліками мобільних застосунків є:

- мобільні застосунки не замінюють вебзастосунки – для успішного розвитку бізнесу також потрібен свій сайт для більшого охоплення аудиторії, важко лише за допомогою одного мобільного застосунка охопити велику аудиторію;

- потреба в розробці на декілька платформ – не слід забувати про те, що на сьогодні не існує єдиної мобільної платформи. Їх як мінімум дві: Android та iOS, тому задля охоплення великої аудиторії слід орієнтуватися на два типи користувачів;

- часте оновлення – для стабільної роботи мобільного застосунка є необхідність в частих оновленнях та виправленнях помилок у попередніх версіях;

- проблеми з поширенням – публікація мобільного застосунку у магазинах мобільних платформ потребує фінансових затрат та подальшого просування продукту [38].

Отже, незважаючи на доступність, мобільні застосунки мають ряд недоліків, що затрудняють просування їх серед користувачів.

Віконний застосунок – особливий тип десктопного застосунку. Призначений для роботи на комп'ютерах з операційною системою Windows. По

суті є локальною версією вебзастосунків, для якої не завжди потрібне підключення до Інтернет. Якщо десктопні застосунки можуть бути в загальному кросплатформенні, то конкретно віконний застосунок є розроблений спеціально для ОС Windows [39].

Серед переваг віконних застосунків варто відзначити:

- потужність – десктопні комп'ютери на базі Windows дають можливість виконувати велику кількість обчислень та запускати багато процесів одночасно;
- поширеність – Windows є найпоширенішою операційною системою для комп'ютерів у світі;
- стабільна та безпечна робота – фактично робота застосунку залежить лише від конфігурації комп'ютера на якому вона працює;
- безпека користувача – скачуючи офіційні застосунки з офіційних джерел мінімізується ймовірність бути взламаним хакерами чи піддаватися хакерським атакам на відміну від вебзастосунків;
- персоналізація – часто віконні застосунки використовують звичні для користувача налаштування інтерфейсу [40]

Недоліками віконних застосунків є:

- орієнтованість на одну ОС – в даному випадку на операційну систему Windows, для використання на інших операційних системах необхідна або адаптація застосунку під конкретну ОС, або встановлення користувачем спеціально ПЗ для запуску застосунку на іншій ОС;
- залежність від версій ОС – деякі застосунки, що були розроблені під конкретну версію ОС можуть не запускатися на новіших версіях.

– проблеми із встановленням – у користувачів можуть виникати проблеми із встановленням таких застосунків, іноді є потреба звернутися до спеціалістів [41].

Отже, віконні застосунки є найбільш стабільним та потужним програмним забезпеченням в порівнянні з іншими.

Для розробки віконних застосунків використовуються декілька платформ. Наведемо деякі з них.

Платформа Java – незалежна платформа, що дозволяє створювати кросплатформенні додатки. Досягається це тим, що початковий код компілюється у байт-код, який потім можна запустити на будь-якій операційній системі. Це є головною перевагою цієї платформи. Проте головним недоліком є те, що для цього необхідно встановлювати додаткове програмне забезпечення, а саме віртуальну машину Java. Хоч на сьогоднішній день, віртуальна машина дозволяє запускати та виконувати код з невеликою затримкою, проте це відбувається повільніше, ніж безпосереднє виконання коду, наприклад у .NET [42].

Платформа Mono – вільна для розповсюдження платформа на основі платформи .NET. Дана платформа реалізована для декількох ОС – Linux, Windows, MacOS, Solaris. Середовище може виконувати код написаний на мовах C#, Java, Python, PHP та ін. Основною перевагою є її доступність та відкритість. Недоліком є обмежена кількість бібліотек та підтримуваних мов програмування (по суті це спрощена версія .NET) [43].

Платформа .NET – платформа, що запропонована компанією Microsoft. Також є кросплатформенною, застосунки розроблені на ній підтримуються на Windows, FreeBSD та Linux. Переваги, які слід виділити, це потужна підтримка

компанією Microsoft – а саме постійні оновлення та великий набір інструментів розробки, які випускаються як компанією Microsoft, так і сторонніми розробниками. Підтримує більше 20 мов програмування, орієнтована як на десктопні застосунки, так і має застосування у розробці веб-застосунків. Серед недоліків варто відзначити складність для новачків та складність синтаксису [44].

.NET надає бібліотеку класів під назвою Framework Class Library, яка надає можливість створювати інтерфейси користувача, підключення до бази даних, розробку веб-додатків, керування пам'яттю, мережевий зв'язок тощо [45]. Це дає можливість створювати широкофункціональні віконні застосунки за допомогою UWP, WPF та WinForms, що дуже актуально для обраного напрямку дослідження та розробки програмного забезпечення.

Підсумовуючи вищевикладене, для розробки програмного забезпечення, обраного у пункті 1.2, доцільним є зупинити вибір на розробці саме віконного застосунку на платформі .NET.

1.5 Мета, задачі та вимоги до реалізації інформаційної системи

Метою кваліфікаційної роботи бакалавра є розробка методу автоматизованого визначення семантичної подібності цифрових текстів для виявлення запозичень, а також його програмна реалізація.

Для прикладного тестування методу слід виконати його програмну реалізацію у вигляді додатку на платформі .NET, що виконує наступні основні функції:

- створення послідовного вектору слів для кожного цифрового тексту авторського базової вибірки;
- створення послідовного вектору слів для тестового тексту;
- створення частотного вектору семантичного ядра для кожного авторського цифрового тексту;
- створення частотного вектору семантичного ядра тестового цифрового тексту;
- очищення частотних векторів семантичних ядер текстів від стоп-слів;
- обрахунок оцінок рівня запозичень тестового цифрового тексту щодо кожного цифрового тексту базової вибірки;
- сортування результуючого переліку за оцінкою рівня запозичень та формування висновку щодо визначення семантичної подібності цифрових текстів.

Для досягнення мети потрібно виконати наступні завдання:

- Провести аналіз предметної області, а саме: аналіз інформаційних моделей, виконати огляд теоретичних аспектів до розв'язку задачі автоматизованого визначення семантичної подібності цифрових текстів для виявлення запозичень, виконати аналіз існуючих програмних рішень та виконати аналіз сучасних засобів створення програмного забезпечення.
- Виконати проєктування інформаційної системи, у рамках якого створити метод автоматизованого визначення семантичної подібності цифрових текстів для виявлення запозичень; побудувати інформаційну структуру системи; обрати засоби розробки системи на основі методу автоматизованого визначення семантичної подібності цифрових текстів для виявлення запозичень.

– Створити програмну реалізацію інформаційної системи на базі методу автоматизованого визначення семантичної подібності цифрових текстів для виявлення запозичень та провести апробацію результатів.

– Для зручності користування створити інструкцію користувача.

Розділ 2 Проєктування інформаційної системи

2.1 Метод автоматизованого визначення семантичної подібності цифрових текстів для виявлення запозичень

Для автоматизованого визначення семантичної подібності цифрових текстів для виявлення запозичень було розроблено метод, який схематично зображений на рисунку 1.1.

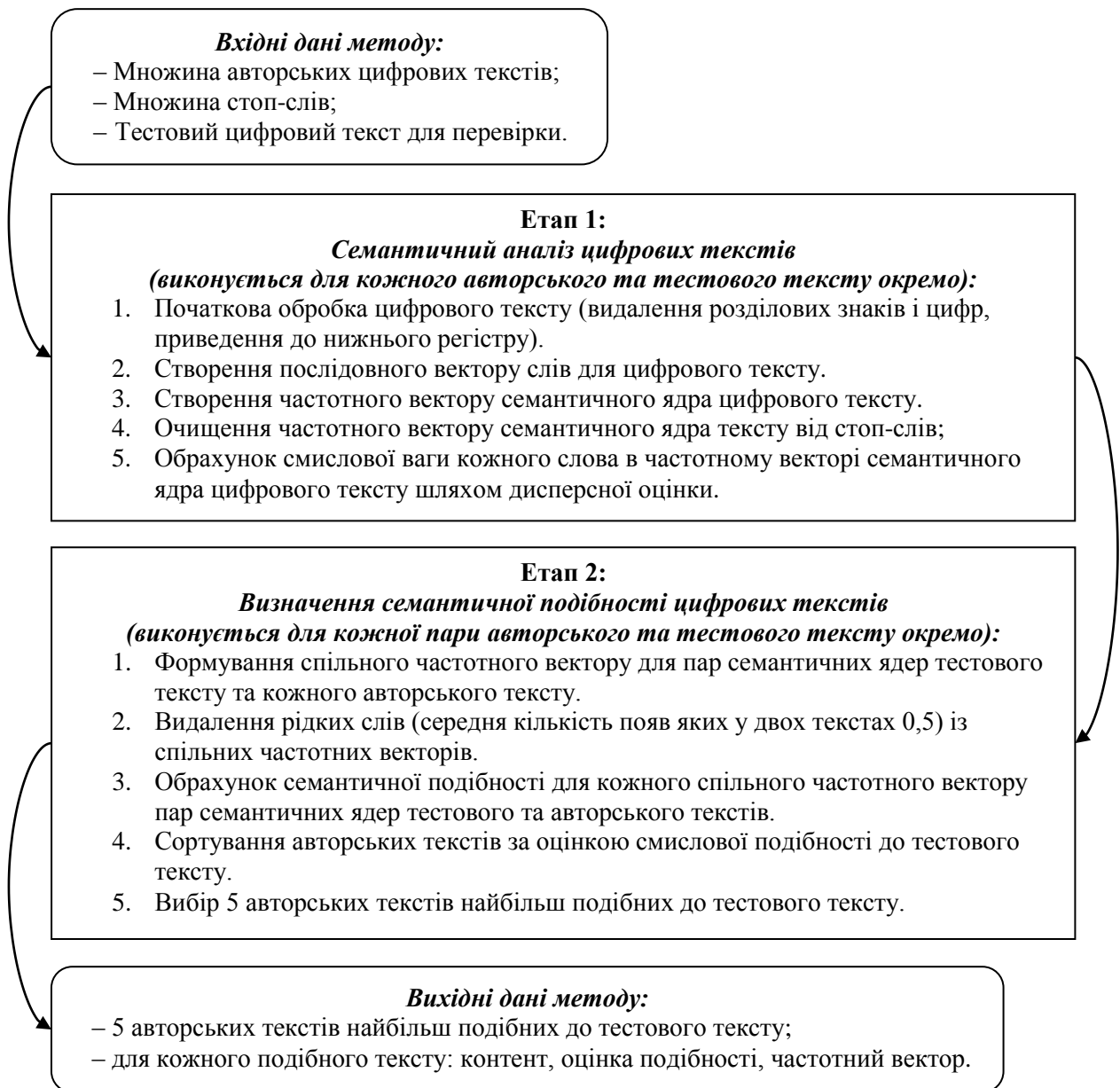


Рисунок 2.1 – Схема методу

Вхідними даними методу є: множина авторських цифрових текстів, множина стоп-слів, тестовий цифровий текст для перевірки.

На першому етапі методу визначення семантичної подібності цифрових текстів для виявлення запозичень для кожного тексту виконується семантичний аналіз цифрових текстів, що включає в себе процес семантичного аналізу цифрових текстів, який проілюстровано на рисунку 2.2.

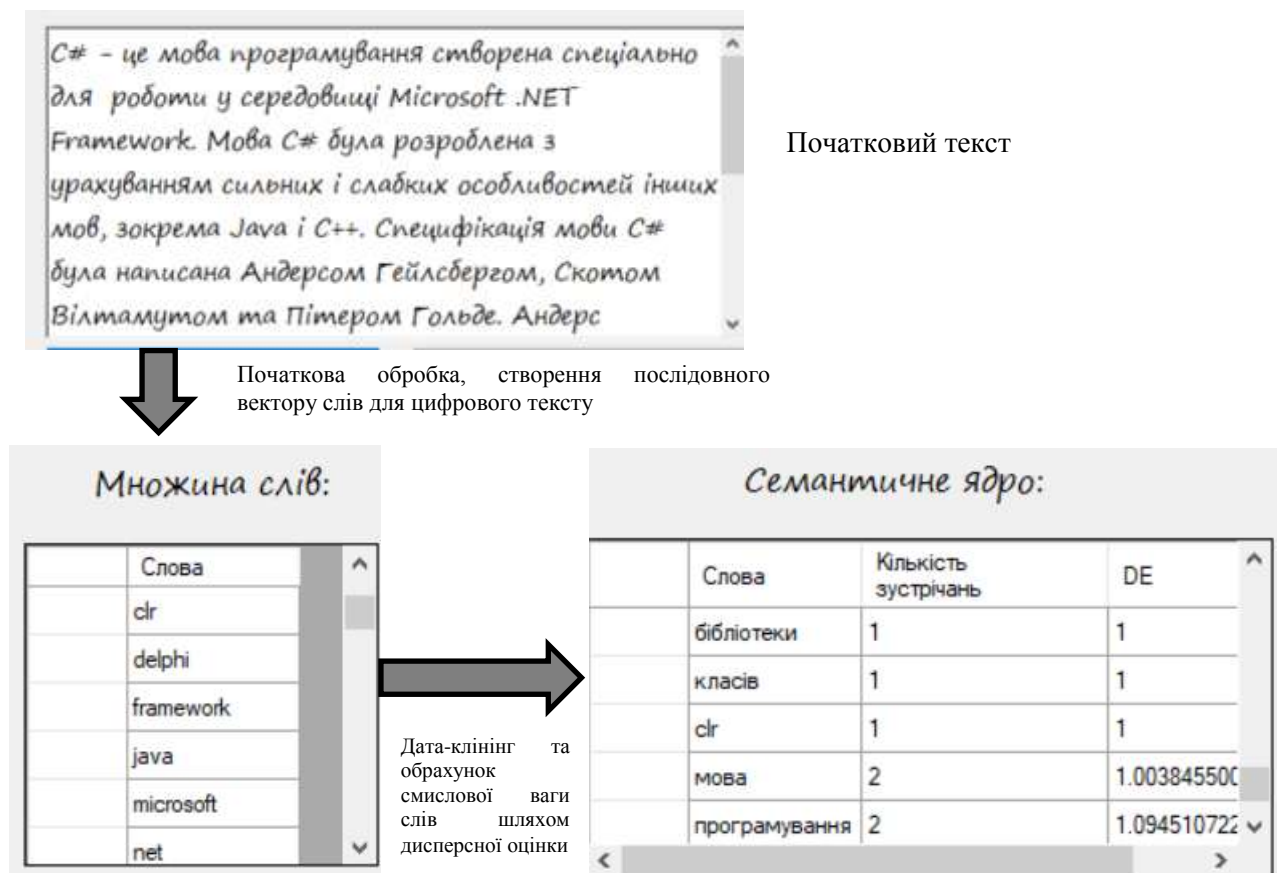


Рисунок 2.2 – Схема процесу семантичного аналізу цифрових текстів

Під час початкової обробки цифрового тексту відбувається видалення спеціальних символів і цифр, приведення тексту до нижнього регістру. Далі відбувається створення послідовного вектору слів для цифрового тексту. Текст

представляється у формі вектору (впорядкованої множини слів). З послідовного вектору слів відбувається створення частотного вектору семантичного ядра цифрового тексту. Після чого отриманий частотний вектор очищується від стоп-слів, визначених користувачем і відбувається обрахунок смислової ваги кожного слова в частотному векторі семантичного ядра цифрового тексту шляхом дисперсної оцінки.

На другому етапі відбувається визначення семантичної подібності цифрових текстів. Цей етап виконується для кожної пари текстів (авторського та тестового). В рамках етапу виконується формування спільного частотного вектору для пар семантичних ядер тестового тексту та кожного авторського тексту, видалення рідких слів із спільних частотних векторів. Рідкісними вважаються ті слова, середня кількість появ яких у двох текстах 0,5. Далі виконується обрахунок семантичної подібності для кожного спільного частотного вектору пар семантичних ядер тестового та авторського текстів та сортування авторських текстів за оцінкою смислової подібності до тестового тексту. На основі оцінок виконується вибір 5 авторських текстів найбільш подібних до тестового тексту з відповідними значеннями відсотку подібності.

Вихідними даними зображеного схематично методу є: 5 авторських текстів найбільш подібних до тестового тексту; для кожного подібного тексту: контент, оцінка подібності, частотний вектор.

Далі більш детально розписано визначення оцінка семантичної подібності для третього кроку другого етапу.

Оцінка семантичної подібності цифрових текстів S_{Semant} тестового тексту $Test$ до авторського тексту $Auth$ обраховується наступним чином:

$$S_{Semant} = \frac{\sum_{i=1}^n DE_{Dif}}{\sum_{i=1}^n DE_{Auth}}, \quad (2.1)$$

де n – кількість слів у спільному частотному векторі для пар семантичних ядер тестового тексту та авторського тексту, DE_{Dif} – обрахована різниця важливості слова i у частотних векторах для семантичних ядер тестового тексту та авторського тексту, DE_{Auth} – оцінка семантичної важливості слова i у частотному векторі для семантичного ядра авторського тексту.

Різниця важливості слова i у частотних векторах для семантичних ядер тестового тексту та авторського тексту в (2.1) обраховується наступним чином:

$$DE_{dif} = DE_{Auth} - |DE_{Auth} - DE_{Test}|, \quad (2.2)$$

де DE_{Auth} – оцінка семантичної важливості слова i у частотному векторі для семантичного ядра авторського тексту, DE_{Test} – оцінка семантичної важливості слова i у частотному векторі для семантичного ядра тестового тексту.

При цьому, одержане в (2.2) значення перш ніж використатись у (2.1) модифікується так:

$$DE_{dif} = \begin{cases} DE_{dif}, & \text{if } DE_{dif} > 0 \\ 0, & \text{if } DE_{dif} \leq 0 \end{cases}, \quad (2.3)$$

Таким чином пропонується автоматизовано визначати семантичну подібність цифрових текстів з метою виявлення запозичень.

2.2 Інформаційна структура системи

2.2.1 Проектна архітектура системи та взаємозв'язок компонентів

Для проєктування інформаційної системи реалізації методу автоматизованого визначення семантичної подібності цифрових текстів для виявлення запозичень було побудовано схему компонентів визначення рівня запозичень (рисунок 2.3).

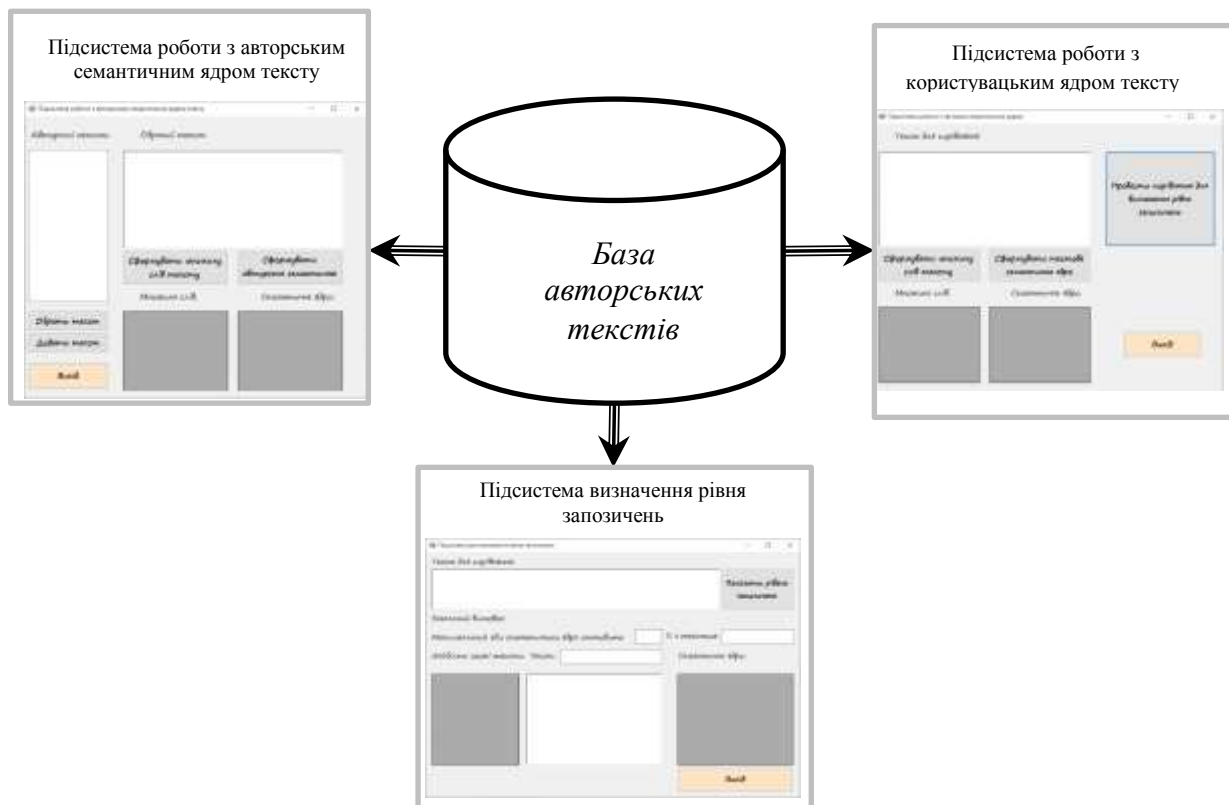


Рисунок 2.3 – Схема компонентів системи визначення рівня запозичень

База даних призначена для збереження авторських текстів та інформації по них, такої як: ключові слова, стоп-слова, тип тексту, сам текст та інші.

Для реалізації методу автоматизованого визначення семантичної подібності цифрових текстів для виявлення запозичень запропоновано виконати

три підсистеми, таких як: підсистема роботи з авторським семантичним ядром тексту, підсистема роботи з користувацьким семантичним ядром тексту та підсистема визначення рівня запозичень.

У підсистемі роботи з авторським семантичним ядром тексту заплановано реалізацію таких функцій:

- Формування та відображення впорядкованої множини слів обраного тексту.
- Додавання нового тексту до бази даних.
- Формування авторського семантичного ядра тексту.
- Збереження в базі даних авторського семантичного ядра тексту.

У підсистемі роботи з користувацьким семантичним ядром тексту заплановано реалізацію таких функцій:

- Формування та відображення впорядкованої множини слів користувацького тексту.
- Формування користувацького семантичного ядра тексту.
- Збереження в базі даних користувацького семантичного ядра тексту.
- Перехід до підсистеми визначення рівня запозичень зі збереженням параметрів.

У підсистемі визначення рівня запозичень заплановано реалізацію таких функцій:

- Виявлення максимального збігу між користувацьким текстом та авторськими та виведення результату користувачу.
- Виведення 5 найбільш схожих тексти на користувацький із оцінками схожості за описаним в 2.1 методом.
- Можливість введення користувацького тексту для порівняння.

Відповідна вищеописана схема компонентів системи визначення рівня запозичень дозволить вирішити поставлені завдання з пункту 1.5.

2.2.2 Інформаційна модель

Для проектування інформаційної системи згідно з темою КРБ було побудовано даталогічну модель бази даних (рисунок 2.4).

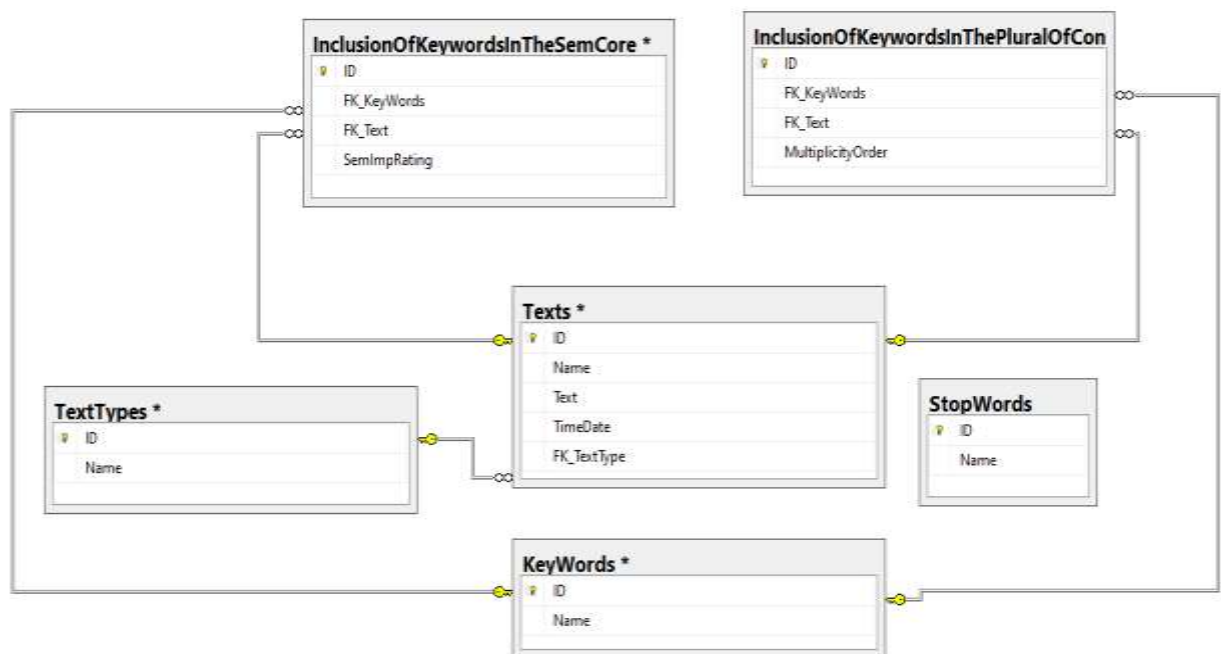


Рисунок 2.4 – Даталогічна модель бази даних інформаційної системи визначення семантичної подібності цифрових текстів для виявлення запозичень

База даних методу автоматизованого визначення семантичної подібності цифрових текстів для виявлення запозичень» містить наступні таблиці: «TextTypes», «KeyWords», «StopWords», «Texts», «InclusionOfKeywordsInTheSemCore» та «InclusionOfKeywordsInThePluralOfCon».

Таблиця «KeyWords» (таблиця 2.1) зберігає ключові слова, що зустрічаються у текстах.

Таблиця 2.1 – Атрибути таблиці «KeyWords»

№ п/п	Назва атрибуту	Тип даних	Опис
1.	ID	int	Первинний ключ, числовий ідентифікатор
2.	Name	Varchar(50)	Ключове слово

У таблиці «StopWords» (таблиця 2.2) містяться стоп-слова, що необхідні при семантичному аналізі тексту.

Таблиця 2.2 – Атрибути таблиці «StopWords»

№ п/п	Назва атрибуту	Тип даних	Опис
1.	ID	int	Первинний ключ, числовий ідентифікатор
2.	Name	Varchar(50)	Стоп-слово

У таблиці «TextTypes» (таблиця 2.3) містяться типи текстів, за якими їх можна класифікувати.

Таблиця 2.3 – Атрибути таблиці «TextTypes»

№ п/п	Назва атрибуту	Тип даних	Опис
1.	ID	int	Первинний ключ, числовий ідентифікатор
2.	Name	Varchar(50)	Назва типу тексту

Таблиця «Texts» (таблиця 2.4) містить інформацію щодо цифрових текстів, а саме: його назву, контент тексту, дату й час, коли текст було додано до БД, а також вторинний ключ, посилання на запис в таблиці «TextType», що вказує на тип тексту.

Таблиця 2.4 – Атрибути таблиці «Texts»

№ п/п	Назва атрибуту	Тип даних	Опис
1.	ID	int	Первинний ключ, числовий ідентифікатор
2.	Name	Varchar(50)	Назва тексту
3.	Text	text	Контент, наповнення тексту
4.	TimeDate	datetime	Дата й час додання тексту до БД
5.	FK_TextType	int	Вторинний ключ, посилання на запис в таблиці «TextType», що вказує на тип тексту

Таблиця «InclusionOfKeywordsInThePluralOfConsecWordsOfTheText» (таблиця 2.5) зберігає інформацію щодо текстів та слів, входжень ключових слів

до семантичного ядра. Таблиця містить наступні поля: посилання на ключове слово, посилання на текст та оцінку семантичної важливості.

Таблиця 2.5 – Атрибути таблиці

«InclusionOfKeywordsInThePluralOfConsecWordsOfTheText»

№ п/п	Назва атрибуту	Тип даних	Опис
1.	ID	int	Первинний ключ, числовий ідентифікатор.
2.	FK_KeyWords	Int	Вторинний ключ, посилання на відповідний запис в таблиці «KeyWords» із записом відповідного ключового слова.
3.	FK_Text	int	Вторинний ключ, посилання на відповідний запис в таблиці «Texts» із записом відповідного тексту.
4.	SemImpRating	Varchar(50)	Оцінка семантичної важливості.

Таблиця «InclusionOfKeywordsInThePluralOfCon» (таблиця 2.6) містить інформацію щодо входження ключових слів у множину послідовних слів тексту та містить наступні поля: посилання на відповідне ключове слово, посилання на відповідний текст, та порядковий номер у множині слів.

Таблиця 2.6 – Атрибути таблиці

«InclusionOfKeywordsInThePluralOfConsecWordsOfTheText»

№ п/п	Назва атрибуту	Тип даних	Опис
1.	ID	int	Первинний ключ, числовий ідентифікатор
2.	FK_KeyWords	int	Вторинний ключ, посилання на відповідний запис в таблиці «KeyWords» із записом відповідного ключового слова.
3.	FK_Text	int	Вторинний ключ, посилання на відповідний запис в таблиці «Texts» із записом відповідного тексту.
4.	MultiplicityOrder	Varchar(50)	Порядковий номер у множині

Таким чином, було реалізовано базу даних для забезпечення справної роботи методу автоматизованого визначення семантичної подібності цифрових текстів для виявлення запозичень.

2.3 Вибір засобів розробки інформаційної системи

Згідно поставленого завдання та описаної вище структури, а саме розробка інформаційної системи автоматизованого визначення семантичної подібності цифрових текстів для виявлення запозичень було прийнято рішення використати наступні засоби для створення відповідних складових:

- мова програмування;
- середовище програмування;

– систему керування базами даних.

Детальніше вибір кожного засобу розглянуто у наступних підрозділах.

2.3.1 Вибір мови програмування

Для розробки віконного застосунку на платформі .NET обрано мову програмування C#.

C# – це об'єктно-орієнтована мова програмування, спеціально розроблена для платформи .NET, тому ідеально підходить для розробки віконних додатків на цій платформі. Вона має багато вбудованих методів та бібліотек, що дозволяють швидко і просто писати програмні продукти. Вона є відмінним вибором для реалізації алгоритмів і структур даних [46]

Завдяки об'єктно-орієнтованому підходу C# добре підходить для роботи і з базами даних, адже за допомогою нього можна легко створювати сутності, що є в базі даних та описувати їх.

Здебільшого C# орієнтований саме на віконні додатки для Windows, що може викликати проблеми із застосуванням на інших ОС, так як частина функцій застосунку може не працювати або потребувати зміни чи доповнення. Це може бути проблемою, адже додаток, що розроблюється для вирішення поставленої задачі орієнтований на пересічного користувача, що може використовувати будь-яку з існуючих операційних систем.

Отже, мова програмування C# є найбільш вдалим вибором для поставленої задачі, тому що вона орієнтована на роботу з віконними додатками для платформи .NET, а також завдяки своєму об'єктно орієнтованому підходу надає можливість описувати сутності бази даних.

2.3.2 Вибір фреймворку

Необхідним для створення програмного забезпечення є і середовище програмування.

Середовище програмування – простими словами це редактор програмного коду, за допомогою якого створюються програми. Може мати різний набір функцій. Для створення програм необхідною складовою є не тільки редактор програмного коду, а й компілятор, інтерпретатор, засоби автоматизації збірки, засоби для створення візуального вигляду ПЗ [47].

Таким середовищем, що забезпечує наявність усіх компонентів, що необхідні для створення ПЗ є Microsoft Visual Studio. Це середовище програмування, розроблене компанією Microsoft та забезпечує доступ до численних функцій та бібліотек .NET. Допомогає швидко та просто створювати програмний код з максимально необхідним для розробника функціоналом. Також забезпечує зв'язок з базою даних для використання даних з неї, що є важливим фактором в рамках поставленого завдання [48].

Таким чином, визначено, що Microsoft Visual Studio – є потужним середовищем програмування, що спрощує роботу розробника під час створення програмних продуктів.

2.3.3 Вибір СКБД

Головну функцію із збереження та роботи з даними в базах даних виконує система керування базами даних. Серед великої кількості таких СКБД обрано Microsoft SQL Server.

Microsoft SQL Server – це ще один програмний продукт, який є необхідним засобом розробки ПЗ. Як мову запитів, використовує SQL.

Microsoft SQL Server використовує вбудовану підтримку .NET. Саме тому, процедури бази даних, що зберігаються, можуть бути реалізовані на будь-якій мові програмування, що підтримується .NET та можуть використовувати його бібліотеки [49].

Серед переваг відзначається простота роботи, безпечність, механізм відновлення даних та продуктивність [50].

Також перевагою Microsoft SQL Server є те, що він може використовуватись і на операційній системі Linux Ubuntu, що певним чином розширює аудиторію користувачів розроблюваного додатку.

Отже, для розробки інформаційної системи прийнято використати ансамбль засобів розробки програмного забезпечення від Microsoft, а саме: платформу .NET, мову програмування C#, середовище розробки Microsoft Visual Studio та систему керування базами даних Microsoft SQL Server. Такий вибір забезпечить злажену роботу компонентів, адже вони створені однією компанією для вирішення конкретних задач.

Розділ 3 Програмна реалізація інформаційної системи

3.1 Структура та функціональне призначення програмних складових системи

Діаграма класів для розроблюваного програмного застосунку по реалізації методу автоматизованого визначення семантичної подібності цифрових текстів для виявлення рівня запозичень зображено на рисунку 3.1.

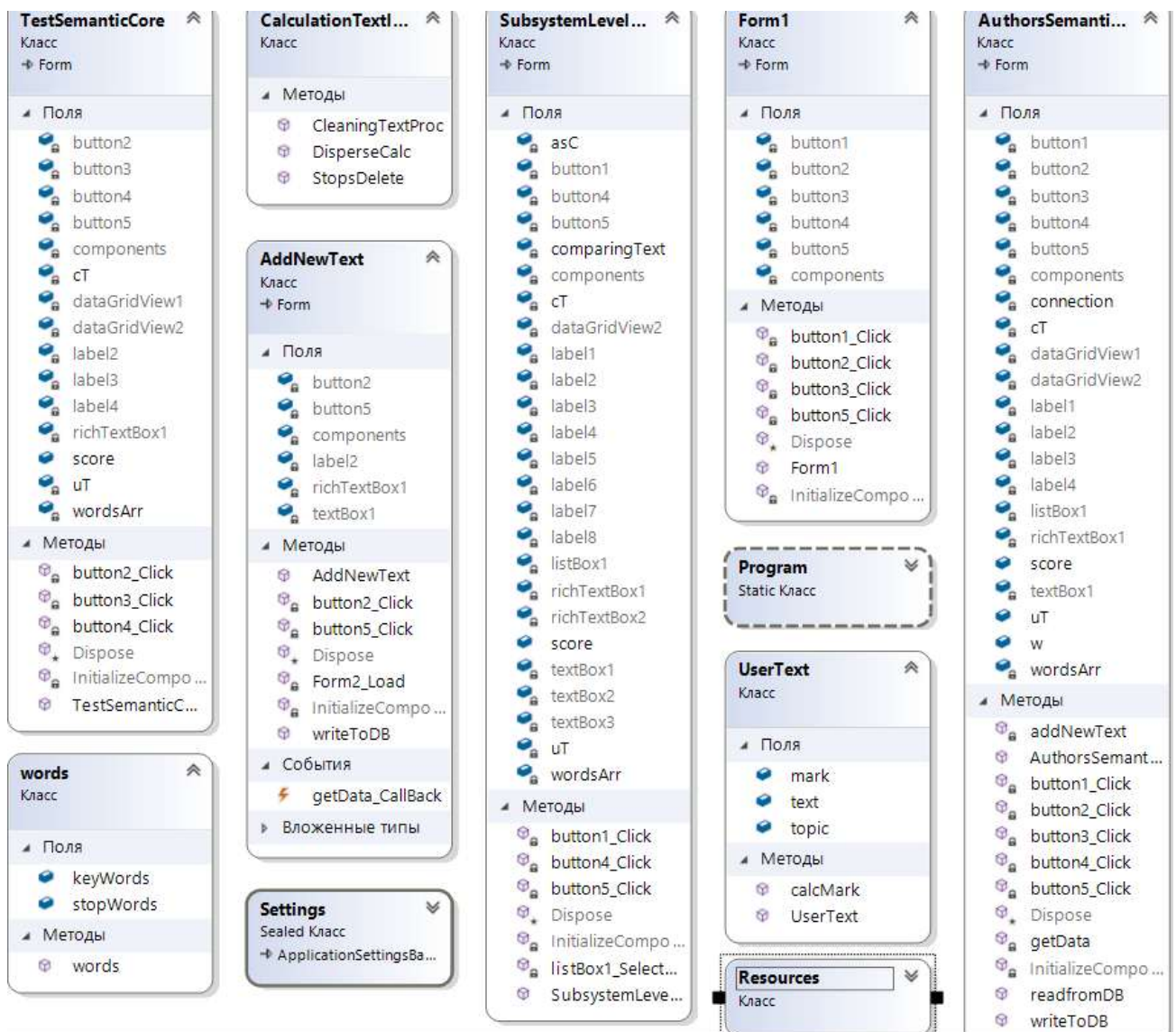


Рисунок 3.1 – Діаграма класів

У класі «words» зберігається інформація про слова текстів. Використовується як для аналізу користувацького тексту так і авторського. Містить відповідні переліки ключових слів та стоп-слів.

Клас «UserText» містить інформацію про сам текст: тема, сам текст та його оцінка схожості на вказаний (оцінка обраховується методом calcMark()).

Клас «AddNewText» представляє собою форму для збереження нового тексту в базу авторських текстів. Для передачі даних у батьківську форму використовується подія для делегату «get_Data_Callback».

У класі «CalculationTextInformation» реалізовано основний функціонал по первинній обробці тексту – процедури очищення від розділових знаків та спец-символів, розбивка тексту на слова та видалення стоп-слів. Також обраховується дисперсна оцінка для кожного ключового слова у векторі.

Клас «SubsystemLevelOfPlagiat» призначена для виведення інформації про рівень плагіату шляхом дисперсного оцінювання семантичної важливості слів.

Функціонал класів «TestSemanticCore» та «AuthorsSemanticCore» дещо схожі між собою за винятком того що для «AuthorsSemanticCore» є зчитування та запис у базу даних текстів.

Особливості реалізації описаних класів більш детально описані в 3.2.

3.2 Особливості реалізації програмних складових системи

Під час програмної реалізації складових системи було розроблено три підсистеми: підсистема роботи з авторським семантичним ядром тексту, підсистема роботи з користувацьким семантичним ядром тексту та підсистема

визначення рівня запозичень. Під час розробки підсистеми роботи з авторським семантичним ядром тексту було розроблено ряд функцій, серед яких є функція додавання нового авторського тексту до бази `addNewText()`. Її програмний код зображено нижче:

```
private void addNewText() {  
    AddNewText addT = new AddNewText();  
    addT.Show();  
    addT.getData_CallBack += getData;  
}
```

Дана функція створює нову форму (рисунок 3.2) та підписує її на функцію зворотного виклику, оскільки є потреба отримати результат, який уведе користувач назад.

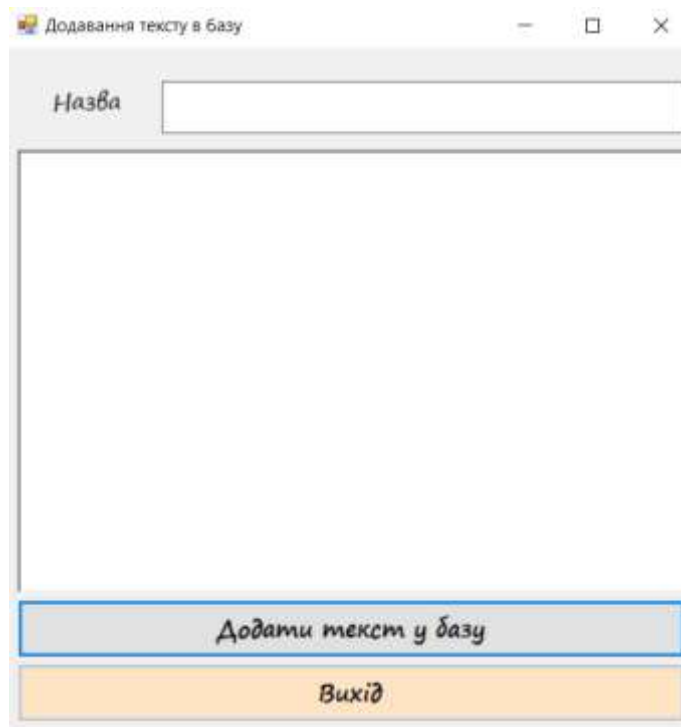


Рисунок 3.2 – Форма додавання авторського тексту до бази

Код методу `getData(String topic, String text)` проілюстровано нижче:

```
private void getData(String topic, String text)
{
    writeToDB(topic, text);
    readfromDB(connection);

    listBox1.Items.Clear();
    for (int i = 0; i < uT.Length; i++)
    {
        listBox1.Items.Add(uT[i].topic);
    }
}
```

Він містить частину запису до бази та зчитування даних, так як потрібно образу у формі користувача оновити інформацію для відображення нового доданого тексту. Відповідно, у класі з формою заповнення нового тексту та інформації про нього (рисунок 3.2) оголошено делегат та відповідну подію:

```
public delegate void callback_data(String topic, String text);
public event callback_data getData_CallBack;
```

При натисненні на кнопку «Додати текст у базу» викличеться функція зчитування даних та передача їх через `CallBack` на початку форму:

```
String text = richTextBox1.Text;
String topic = textBox1.Text;
```

```
getData_CallBack(topic, text);
```

```
MessageBox.Show("Текст додано у базу");
```

Доданий текст відображається у користувача та виводиться відповідне повідомлення що текст додано (рисунок 3.3).

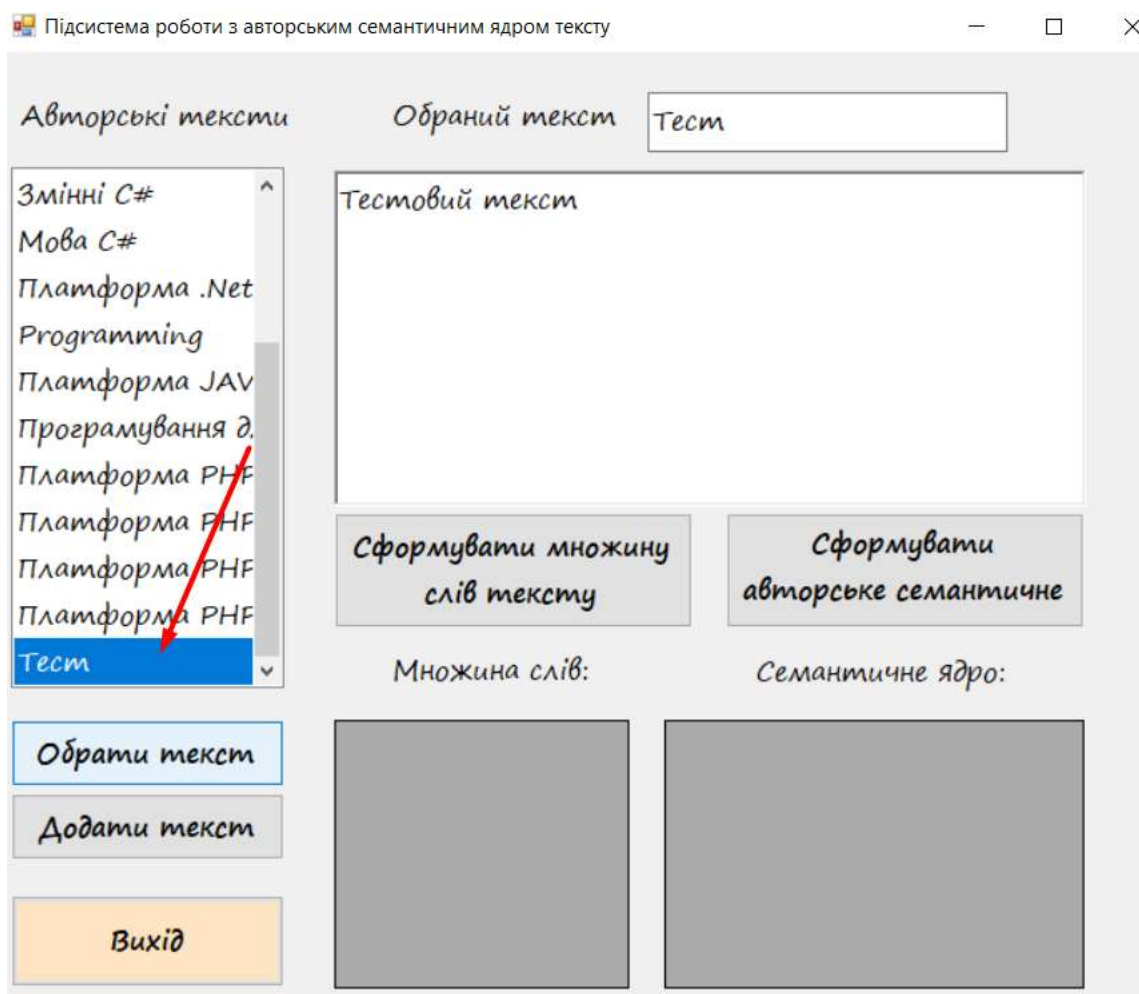


Рисунок 3.3 – Відображення доданого тексту до бази на формі

Для реалізації видалення стоп-слів було реалізовано метод `StopsDelete()`, який видаляє всі слова довжина яких менше двох та видаляє всі вказані стоп-слова з користувацького переліку. Метод повертає перерахування слів без стоп-слів. Код методу подано нижче:

```

public IEnumerable<String> StopsDelete(IEnumerable<String> t)
{
    IEnumerable<String> words = t;

    words = words.Where(x => x.Length > 2);

    foreach (var w1 in AuthorsSemanticCore.w.stopWords) {
        words = words.Where(x => !x.Equals(w1));
    }

    return words;
}

```

Таким чином, було реалізовано програмні компоненти системи автоматизованого визначення рівня запозичень цифрових текстів шляхом дисперсійної оцінки.

3.3 Тестування інформаційної системи

Для перевірки можливості створеного програмного продукту виконувати поставлені у 1.5 функції було виконано тестування ПЗ шляхом створення Unit Test, Test Case та проведено функціональне тестування.

Перший метод, який буде протестовано – метод очистки від спеціальних символів та розбиття тексту на слова. Код тестового методу подано нижче:

```
[TestMethod()]
```

```

public void CleaningTextProcTestTrue()
{
    CalculationTextInformation c = new
CalculationTextInformation();
    var text = "Python (найчастіше вживане прочитання – «Пайтон»,
запозичено назву[7] " +
        "з британського шоу Монті Пайтон)" +
        " – інтерпретована " +
        "мова програмування високого рівня зі строгою динамічною
типізацією.";

    var result = new List<string>()
    { "python", "найчастіше","вживане", "прочитання", "пайтон",
"запозичено","назву","7",
        "з","британського","шоу", "монті","пайтон","інтерпретована",
"мова","програмування","високого","рівня","зі","строгою","динамічною","тип
ізацією"};

    var wordsrez = c.CleaningTextProc(text);
    Assert.IsTrue(result.SequenceEqual(wordsrez));
}
}

```

Після запуску тесту на виконання з'явиться результат виконання – пройдено. Результат проілюстровано на рисунку 3.4

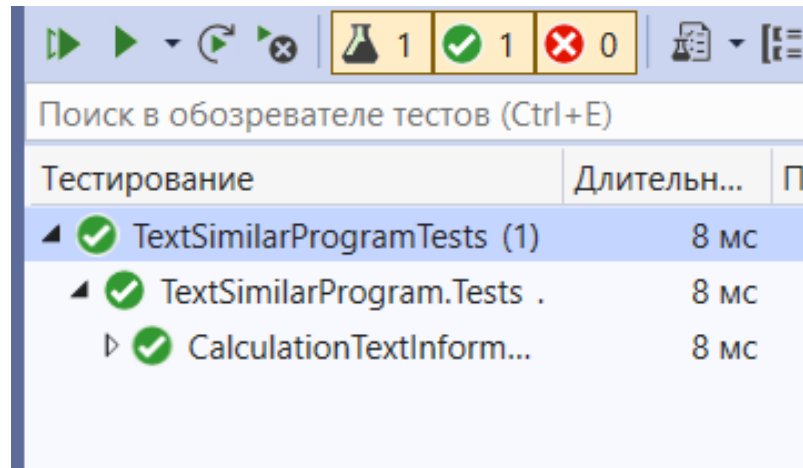


Рисунок 3.4 – Успішне проходження тесту на очищення тексту від спеціальних символів та розбиття на слова

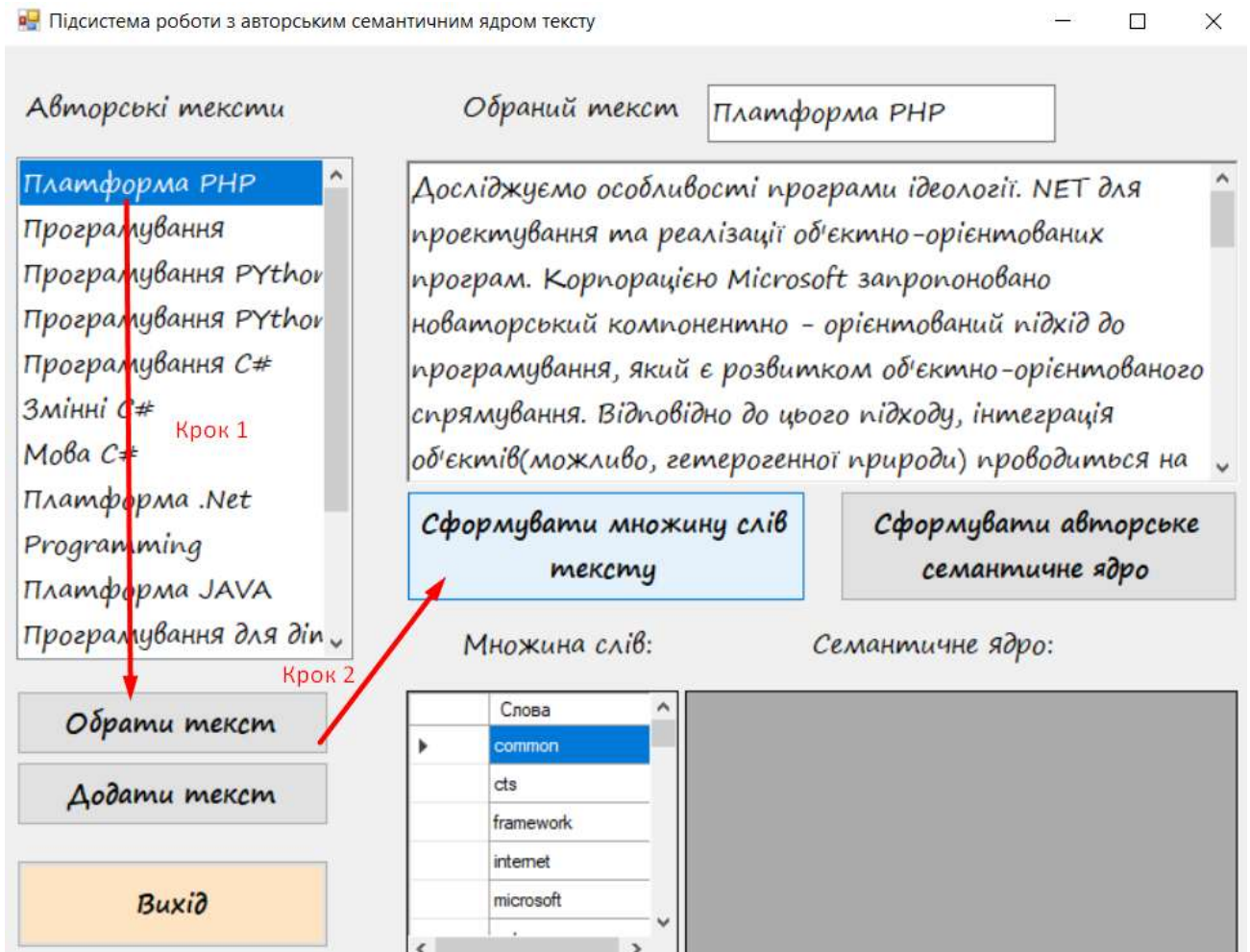


Рисунок 3.5 – Відображення у програмі очищення тексту від спеціальних символів та розбиття на слова

При запуску програми обравши першу форму «Робота з авторськими текстами» у лівій колонці можна обрати будь-який із текстів та при натисканні кнопки «Обрати текст» обраний текст буде відображено у правому верхньому вікні. При натисненні кнопки «Сформувати множину слів тексту» буде сформовано множину слів без спеціальних символів та відображено на екрані користувача (рисунок 3.5)

Наступним Unit-тестом буде тестовий метод на видалення стоп-слів, у які входять усі слова довжиною до 3-х літер не включно та видалення слів зі списку стоп-слів, заданого користувачем. Код тестового методу буде таким:

```
[TestMethod()]
public void StopsDeleteTestRezTrue()
{
    CalculationTextInformation c = new
CalculationTextInformation();

    var test = new List<string>()
    { "python", "найчастіше","вживане", "прочитання", "пайтон",
"запозичено","назву","7",
    "з","британського","шоу", "монті","пайтон","інтерпретована",
"мова","програмування","високого","рівня","зі","строгою","динамічною","тип
ізацією"};

    var rez = new List<string>()
    { "python","вживане", "прочитання", "пайтон",
"запозичено","назву","британського","шоу",
"монті","пайтон","інтерпретована",
```

```
"мова", "програмування", "високого", "рівня", "строгою", "динамічною", "типізаці  
єю"};
```

```
var wordsrez = c.StopsDelete(test);  
Assert.IsTrue(rez.SequenceEqual(wordsrez));  
}
```

Після запуску тестового методу на виконання з'явиться результат виконання – пройдено. Результат проілюстровано на рисунку 3.6

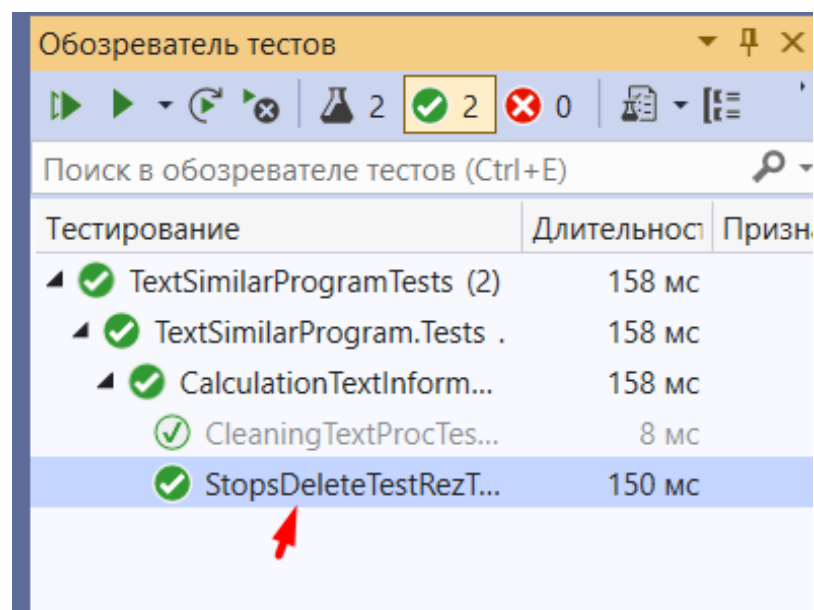


Рисунок 3.6 – Відображення у програмі очищення тексту від стоп-слів

Також для тестування програмного забезпечення було проведено тестування методом тест-кейсів. Перший тестовий випадок – перевірка на збереження авторського тексту у базі та відображення доданого тексту на екрані користувача (таблиця 3.1).

Таблиця 3.1 Тест-кейс ТК0001

Тест-кейс ID: ТК0001	Приоритет: 1	Створено: 5.05.2022, Раєва А.В.
<p>Назва: Перевірка збереження авторського тексту у базі та відображення доданого тексту на екрані користувача.</p> <p>Вхідні дані: Додати текст із назвою «Тестовий текст» та текстом: «Цей текст додано як тестовий».</p>		
Кроки		Очікуваний результат
<ol style="list-style-type: none"> 1. Запустити програму 2. На формі обрати пункт «Підсистема роботи з авторським семантичним ядром» 3. Натиснути кнопку «Додати текст» 4. У поле «Назва» вписати назву «Тестовий текст» 5. У поле тексту помістити текст «Цей текст додано як тестовий» 6. Натиснути кнопку «Додати текст у базу» 7. Закрити форму додавання нового тексту 8. Перейти до «Підсистеми роботи з авторським семантичним ядром» та пролистати список «Авторські тексти» до кінця. 9. Перевірити наявність доданого тексту. 		<p>Доданий текст успішно відображається в переліку</p>
<p>Результат виконання тест-кейсу: перевірку пройдено успішно.</p>		

Після запуску програми необхідно виконати кроки, вказані у таблиці 3.1. Після чого у програмі користувач побачить результат про успішне додавання тексту до бази та відображення тексту у переліку (рисунок 3.8).

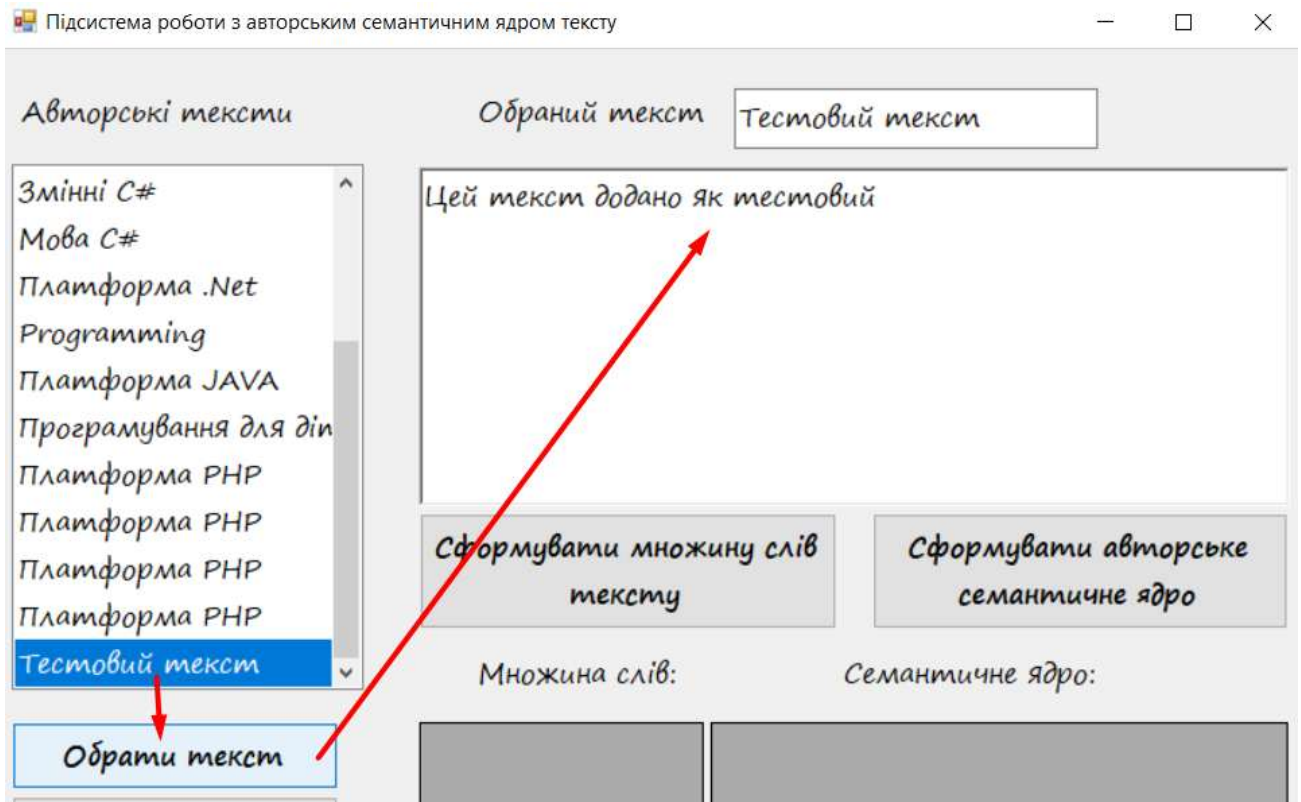


Рисунок 3.7 – Додавання тексту до авторських

Другим тестовим випадком буде формування множини слів для підсистеми роботи з тестовим семантичним ядром (таблиця 3.2).

Після запуску програми необхідно виконати кроки, вказані у таблиці 3.2. Після чого у програмі користувач побачить побудовану множину слів (рисунок 3.8).

Таблиця 3.2 Тест-кейс ТК0002

Тест-кейс ID: ТК0002	Пріоритет: 1	Створено: 6.05.2022, Раєва А.В.
<p>Назва: Формування множини слів</p> <p>Вхідні дані: В текстове поле «Текст для порівняння» ввести: «З цього тексту буде побудовано множину слів»</p>		
Кроки		Очікуваний результат
<ol style="list-style-type: none"> 1. Запустити програму 2. На формі обрати пункт «Підсистема роботи з тестовим семантичним ядром» 3. В текстове поле «Текст для порівняння» ввести: «З цього тексту буде побудовано множину слів» 4. Натиснути кнопку «Сформувати множину слів тексту» 5. Перевірити результат 		<p>Множина слів відображена у таблиці «Множина слів».</p>
<p>Результат виконання тест-кейсу: перевірку пройдено успішно.</p>		

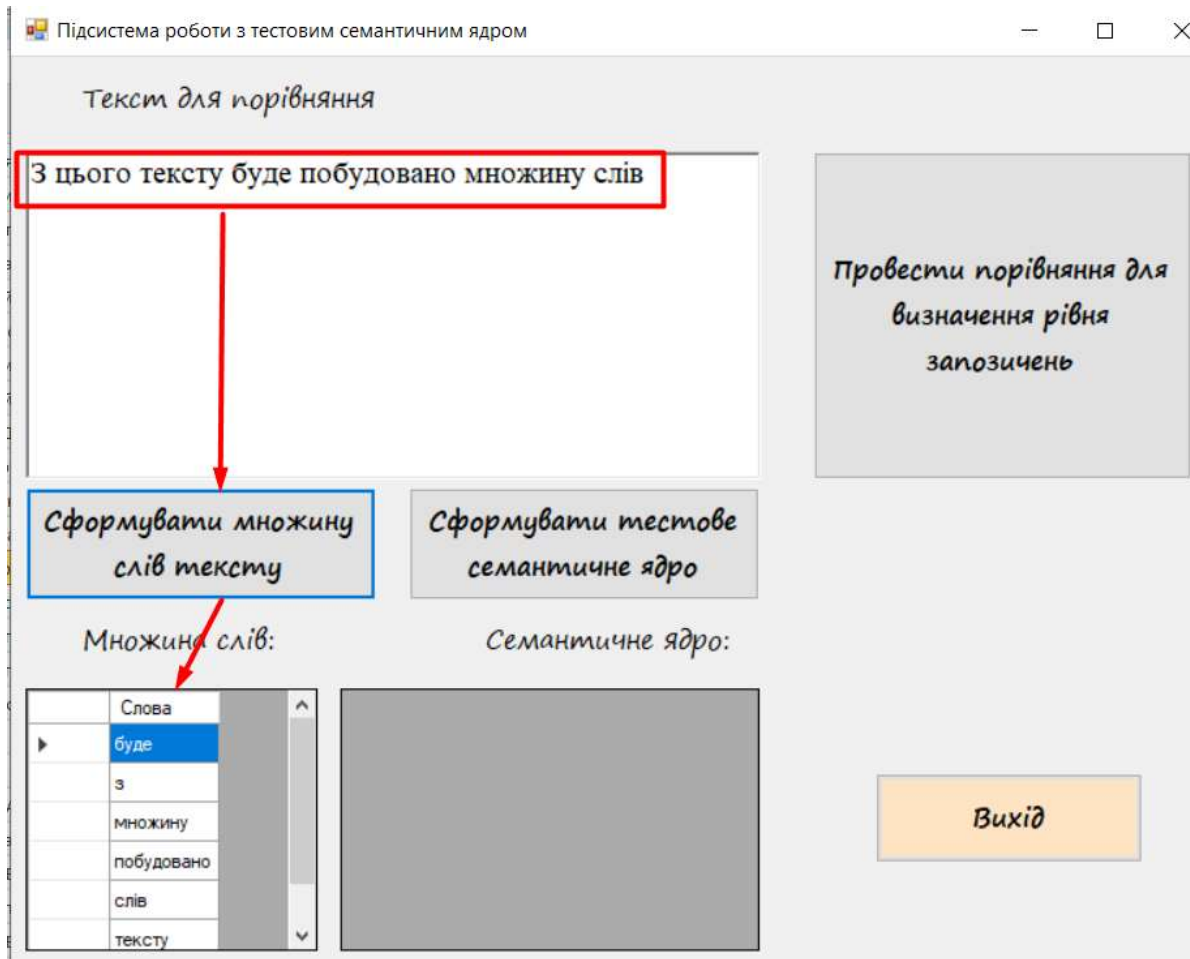


Рисунок 3.8 – Відображення множини слів

При тестуванні створеного програмного продукту некоректно працюючих функцій виявлено не було. У результаті проведеного тестування можна зробити висновок, що застосування працює коректно згідно поставленої задачі.

3.4 Інструкція користувача

Для зручності використання створеного програмного забезпечення було створено інструкцію користувача. При запуску застосування користувач бачить стартовий екран (рисунок 3.9)

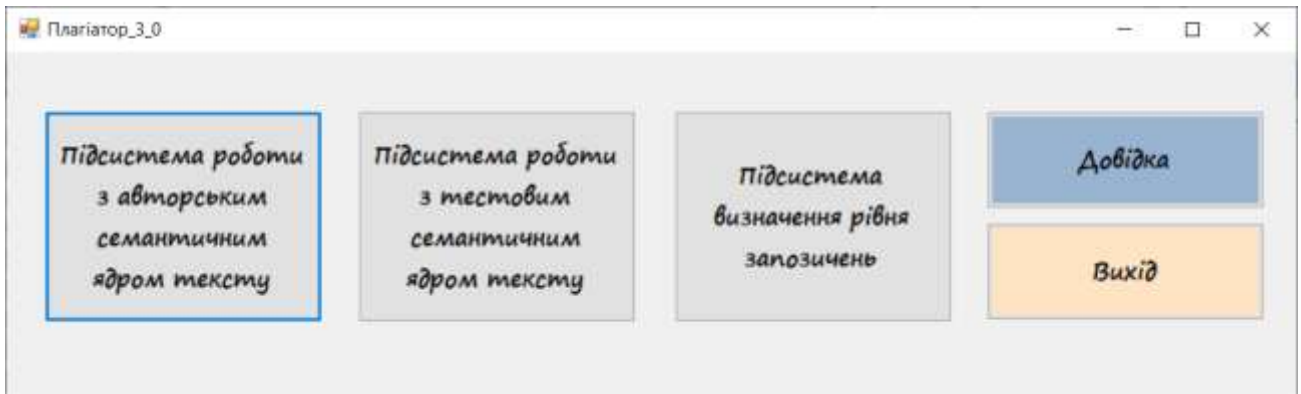


Рисунок 3.9 – Стартовий екран програми

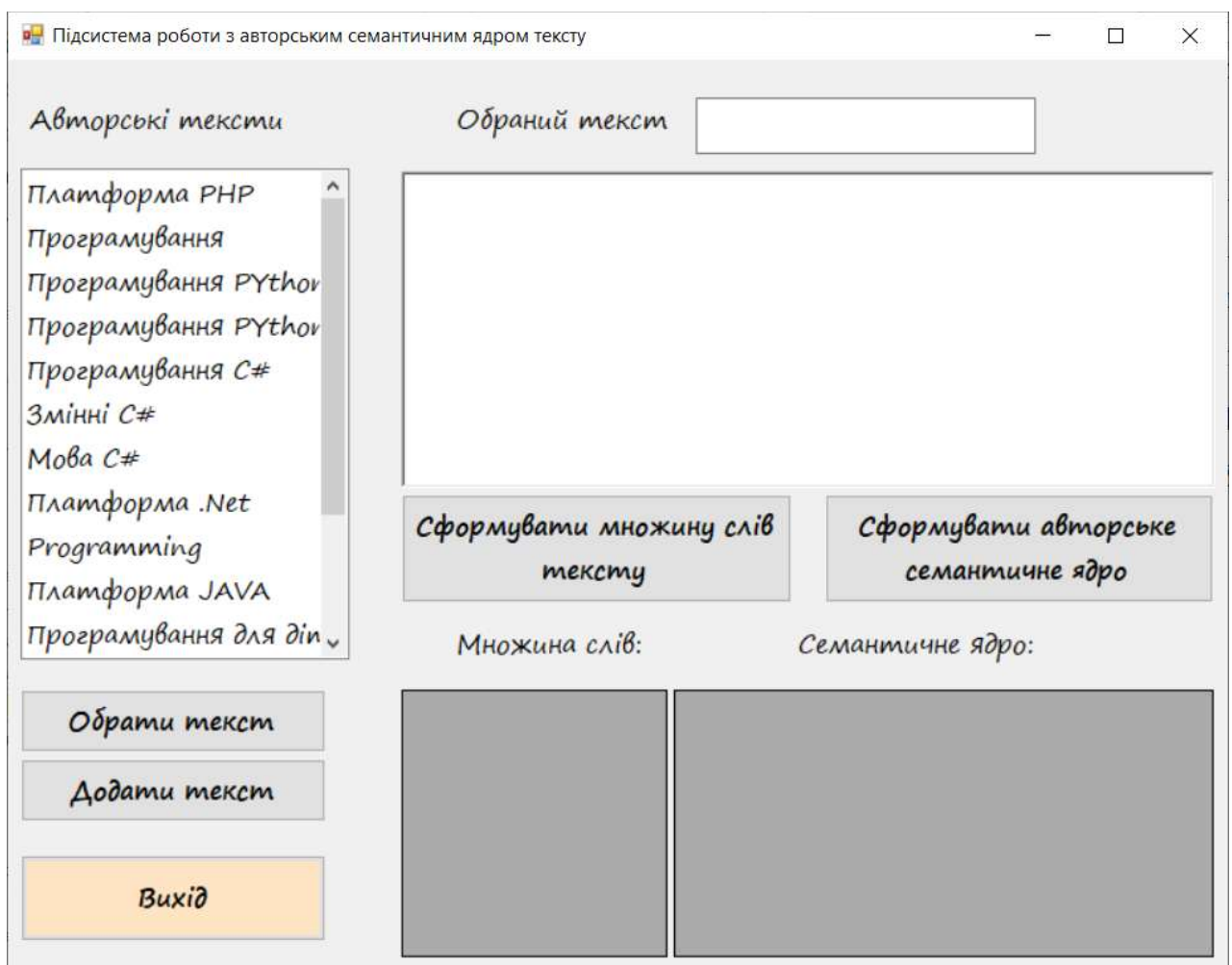


Рисунок 3.10 – Екран підсистеми роботи з авторським семантичним ядром

Зі стартового екрану програми є можливість переходів на такі підсистеми: «Підсистема роботи з авторським семантичним ядром тексту», «Підсистема роботи з тестовим семантичним ядром тексту», «Підсистема визначення рівня

запозичень». Також з головної форми доступні кнопки «Довідка» та «Вихід». При натисненні на кнопку «Підсистема роботи з авторським семантичним ядром тексту» користувач перейде на відповідну форму (рисунок 3.10).

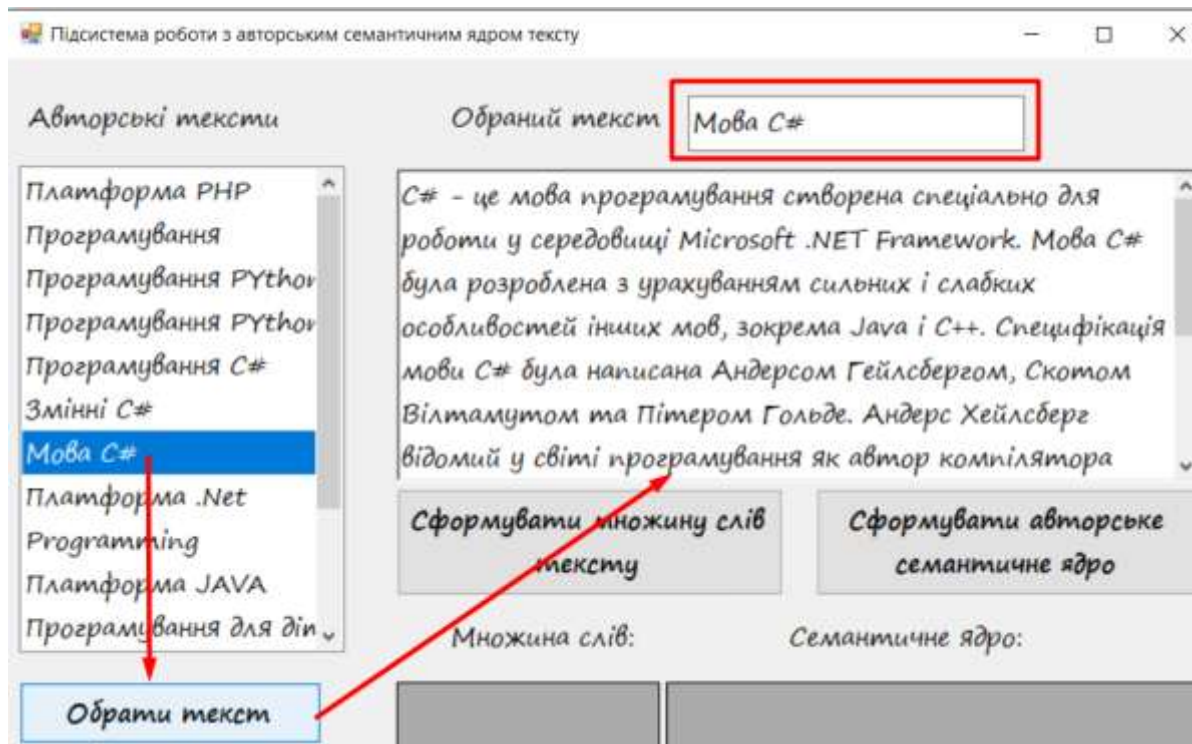


Рисунок 3.11 – «Підсистема роботи з авторським семантичним ядром».

Підсистема роботи з авторським семантичним ядром

У списку «Авторські тексти» можна обрати текст для аналізу. Для цього на потрібний текст потрібно натиснути мишею та натиснути кнопку «Обрати текст». Після чого обраний текст відобразиться у відповідних полях (рисунок 3.11).

Для формування множини слів тексту потрібно натиснути на кнопку «Сформувати множину слів тексту», а для формування авторського семантичного ядра потрібно натиснути на кнопку «Сформувати авторське семантичне ядро» (рисунок 3.12).

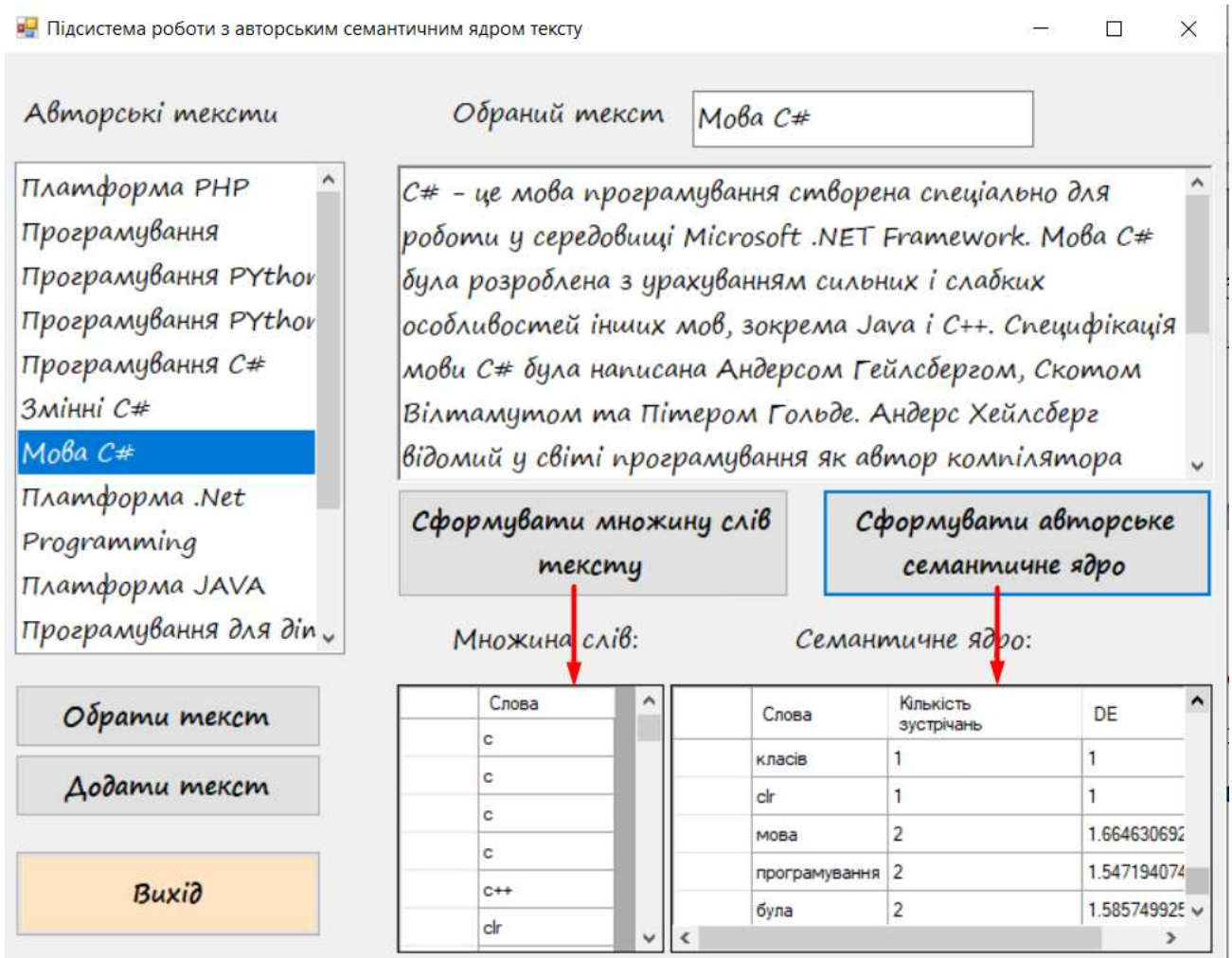


Рисунок 3.12 – «Підсистема роботи з авторським семантичним ядром».

Відображення множини слів та семантичного ядра

Для додавання нового тексту до бази авторських потрібно натиснути на кнопку «Додати текст», після чого відкриється форма для додавання нового тексту з відповідним параметрами (рисунок 3.13). Після вказування назви та самого тексту потрібно натиснути кнопку «Додати текст у базу». Для виходу з цієї форми потрібно натиснути кнопку «Вихід».

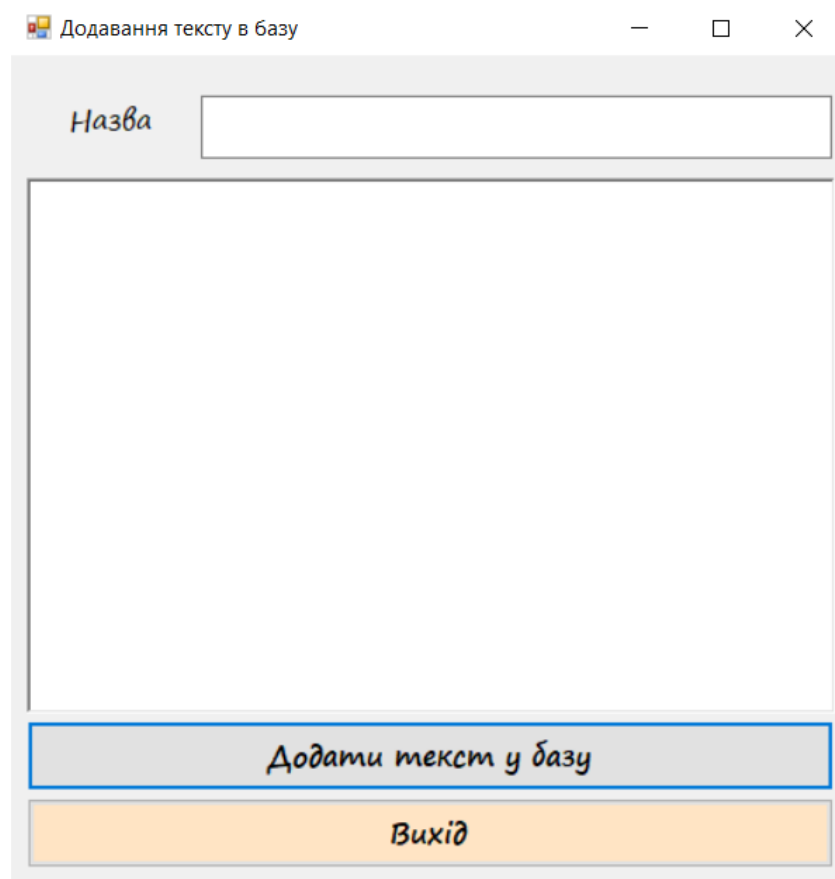


Рисунок 3.13 – «Підсистема роботи з авторським семантичним ядром».

Додавання нового тексту

При обиранні з головної форми «Підсистема роботи з тестовим семантичним ядром тексту» відкриється форма роботи з тестовим текстом (рисунок 3.14).

З даної форми можна виконати аналіз користувацького тексту, сформуванати для нього множину слів тексту та семантичне ядро. Відповідно, для цього потрібно вставити текст який цікавить користувача для аналізу у поле «Текст для порівняння» та натиснути кнопки «Сформуванати множину слів тексту» та «Сформуванати тестове семантичне ядро» (рисунок 3.15).

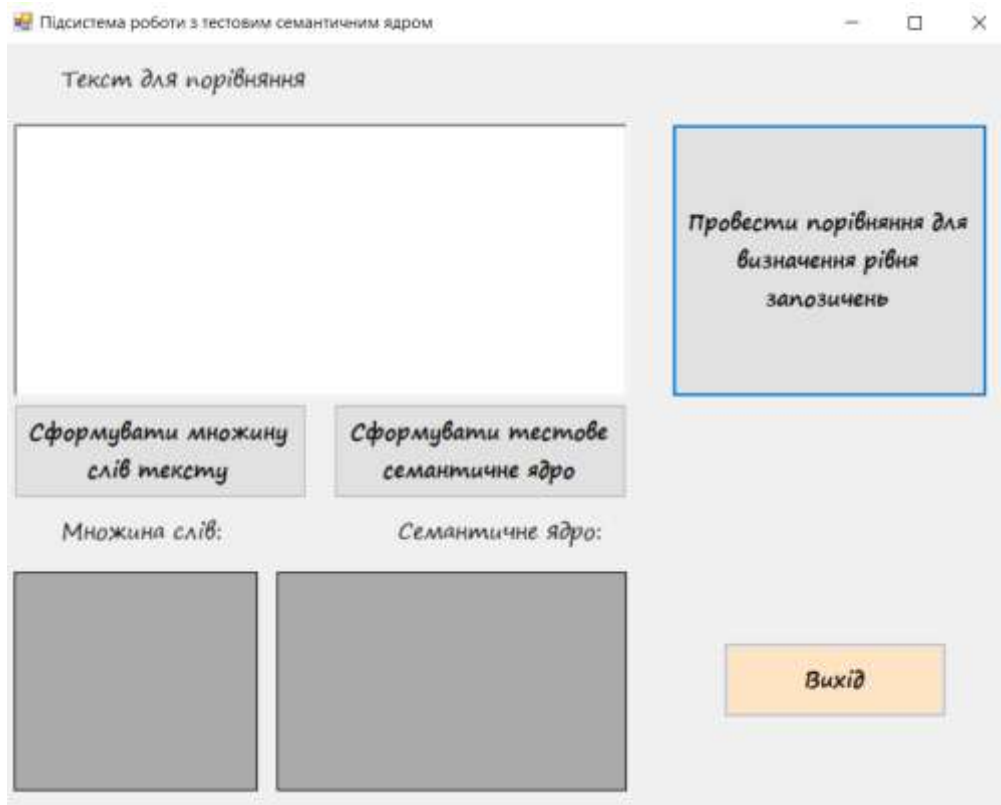


Рисунок 3.14 – «Підсистема роботи з тестовим семантичним ядром»

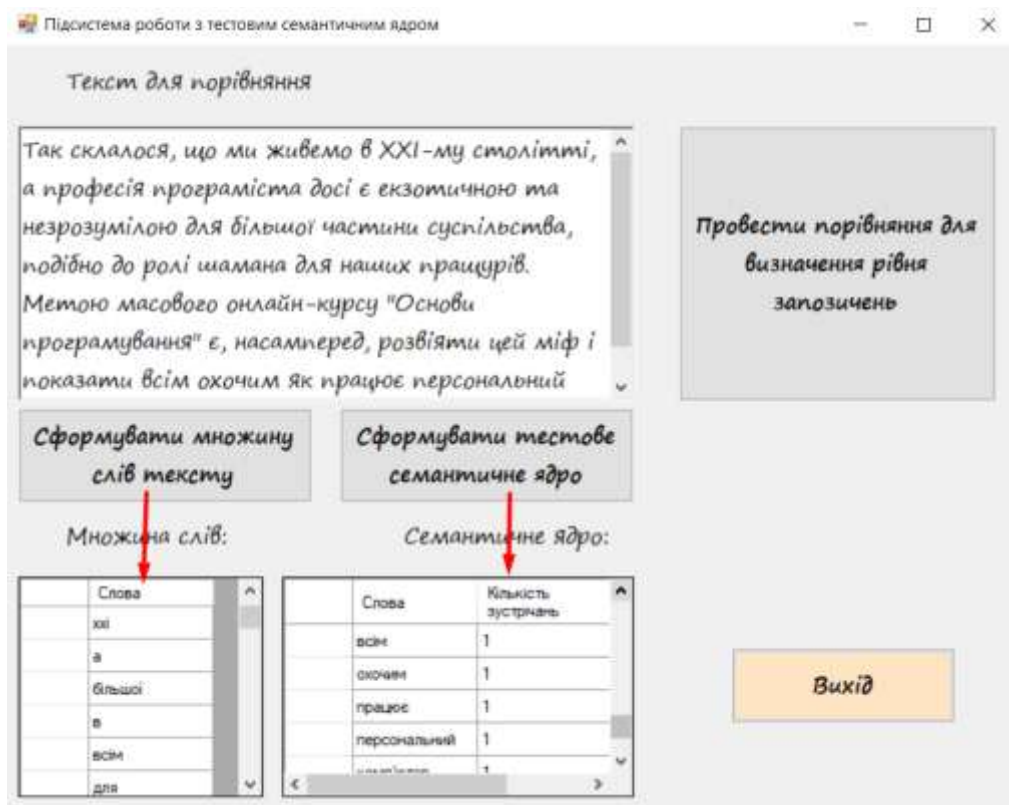


Рисунок 3.15 – «Підсистема роботи з тестовим семантичним ядром».

Формування множини слів та семантичного ядра

Якщо натиснути на кнопку «Провести порівняння для визначення рівня запозичень», користувач перейде на форму підсистеми «Підсистема визначення рівня запозичень» (рисунок 3.16), де вже буде завантажено користувацький текст для аналізу.

Рисунок 3.16 – «Підсистема визначення рівня запозичень»

Також на цю форму користувач може перейти безпосередньо з головного меню (рисунок 3.9), натиснувши кнопку «Підсистема визначення рівня запозичень». У такому випадку поле «Текст для порівняння» буде порожнім (рисунок 3.17).

Підсистема для визначення рівня запозичень

Текст для порівняння:

Показати рівень запозичень

Загальний висновок: Максимальний збіг семантичного ядра становить %

з текстом:

Найбільш схожі тексти: Текст: Семантичне ядро:

Показати обраний текст

Вихід

Рисунок 3.17 – «Підсистема визначення рівня запозичень». Стартовий екран

Для визначення рівня плагіату потрібно завантажити текст який цікавить користувача та натиснути кнопку «Показати рівень запозичень». Після чого буде сформовано відповідний висновок (рисунок 3.18)

Також формується список із 5-и текстів, найбільш схожих на тестовий. Щоб проглянути семантичне ядро та сам текст, на який видало плагіат потрібно обрати з переліку текст, який цікавить користувача та натиснути кнопку «Показати обраний текст». Після цих дій буде сформовано відповідне семантичне ядро (рисунок 3.19).

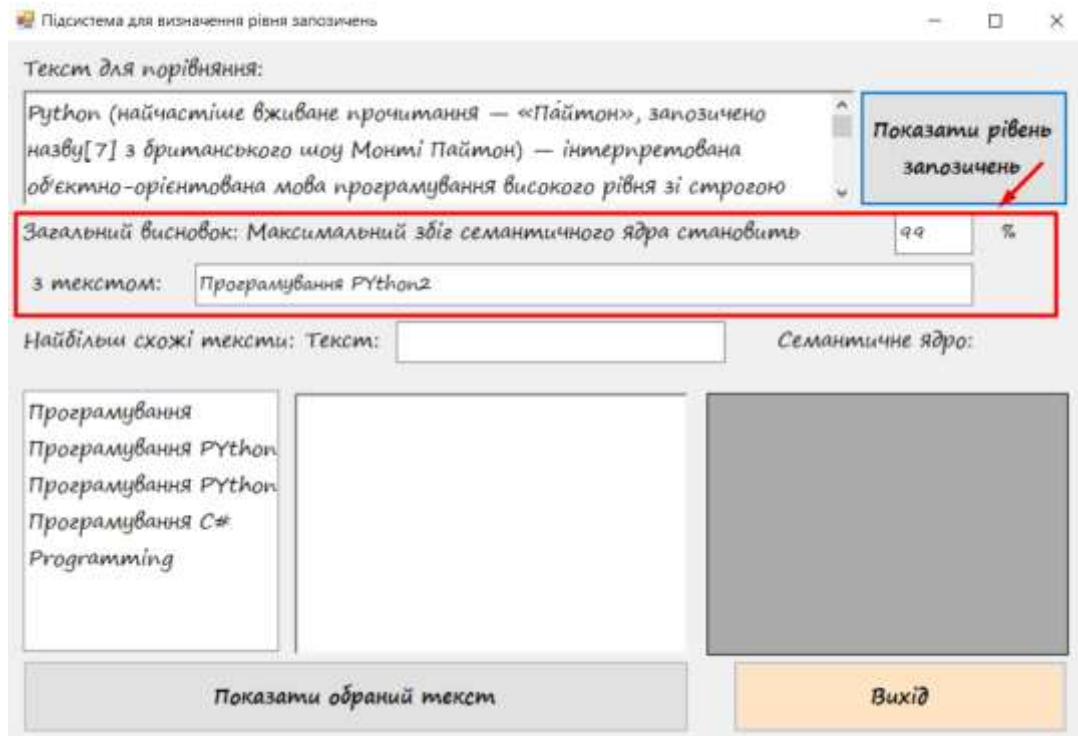


Рисунок 3.18 – «Підсистема визначення рівня запозичень». Показ рівня запозичень

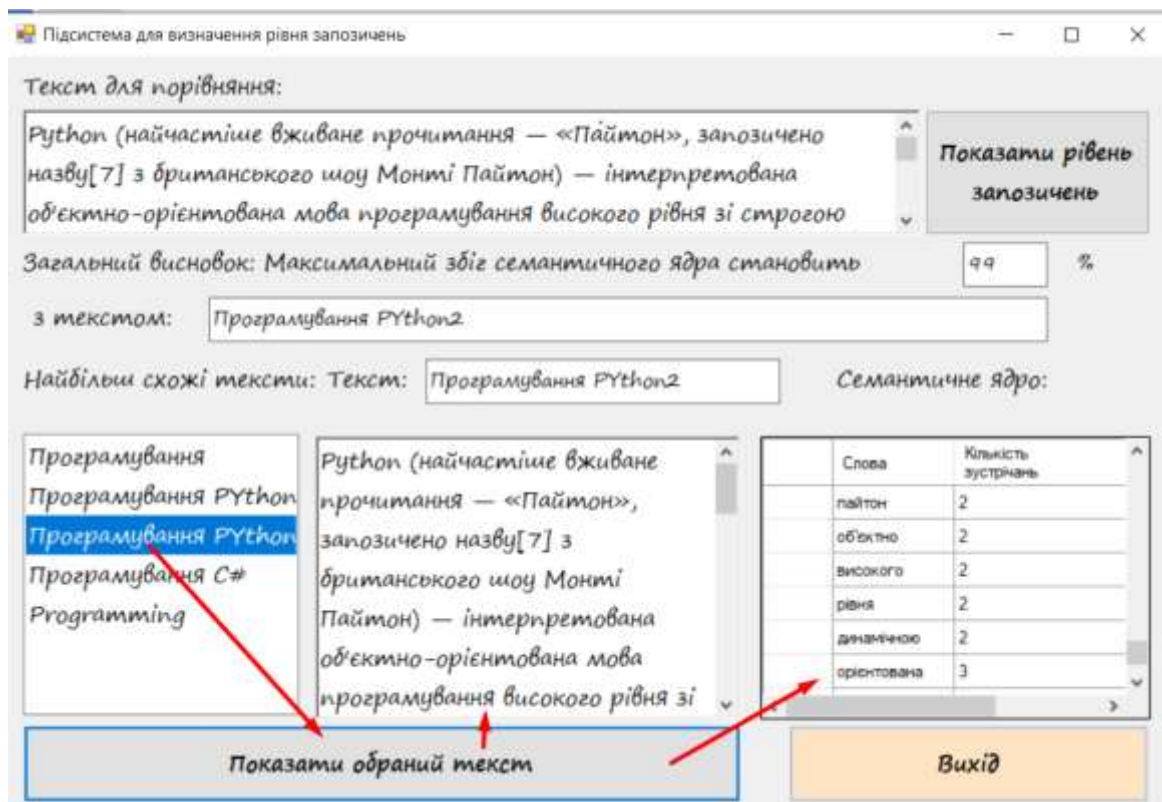


Рисунок 3.19 – «Підсистема визначення рівня запозичень». Відображення семантичного ядра тексту, який видало як плагіат

Таким чином, створену програмну реалізацію можна використовувати не тільки у цілях виявлення плагіату, а й у цілях проведення семантичних досліджень.

Висновки

Відповідно до мети КРБ, було розроблено метод автоматизованого визначення семантичної подібності цифрових текстів для виявлення запозичень, а також його програмну реалізацію. Для прикладного тестування методу було виконано його програмну реалізацію у вигляді застосування на платформі .NET, що виконує наступні основні функції:

- створення послідовного вектору слів для кожного цифрового тексту авторського базової вибірки;
- створення послідовного вектору слів для тестового тексту;
- створення частотного вектору семантичного ядра для кожного авторського цифрового тексту;
- створення частотного вектору семантичного ядра тестового цифрового тексту;
- очищення частотних векторів семантичних ядер текстів від стоп-слів;
- обрахунок оцінок рівня запозичень тестового цифрового тексту щодо кожного цифрового тексту базової вибірки;
- сортування результуючого переліку за оцінкою рівня запозичень та формування висновку щодо визначення семантичної подібності цифрових текстів.

Отже, для досягнення мети були виконані усі поставлені завдання:

1. Проведено аналіз предметної області, а саме: аналіз інформаційних моделей, виконано огляд теоретичних аспектів до розв'язку задачі автоматизованого визначення семантичної подібності цифрових текстів для

виявлення запозичень, виконати аналіз існуючих програмних рішень та виконано аналіз сучасних засобів створення програмного забезпечення.

2. Виконано проєктування інформаційної системи, у рамках якого було створено метод автоматизованого визначення семантичної подібності цифрових текстів для виявлення запозичень; побудовано інформаційну структуру системи; обрано у якості засобів розробки системи на основі методу автоматизованого визначення семантичної подібності цифрових текстів для виявлення запозичень платформи .NET, мови програмування C#, СКБД MS SQL Server та редактор програмного коду VisualStudio.

3. Створено програмну реалізацію інформаційної системи на базі методу автоматизованого визначення семантичної подібності цифрових текстів для виявлення запозичень та проведено апробацію результатів шляхом програмного та функціонального тестування.

4. Для зручності користування інформаційною системою на базі методу автоматизованого визначення семантичної подібності цифрових текстів для виявлення запозичень було створено інструкцію користувача.

Перелік посилань

1. Wikipedia. Результат інтелектуальної діяльності. URL: https://uk.wikipedia.org/Результат_інтелектуальної_діяльності
2. WikiLegalAid. Авторське право. URL: https://wiki.legalaid.gov.ua/index.php/Авторське_право
3. Науковий відділ Львівського національного медичного університету імені Данила Галицького. Плагіат і академічний плагіат: причини виникнення та шляхи подолання. URL: <http://nauka.meduniv.lviv.ua/wp-content/uploads/simple-file-list/Academ-dobrochesnist/plagiat-i-akademichnyj-plagiat.pdf>
4. J. B. Williams The Plagiarism Problem: Are Students Entirely To Blame? 2022. URL: https://www.researchgate.net/publication/27481233_The_Plagiarism_Problem_Are_Students_Entirely_To_Blame
5. M. Knight The Problem of Plagiarism: Finding Ways to Promote Academic Integrity in the Classroom. Business and Professional Communication Quarterly 2018, Vol. 81(3) 279–280 cc. URL: <https://journals.sagepub.com/doi/pdf/10.1177/2329490618794188>
6. Academic Writing Help Centre (AWHC). Plagiarism. URL: <https://sass.uottawa.ca/sites/sass.uottawa.ca/files/awhc-plagiarism.pdf>
7. Wikipedia. Плагіат. URL: <https://uk.wikipedia.org/Плагіат>
8. Кіллер атиплагіат. Обхід антиплагіата, 10 способів. URL: <https://killer-antiplagiat.ru/blog/obxod-antiplagiata-10-sposobov-obmanut-onlajn>
9. Wikipedia. Репозитарій. URL: <https://uk.wikipedia.org/Репозитарій>
10. Wikipedia. Національний репозитарій академічних текстів. URL: https://uk.wikipedia.org/Національний_репозитарій_академічних_текстів

11. Репозитарій Харківського національного медичного університету. Як запобігти академічного плагіату: стратегія науковця. URL: https://repo.knmu.edu.ua/bitstream/Павленко_Як%20запобігти%20АП_текст.pdf
12. Stanford Encyclopedia of Philosophy. Computational Linguistics. URL: <https://plato.stanford.edu/entries/computational-linguistics/>
13. Wikipedia. Semantic analysis (computational). URL: [https://en.wikipedia.org/wiki/Semantic_analysis_\(computational\)](https://en.wikipedia.org/wiki/Semantic_analysis_(computational))
14. Сайт для філологів-українців. Лінгвістичний аналіз тексту. URL: http://grinch-home.at.ua/publ/baza_referativ/ukrajinska_mova/lingvistichnij_analiz_tekstu/24-1-0-159
15. Expert AI. What Semantic Analysis Means to Natural Language Processing. URL: <https://www.expert.ai/blog/natural-language-process-semantic-analysis-definition/>
16. Wikipedia. Семантична мережа. URL: https://uk.wikipedia.org/Семантична_мережа
17. Analytics Steps. Semantic Analysis: Working and Techniques. URL: <https://www.analyticssteps.com/blogs/semantic-analysis-working-and-techniques>
18. Wikipedia. Ключове слово. URL: https://uk.wikipedia.org/Ключове_слово
19. Wikipedia. Letter frequency. URL: https://en.wikipedia.org/wiki/Letter_frequency
20. Wikipedia. TF-IDF. URL: <https://uk.wikipedia.org/wiki/TF-IDF>
21. Ряба А.О., Мазурець О.В. Різновиди методу пошуку ключових слів у цифрових текстах за дисперсійним оцінюванням. Актуальні проблеми комп'ютерних технологій 2019 : зб. наук. пр. за матеріалами дванадцятої міжнар. наук.-техн. конф. Хмельницький, 2020. С. 169–172.

22. Unicheck. Широкий спектр зрозумілих функцій. URL: <https://unicheck.com/uk-ua/plagiarism-tool-for-personal-use>
23. Rigorous Themes. Unicheck vs Turnitin – Which Is Better? URL: <https://rigorousthemes.com/blog/unicheck-vs-turnitin-which-is-better/>
24. Наукова бібліотека Національного юридичного університету імені Ярослава Мудрого. Плагіат у студентських роботах: методи виявлення та запобігання. URL: https://library.nlu.edu.ua/POLN_TEXT/POSIBNIKI_2013/Methodichka-plagiat.pdf
25. Wikipedia. Turnitin. URL: <https://en.wikipedia.org/wiki/Turnitin>
26. Wikipedia. Copyscape. URL: <https://en.wikipedia.org/wiki/Copyscape>
27. Laptrinhx. Is Copyscape Premium Still Worth it? We Tried 7 Online Plagiarism Checker Tools. URL: <https://laptrinhx.com/is-copyscape-premium-still-worth-it-we-tried-7-online-plagiarism-checker-tools-2149378131/>
28. EduBirdie. Best Plagiarism Checkers in 2022. URL: <https://edubirdie.com/blog/best-plagiarism-checker#EduBirdie>
29. Plagium. Plagium helps you to ensure the originality of a text by detecting and identifying possible plagiarisms. URL: <https://www.plagium.com/en/plagiarismchecker>
30. GetApp. Plagium. URL: <https://www.getapp.co.uk/software/2037427/plagium>
31. Ayoub Ali M. Saeed, Alaa Yaseen Taqa A proposed approach for plagiarism detection in Article documents. Sinkron : Jurnal dan Penelitian Teknik Informatika Volume 7, Number 2, April 2022. URL: <https://www.polgan.ac.id/jurnal/index.php/sinkron/article/view/11381/898>
32. M. Davoodifard Automatic Detection of Plagiarism in Writing. studies in Applied Linguistics & TESOL at Teachers College, Columbia University, Vol. 21, No.

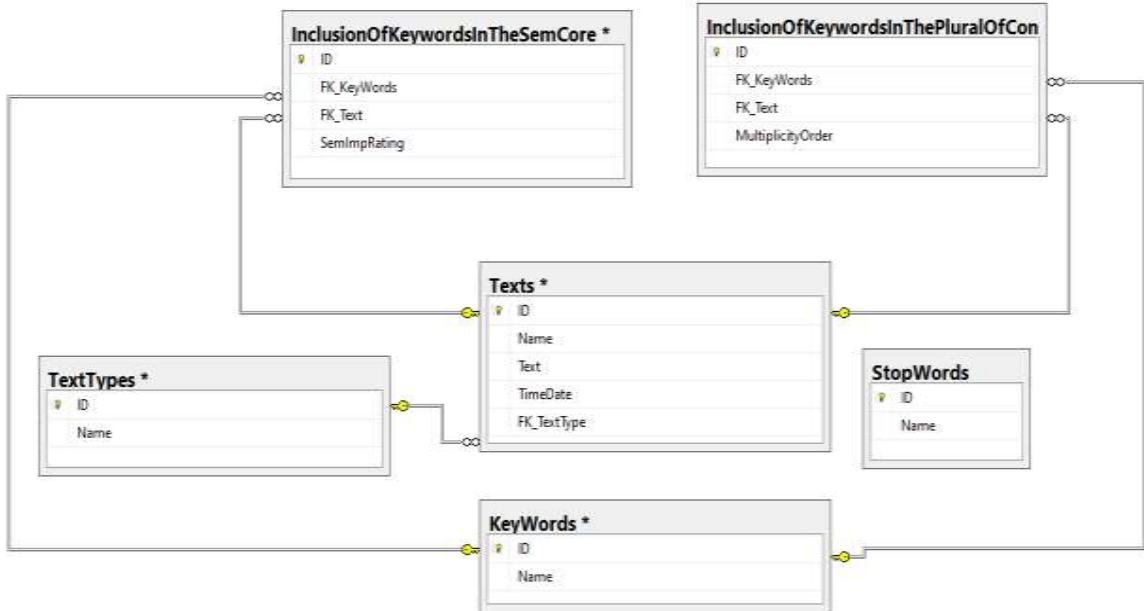
- 2, pp. 54–60. URL:
<https://journals.library.columbia.edu/index.php/SALT/article/view/9058/4651>
33. M. Chong, L. Specia, R. Mitkov Using Natural Language Processing for Automatic Detection of Plagiarism. URL: https://www.researchgate.net/profile/Lucia-Specia/publication/242783426_Using_Natural_Language_Processing_for_Automatic_Detection_of_Plagiarism/links/56179c2908ae0224ebce9956/Using-Natural-Language-Processing-for-Automatic-Detection-of-Plagiarism.pdf
34. The benefits of web-based application. URL:
<https://coresolutions.ca/blogs/core-web/the-benefits-of-web-based-applications>
35. Advantages and Disadvantages of Web Application | Drawbacks & Benefits of Web Application. URL: <https://www.hitechwhizz.com/2021/04/5-advantages-and-disadvantages-drawbacks-benefits-of-web-application.html>
36. What Is A Mobile Application? URL:
<https://www.emizentech.com/blog/what-is-a-mobile-application.html>
37. Advantages of Mobile Apps over responsive eCommerce websites. URL:
<https://anblicks-inc.medium.com/9-advantages-of-mobile-apps-over-responsive-e-commerce-websites-6aed1e6db0d8>
38. Top 9 Disadvantages of Mobile Apps – and What to Do Instead. URL:
<https://www.beezer.com/disadvantages-of-mobile-apps/>
39. Windows application. URL:
<https://encyclopedia2.thefreedictionary.com/Windows+application>
40. What are the Benefits of a Desktop Application? URL:
<https://opendesignct.com/what-are-the-benefits-of-a-desktop-application/>
41. Web Application Vs. Desktop Application: Pros and Cons. URL:
<https://positiwise.com/blog/web-application-vs-desktop-application-pros-and-cons/>

42. Wikipedia. Java. URL:
<https://uk.wikipedia.org/wiki/Java#%D0%9F%D0%BB%D0%B0%D1%82%D1%84%D0%BE%D1%80%D0%BC%D0%B0>
43. Wikipedia. Mono. URL: <https://uk.wikipedia.org/wiki/Mono>
44. Wikipedia. .NET_Framework. URL:
https://uk.wikipedia.org/wiki/.NET_Framework
45. Best Frameworks for Desktop Application Development. URL:
<https://dzone.com/articles/best-frameworks-for-desktop-application-developmen>
46. The C# Programming Language. URL:
<https://www.bairesdev.com/technologies/csharp/>
47. Wikipedia. Візуальне середовище програмування. URL:
https://uk.wikipedia.org/wiki/Візуальне_середовище_програмування
48. Benefits of Using Visual Studio Projects. URL:
<https://www.microfocus.com/documentation/visual-cobol/VC40/VS2015/GUID-802D1D7A-BF87-454F-9808-1227DA0607C9.html>
49. Microsoft SQL Server. URL:
https://uk.wikipedia.org/wiki/Microsoft_SQL_Server
50. MS SQL Server history and advantages. URL:
<https://bytescout.com/blog/2014/09/ms-sql-server-history-and-advantages.html>

ДОДАТКИ

Додаток А

Структура бази даних інформаційної системи автоматизованого визначення семантичної подібності цифрових текстів для виявлення рівня запозичень



Додаток Б

Розгорнута структура класів інформаційної системи автоматизованого визначення семантичної подібності цифрових текстів для виявлення рівня запозичень

The image displays a software development environment showing the class hierarchy for an information system. The interface is organized into several panels, each representing a different class or component. The main panels are:

- TestSemanticCore**: A class containing fields (button2-5, components, cT, dataGridViews, labels, richTextBox1, score, uT, wordsArr) and methods (button2_Click, button3_Click, button4_Click, Dispose, InitializeCompo..., TestSemanticC...).
- CalculationText...**: A class containing methods (CleaningTextProc, DisperseCalc, StopsDelete).
- AddNewText**: A class containing fields (button2-5, components, label2, richTextBox1, textBox1) and methods (AddNewText, button2_Click, button5_Click, Dispose, Form2_Load, InitializeCompo..., writeToDB, and an event getData_CallBack).
- SubsystemLevel...**: A class containing numerous fields (asC, button1-5, comparingText, components, cT, dataGridViews, labels, listBox1, richTextBox1-2, score, textBox1-3, uT, wordsArr) and methods (button1_Click, button4_Click, button5_Click, Dispose, InitializeCompo..., listBox1_Select..., SubsystemLeve...).
- Form1**: A class containing fields (button1-5, components) and methods (button1_Click, button2_Click, button3_Click, button5_Click, Dispose, Form1).
- AuthorsSemanti...**: A class containing fields (button1-5, components, connection, cT, dataGridViews, labels, listBox1, richTextBox1, score, textBox1, uT, w, wordsArr) and methods (addNewText, AuthorsSemant..., button1_Click, button2_Click, button3_Click, button4_Click, button5_Click, Dispose, getData, InitializeCompo..., readfromDB, writeToDB).
- words**: A class containing fields (keyWords, stopWords) and a method (words).
- Settings**: A Sealed Class containing ApplicationSettingsBa....
- Program**: A Static Class.
- UserText**: A class containing fields (mark, text, topic) and methods (calcMark, UserText).
- Resources**: A Class.

Додаток В

Презентаційний матеріал

КВАЛІФІКАЦІЙНА РОБОТА БАКАЛАВРА

МЕТОД АВТОМАТИЗОВАНОГО ВИЗНАЧЕННЯ СЕМАНТИЧНОЇ ПОДІБНОСТІ ЦИФРОВИХ ТЕКСТІВ ДЛЯ ВИЯВЛЕННЯ ЗАПОЗИЧЕНЬ



Виконав:
студентка 4 курсу, групи КН-18-1
Раєва Анна Вячеславівна



Керівник:
к.т.н., доцент кафедри КН
Мазурець Олександр Вікторович

Актуальність

З метою захисту авторських творів створено авторське право. Воно дозволяє авторам захистити свої твори від присвоєння іншими особами, або використання в інших цілях.

Для боротьби з такого роду шахрайством було створено спеціальні системи для боротьби з плагіатом. Суть таких систем полягає в тому, щоб виявити запозичення текстів у творах інших авторів. Для створення бази творів використовуються репозитарії.

Такі системи розроблюються на основі семантичного аналізу текстів, що дозволяє виявити навіть прихований плагіат.

Отже, можна зробити висновок про те, що задача автоматизованого визначення семантичної подібності цифрових текстів для виявлення запозичень є актуальною у всьому світі. Хоч і існує багато методів для розв'язку цієї задачі, проте майже всі вони потребують вдосконалення, адже шахраї вдаються все до нових методів обходу систем антиплагіату.

Мета і задачі роботи

Метою кваліфікаційної роботи бакалавра є розробка методу автоматизованого визначення семантичної подібності цифрових текстів для виявлення запозичень, а також його програмна реалізація.

Для прикладного тестування методу слід виконати його **програмну реалізацію** у вигляді додатку на платформі .NET, що виконує наступні основні функції:

- створення послідовного вектору слів для кожного цифрового тексту авторського базової вибірки;
- створення послідовного вектору слів для тестового тексту;
- створення частотного вектору семантичного ядра для кожного авторського цифрового тексту;
- створення частотного вектору семантичного ядра тестового цифрового тексту;
- очищення частотних векторів семантичних ядер текстів від стоп-слів;
- обрахунок оцінок рівня запозичень тестового цифрового тексту щодо кожного цифрового тексту базової вибірки;
- сортування результуючого переліку за оцінкою рівня запозичень та формування висновку щодо визначення семантичної подібності цифрових текстів.

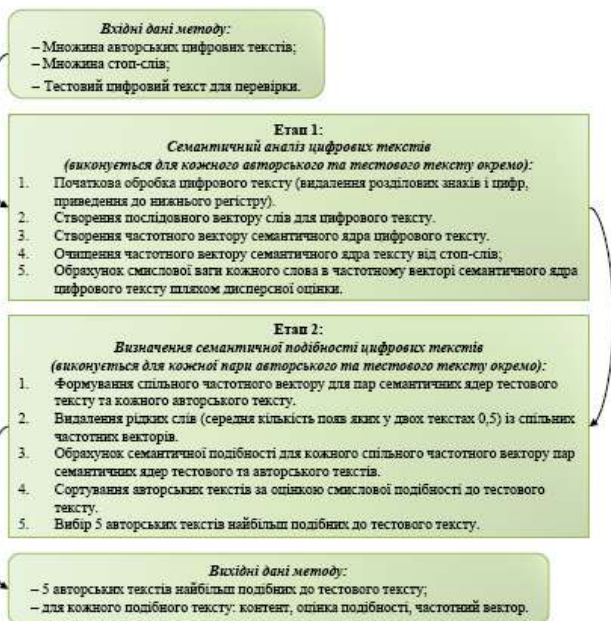
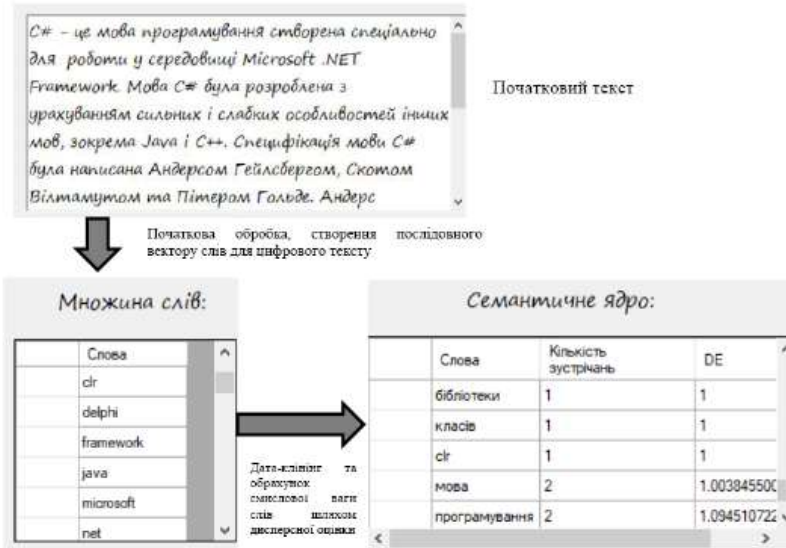


Схема методу
автоматизованого
визначення
семантичної подібності
цифрових текстів для
виявлення запозичень

Схема процесу семантичного аналізу цифрових текстів



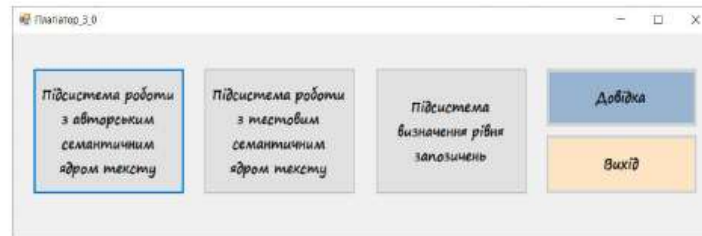
Оцінка семантичної подібності цифрових текстів

Оцінка семантичної подібності цифрових текстів S_{Semant} тестового тексту $Test$ до авторського тексту $Auth$ обраховується наступним чином:

$$S_{Semant} = \frac{\sum_{i=1}^n DE_{Dif}}{\sum_{i=1}^n DE_{Auth}},$$

де n – кількість слів у спільному частотному векторі для пар семантичних ядер тестового тексту та авторського тексту, DE_{Dif} – обрахована різниця важливості слова i у частотних векторах для семантичних ядер тестового тексту та авторського тексту, DE_{Auth} – оцінка семантичної важливості слова i у частотному векторі для семантичного ядра авторського тексту.

Стартовий екран програми



Відображення множини слів та семантичного ядра

Підсистема роботи з авторським семантичним ядром тексту

Авторські тексти

- Платформа PHP
- Програмування
- Програмування Python
- Програмування Python
- Програмування C#
- Мова C#**
- Платформа .Net
- Programming
- Платформа JAVA
- Програмування для дів

Обраний текст: Мова C#

C# - це мова програмування створена спеціально для роботи у середовищі Microsoft .NET Framework. Мова C# була розроблена з урахуванням сильних і слабких особливостей інших мов, зокрема Java і C++. Специфікація мови C# була написана Андерсом Гейлсбергом, Скотом Вілтанутом та Пітером Гольде. Андерс Хейлсберге відомий у світі програмування як автор компілятора

Сформувати множину слів тексту

Сформувати авторське семантичне ядро

Множина слів:

Слова
c
c
c
c
c++
сі

Семантичне ядро:

Слова	Кількість зустрічей	DE
класів	1	1
сі	1	1
мова	2	1.66463065
програмування	2	1.54719407
була	2	1.58574992

Обрати текст

Додати текст

Вихід

Формування множини слів та семантичного ядра

Підсистема роботи з тестовим семантичним ядром

Текст для порівняння

Так склалося, що ми живемо в XXI-му столітті, а професія програміста досі є екзотичною та незрозумілою для більшої частини суспільства, подібно до ролі шамана для наших пращурів. Метою масового онлайн-курсу "Основи програмування" є, насамперед, розвіяти цей міф і показати всім очочим як працює персональний

Провести порівняння для визначення рівня запозичень

Сформувати множину слів тексту

Сформувати тестове семантичне ядро

Множина слів:

Слова
ли
а
більшої
в
всім
для

Семантичне ядро:

Слова	Кількість зустрічей
всім	1
очочим	1
працює	1
персональний	1
супільства	1

Вихід

Визначення рівня запозичень

Підсистема для визначення рівня запозичень

Текст для порівняння:

Python (найчастіше вживане прочитання — «пайтон», запозичено назву[7] з британського шоу Монті Пайтон) — інтерпретована об'єктно-орієнтована мова програмування високого рівня зі строгою

Показати рівень запозичень

Загальний висновок: Максимальний збіг семантичного ядра становить 99 % з текстом: Програмування Python2

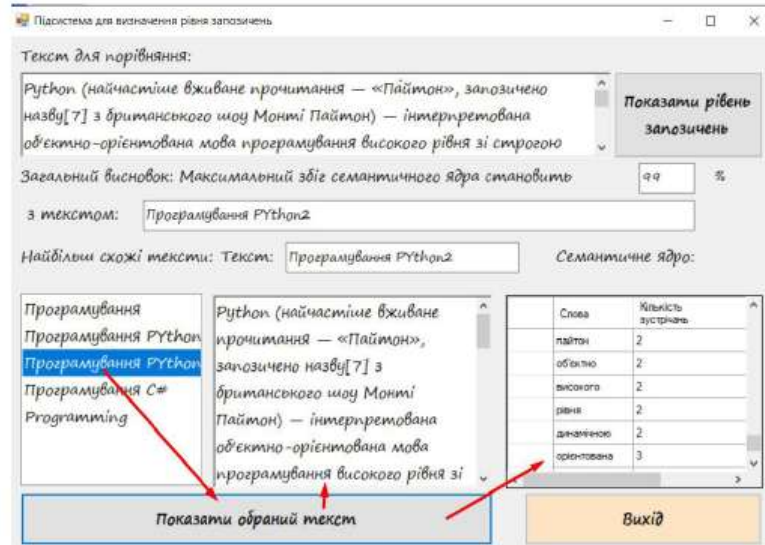
Найбільш схожі тексти: Текст: Семантичне ядро:

Програмування
Програмування Python
Програмування Python
Програмування C#
Programming

Показати обраний текст

Вихід

Відображення тексту, який знайдено як плагіат



Висновки

Відповідно до мети КРБ, було розроблено **метод автоматизованого визначення семантичної подібності цифрових текстів для виявлення запозичень**, а також його **програмну реалізацію** на платформі .NET, що виконує наступні основні функції:

- створення послідовного вектору слів для кожного цифрового тексту авторського базової вибірки;
- створення послідовного вектору слів для тестового тексту;
- створення частотного вектору семантичного ядра для кожного авторського цифрового тексту;
- створення частотного вектору семантичного ядра тестового цифрового тексту;
- очищення частотних векторів семантичних ядер текстів від стоп-слів;
- обчислення оцінок рівня запозичень тестового цифрового тексту щодо кожного цифрового тексту базової вибірки;
- сортування результуючого переліку за оцінкою рівня запозичень та формування висновку щодо визначення семантичної подібності цифрових текстів.

Ім'я користувача:
Кафедра КН

ID перевірки:
1011548894

Дата перевірки:
12.06.2022 09:07:50 EEST

Тип перевірки:
Doc vs Internet + Library

Дата звіту:
12.06.2022 09:14:53 EEST

ID користувача:
100005671

Назва документа: Раєва_ЗАПИСКА_short

Кількість сторінок: 73 Кількість слів: 9759 Кількість символів: 74759 Розмір файлу: 3.11 MB ID файлу: 1011420928

Виявлено модифікації тексту (можуть впливати на відсоток схожості)

6.5% Схожість

Найбільша схожість: 2.4% з джерелом з Бібліотеки (ID файлу: 1011254702)

2.92% Джерела з Інтернету

116

Сторінка 75

4.12% Джерела з Бібліотеки

73

Сторінка 76

0% Цитат

Вилучення цитат вимкнене

Вилучення списку бібліографічних посилань вимкнене

0% Вилучень

Немає вилучених джерел

Модифікації

Виявлено модифікації тексту. Детальна інформація доступна в онлайн-звіті.

Підозріле форматування

22
сторінки

Anti-Plagiarism v-15.257

Максимальное совпадение с одним документом 5.0%

Словари проверки: en_US, ru_RU, ua_UA. **Ошибок в документах: 11%**

ID: 105045 Название: КВАЛІФІКАЦІЙНА РОБОТА БАКАЛАВРА на тему Метод автоматизованого визначення семантичної подібності цифрових текстів для виявлення запозичень Добавлено в БД: 2022-06-12 Авторы: А.В. Раєва Руководители: О.В. Мазурець Консультанты: Оponentы:	Документ		Суммарное совпадение по Базе Данных	
	Символы	Лексемы	Символы	Лексемы
	59678	885	4957 (8%)	77 (9%)

Источник плагиата

ID	Описание	Наличие плагиата в документе	
		Символы	Лексемы

**РІШЕННЯ ЕКСПЕРНОЇ КОМІСІЇ КАФЕДРИ КОМП'ЮТЕРНИХ НАУК
ПРО ДОПУСК КВАЛІФІКАЦІЙНОЇ РОБОТИ ДО ЗАХИСТУ**

Підтверджуємо ознайомлення з результатом звіту подібності щодо роботи, генерованого системою виявлення текстових збігів/ідентичності/схожості:

Назва: Метод автоматизованого визначення семантичної подібності цифрових текстів для виявлення запозичень

Автор: студентка групи КН-18-1 Раєва Анна Вячеславівна

Спеціальність: 122 – Комп'ютерні науки

Освітня програма: освітньо-професійна

Науковий керівник: к.т.н., доцент кафедри КН Мазурець О.В.

Після аналізу звіту подібності зроблено такий висновок:

№	Висновок	Позначка про відповідність
1	Запозичення, виявлені в роботі, є законними і не є плагіатом. Робота приймається до захисту.	<i>відповідає</i>
2	Виявлені запозичення не є плагіатом, розміщені в розділах, які не описують безпосередньо авторське дослідження, але кількість цитат перевищує обсяг, виправданий поставленою метою роботи. Робота приймається до захисту, але має бути відкоригована. Відкоригований варіант має бути поданий на кафедру за 2 дні до захисту, разом із заявою щодо самостійності виконання письмової роботи та ідентичності друкованої та електронної версії роботи	
3	Виявлені запозичення не є плагіатом, але частково розміщені в розділах, які описують безпосередньо авторське дослідження, а кількість цитат перевищує обсяг, виправданий поставленою метою роботи. В зв'язку з цим мета роботи та поставлені завдання не були досягнені. Робота може бути допущена до захисту (наступного року) після того як буде відкоригована та допрацьована і успішно пройде повторну перевірку на академічний плагіат.	
4	Робота містить навмисні текстові спотворення, передбачувані спроби укриття запозичень або інші прояви академічного плагіату. Робота містить фабрикацію або фальсифікацію даних. Робота не допускається до захисту.	

Підтвердження:

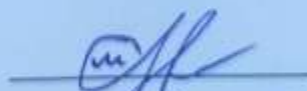
Запозичення, виявлені в роботі Раєвої А.В., не є плагіатом, оскільки: запозичення розміщені в розділі огляду існуючих підходів, не описують безпосередньо авторську роботу і не стосуються її результатів; усі запозичення фрагментарні; до запозичень входять фрагменти програмного коду, що не мають авторства і містять поширені конструкції; серед запозичень знаходяться загальновідомі терміни та скорочення.

Обсяг запозичень, визначений системами виявлення збігів/ідентичності/схожості, складає:

- за системою Anti-Plagiarism: 8%;

- за системою Unicheck: 6,5%.

Керівник роботи



Олександр МАЗУРЕЦЬ

Гарант ОП



Олександр МАЗУРЕЦЬ

Завідувач кафедри КН



Олександр БАРМАК



ВІДГУК НАУКОВОГО КЕРІВНИКА на кваліфікаційну роботу бакалавра

студентки гр. КН-18-1 Раєвої Анни Вячеславівни

за темою Метод автоматизованого визначення семантичної подібності цифрових текстів для виявлення запозичень

1. Актуальність теми

Задача створення систем антиплагіату та довершення методів пошуку запозичень в цифрових текстах є надзвичайно актуальною, адже кожного дня створюються терабайти інформації, яку необхідно досліджувати на предмет використання їх неналежним чином. Для валідації реалізованого методу необхідно створити відповідний програмний застосунок на його основі, який надав би можливість пошуку запозичень у цифрових текстах, виділенні місць запозичень та виведенні загальної статистики за результатами дослідження.

2. Відповідність роботи предметній області Стандарту спеціальності 122 Комп'ютерні науки

За стандартом, а саме описом предметної області, об'єктами вивчення та діяльності є математичні, інформаційні, імітаційні моделі реальних явищ, об'єктів, систем і процесів та методи і технології отримання, зберігання, обробки, передачі та використання інформації. Метою роботи саме є створення методу автоматизованого визначення семантичної подібності цифрових текстів для виявлення запозичень та реалізація програмного застосунку на базі цього методу. При вирішенні поставленої задачі використано математичні моделі, методи та алгоритми розв'язання теоретичних і прикладних задач, що виникають при розробці інформаційних технологій. Тому результати виконання кваліфікаційної роботи бакалавра відповідають стандарту бакалавра спеціальності 122 – Комп'ютерні науки.

3. Професійні та особистісні якості бакалавра

Під час написання кваліфікаційної роботи бакалавра Раєва Анна Вячеславівна проявила себе відповідальним, кваліфікованим фахівцем та дисциплінованим студентом, згідно з встановленим графіком виконувала роботу, досягала поставлені цілі. В процесі створення пояснювальної записки та розробки програмного продукту проявила достатні

для одержання результату знання та навички, сформовані під час навчання. Опанувала професійні навички за напрямком «Комп'ютерні науки».

4. Ступінь самостійності під час виконання кваліфікаційної роботи

Одержані в роботі результати є наслідком особистої діяльності студентки, котра самостійно виконувала всі поставлені, відповідно до плану, задачі.

5. Ступінь оволодіння методами дослідження

При реалізації кваліфікаційної роботи проявила достатній рівень компетентностей та володіння необхідними інструментами та обладнанням, методами, методиками та технологіями предметної області комп'ютерних наук.

6. Повнота та якість розкриття теми роботи

Тема роботи в повній мірі обґрунтована й розкрита, проведено аналіз актуальності та відомих досліджень в межах обраної теми, поставлені завдання, які у роботі виконані, та розроблено програмне забезпечення для валідації та верифікації запропонованого метода.

7. Логічність, послідовність, аргументованість, літературна грамотність викладення матеріалу

Структура роботи та послідовність викладення логічні та відповідають поставленій меті. Викладення матеріалу послідовне, аргументоване, літературно грамотне.

8. Можливість практичного застосування кваліфікаційної роботи бакалавра, окремих її частин

Розроблений у роботі метод автоматизованого визначення семантичної подібності цифрових текстів може бути використаний в установах, що стикаються із задачами аналізу тексту на предмет збігів з іншими існуючими цифровими текстами, як-от вищі навчальні заклади, тощо.

9. Висновок про можливість допуску кваліфікаційної роботи бакалавра до захисту, на яку оцінку заслуговує робота

Враховуючи високий рівень виконання та забезпечення усіх необхідних вимог, робота може бути допущена до захисту. Рекомендована оцінка «добре».

Керівник



к.т.н., доцент каф. КН Олександр МАЗУРЕЦЬ

РЕЦЕНЗІЯ

на кваліфікаційну роботу бакалавра

студентки гр. КН-18-1 Расвої Анни Вячеславівни

за темою Метод автоматизованого визначення семантичної подібності цифрових текстів для виявлення запозичень

1. Актуальність і значення Задача створення систем антиплагіату та ефективних методів пошуку запозичень в цифрових текстах є надзвичайно актуальною, оскільки кожного дня створюються терабайти інформації, яку необхідно досліджувати на предмет використання їх неналежним чином. Для валідації реалізованого методу необхідно створити відповідний програмний застосунок на його основі, який надав би можливість пошуку запозичень у цифрових текстах, виділенні місць запозичень та складенні загальної статистики за результатами дослідження.

2. Оцінка запропонованих моделей, підходів, алгоритмів, інформаційної складової та засобів розробки Реалізовані авторкою схема і компоненти методу автоматизованого визначення семантичної подібності цифрових текстів для виявлення запозичень дозволяють ефективно знаходити збіги в цифрових текстах, вказувати на місця, де було використано частину існуючого тексту та оцінити унікальність тексту.

3. Оцінка розробленої інформаційної системи, її практична цінність та економічна доцільність Створений метод дозволяє ефективно вирішувати актуальну задачу інформаційних технологій, а саме визначення запозичень та плагіату в цифрових текстах.

4. Загальний висновок та оцінка Всі поставлені завдання виконані повністю та змістовно. За своєю структурою, практичними цінностями, поставленій меті та вирішеними завданнями робота відповідає вимогам вищого навчального закладу і вимогам, що пред'являються до освітньо-кваліфікаційного рівня «бакалавр», а її автор Расва А.В. заслуговує присвоєння кваліфікації бакалавра з комп'ютерних наук.

Робота заслуговує на оцінку « добре ».

Рецензент Берратюк Л.П. д.ф.м.н.