

NLP: опрацювання омонімів в українськомовних текстах

Доповідач: Парнак Д. А.
Керівник: ст. в. Смиш О. Р.

Обґрунтування вибору теми дослідження



Брак готових
рішень

Неоднозначність
омонімів



Неправильне
опрацювання на
морфологічному
рівні



Помилки в
результатах
роботи систем із
використанням
NLP

Мета роботи



уможливити підвищення точності розомонімізації в українській мові, шляхом створення системи для розпізнавання та розбору омонімів в текстах, написаних українською мовою.

Об'єкт дослідження

Оброблення текстових даних, написаних природною українською мовою.

Предмет дослідження

Розроблення системи для опрацювання омонімів, а також тренування моделі розпізнавання омонімів з метою використання в проєктах з обробки текстових даних.

Завдання роботи

- **аналіз** поточних рішень для обробки природної української мови, зокрема розпізнавання омонімів в текстах;
- **формування** методів та тренування моделі машинного навчання для розпізнавання та аналізу різних груп омонімів та паронімів;
- **апробація** результатів застосування реалізованих методів.

Наукова новизна

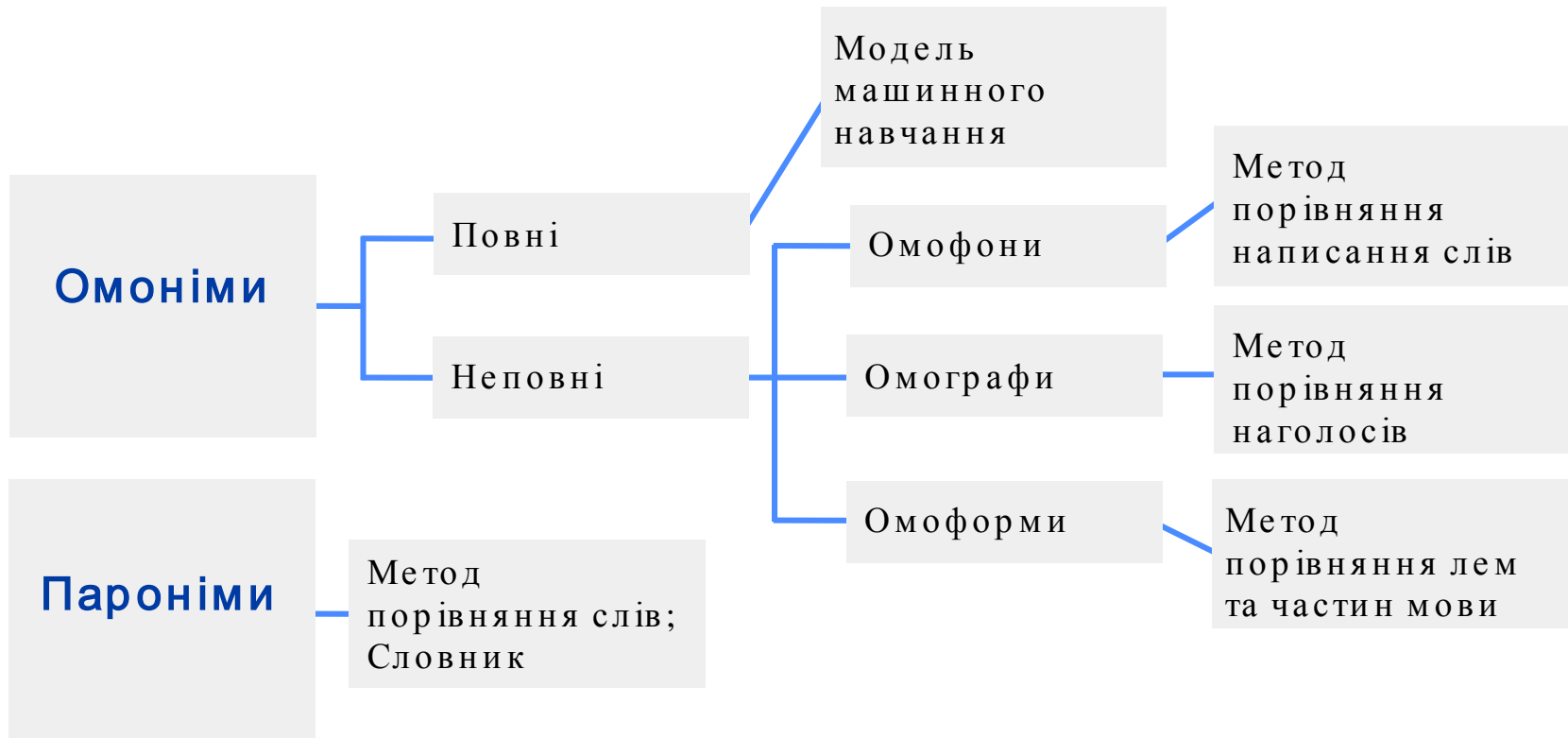


уперше проаналізовано наявні рішення для автоматизованої роботи з українськомовними омонімами в тексті;



уперше розроблено комплексну систему для автоматизованого виявлення та розбору всіх груп омонімів із використанням методів на основі правил та моделі машинного навчання.

Аналіз омонімів та паронімів



Аналіз словника паронімів

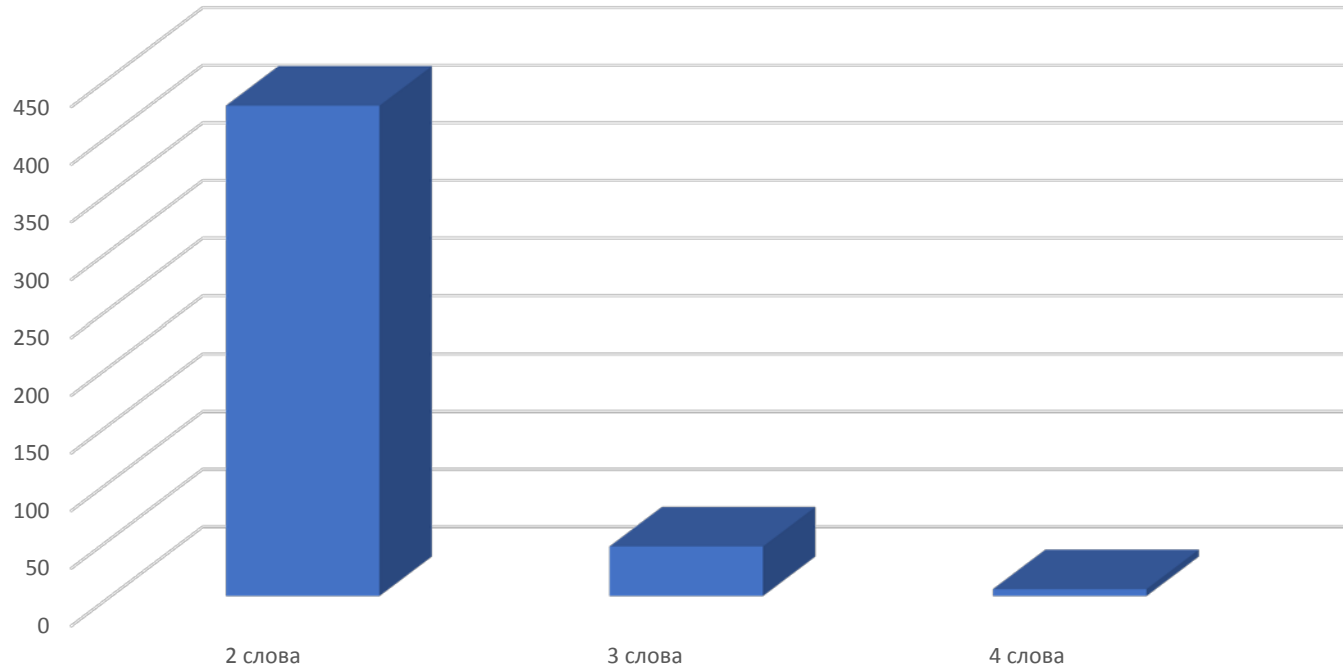


paronyms_json

```
1 {  
2   "group": "КОНТАКТ // КОНТРАКТ",  
3   "paronym1": "контакт",  
4   "def1": "Діловий зв'язок, тісні стосунки, взаємна узгодженість дій.",  
5   "paronym2": "контракт",  
6   "def2": "Письмова угода, договір із взаємними зобов'язаннями його  
7 учасників."  
8 }
```

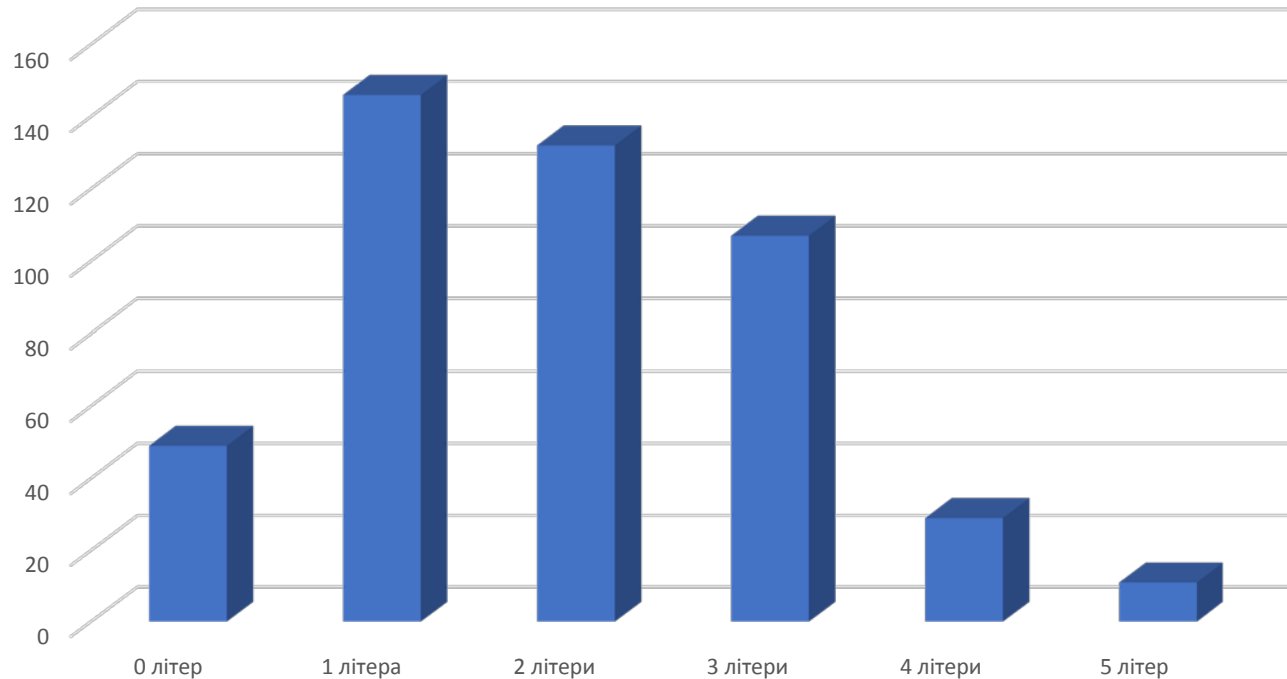
Аналіз словника паронімів

Кількість слів у паронімічних групах



Аналіз словника паронімів

Різниця довжини слів у групах (у літерах)



Розроблені методи

Опрацювання омоформ



- Очищення тексту від пунктуації
- Встановлення лем, та частин мови лексем та їхнє паралельне порівнювання

Опрацювання омографів



- Встановлення подвійних наголосів
- Підтвердження істинності омографів за врахування розбіжності тлумачень

Опрацювання омофонів



- Очищення тексту від пунктуації
- Пошук омофонів за регулярними виразами

Опрацювання паронімів



- Пошук паронімічних груп за регулярними виразами
- Виокремлення окремих паронімів та тлумачень зі словника

Приклад роботи вебінтерфейсу програми

Результати аналізу

Оригінальний текст:

Скоро розпочинається вступна кампанія. Ці брати були дружніми. Цей замок обабіч місця, де ми часто збиралися разом; там завжди яскраво світило сонце. Не можна брати цей атлас, він не мій. Якби ж я міг мати щось більше, ніж просто те, що дала мені дорога мати. Сон - це найкращий відпочинок. Я би не хотів працювати у цій компанії.

Аналіз:

Омоформи:

брати - леми: ['брат', 'брати'], частини мови: ['NOUN', 'VERB']

мати - леми: ['мати', 'мати'], частини мови: ['VERB', 'NOUN']

Слова з різним наголосом:

замок

1: ЗАМОК, мка, ч. Укріплене житло феодала доби середньовіччя з оборонними, господарськими, культовими і т. ін. будівлями, звичайно оточене високим кам'яним муром з кількома вежами.

2: ЗАМОК, мка, ч. Пристрій для замикання дверей у приміщеннях, дверцят шафи, скринь, шухляд і т. ін.

атлас

1: АТЛАС, у, ч. Укладений за певною системою і виданий у формі альбома або книжки збірник географічних, історичних та ін. карт.

2: АТЛАС, у, ч. Шовкова або напівшовкова тканина, блисуча і гладенька з лиця.

Омофони:

Сон - це - сонце.

Пароніми (зі словника):

кампанія паронім до компанія

Пароніми (в тексті):

кампанія - компаніі

Тренування моделі

- Омонім: "ручка"
- 504 речення в датасеті

Формування
датасету: GPT4



Препроцесинг
UDPipe



Векторизація TF-
IDF



Тренування
моделі логістичної
регресії

Результати тренування моделі

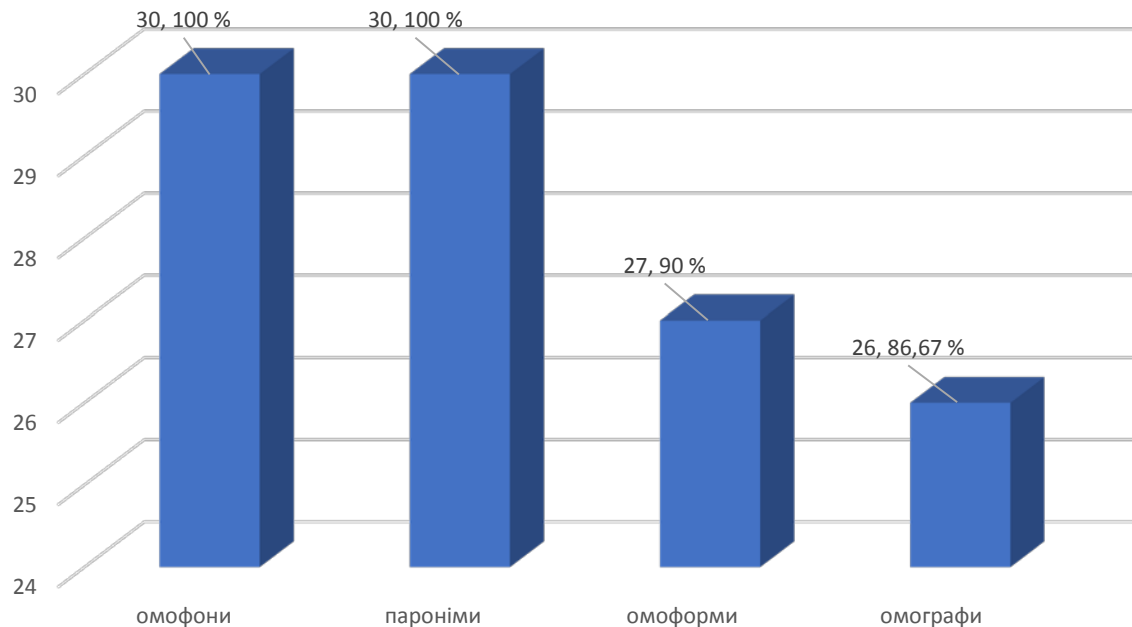


classification_report

	precision	recall	f1-score	support	
1					
2					
3					
3	Зменшувальна до рука	0.87	0.93	0.90	50
4	Прилад для писання	0.93	0.93	0.93	50
5	Частина предмета	0.92	0.86	0.89	50
6					
7	accuracy		0.90	150	
8	macro avg	0.91	0.90	0.90	150
9	weighted avg	0.91	0.90	0.90	150
10					
11	Accuracy:	0.9047619047619048			

Апробація результатів: аналіз точності методів на основі правил

Правильно опрацьовані лексеми у реченнях



Апробація результатів: аналіз точності моделі машинного навчання

Омонім: "вид"

```
classification_report
  precision  recall  f1-score  support
1
2
3      Вигляд, краєвид      0.92      0.88      0.90      13
4 Окрема галузь, категорія      0.93      0.90      0.91      13
5
6      accuracy      0.90      26
7      macro avg      0.92      0.89      0.90      26
8      weighted avg      0.92      0.89      0.90      26
9
10 Accuracy: 0.90354630104
```

Омонім: "нота"

```
classification_report
  precision  recall  f1-score  support
1
2
3      Музичний знак      0.93      0.91      0.92      13
4 Офіційне звернення      0.92      0.91      0.91      13
5
6      accuracy      0.92      26
7      macro avg      0.93      0.91      0.92      26
8      weighted avg      0.93      0.91      0.92      26
9
10 Accuracy: 0.923890416428
```

Омонім: "дисципліна"

```
classification_report
  precision  recall  f1-score  support
1
2
3      Підтримання порядку      0.91      0.90      0.90      13
4 Предмет, розділ науки      0.88      0.89      0.88      13
5
6      accuracy      0.89      26
7      macro avg      0.89      0.90      0.89      26
8      weighted avg      0.89      0.90      0.89      26
9
10 Accuracy: 0.8924538022519
```

Висновки

- Здійснено аналіз ресурсів, які застосовуються для автоматичного опрацювання омонімів в обробці природної української мови;
- Проведено аналіз проблематики омонімії та паронімії в контексті українськомовного NLP;
- Розроблено методи на основі правил для розпізнавання часткових омонімів;
- Натреновано модель розпізнавання повного омоніма з точністю 90,47 %, що може використовуватись як попередньо тренована
- Здійснено апробацію результатів роботи, встановлено точність розробленої системи на рівні понад 86 %

**Дякую
за
увагу!**

