

Як ШІ може допомагати боротися з вандалізмом у Вікіпедії та Вікіданих

Микола Трохимович
mykola.trokhymovych@upf.edu

Не кожна зміна покращує Вікіпедію – деякі її спотворюють.

Example of edit (a.k.a. revision) reverting a bad-faith one (revision_id = 1149625753).

```
Lvov emerged as the centre of the historical regions of [[Red Ruthenia]] and [[Galicia (Eastern Europe)|Galicia]] in the [[14th
```

```
Lviv emerged as the centre of the historical regions of [[Red Ruthenia]] and [[Galicia (Eastern Europe)|Galicia]] in the [[14th
```

Хоча існували моделі, що допомагали патрульним (наприклад, ORES), залишалися відкриті проблеми, такі як:

- Продуктивність моделей;
- Покриття мов;
- Неупередженість.



Вікіпедія:Вандалізм





Вандалізм у Вікіпедії — явно шкідницьке додання, видалення чи зміна вмісту, скоєне зумисне з метою скомпрометувати достовірність і авторитетність...

Wikipedia

RESEARCH-ARTICLE



Fair Multilingual Vandalism Detection System for Wikipedia

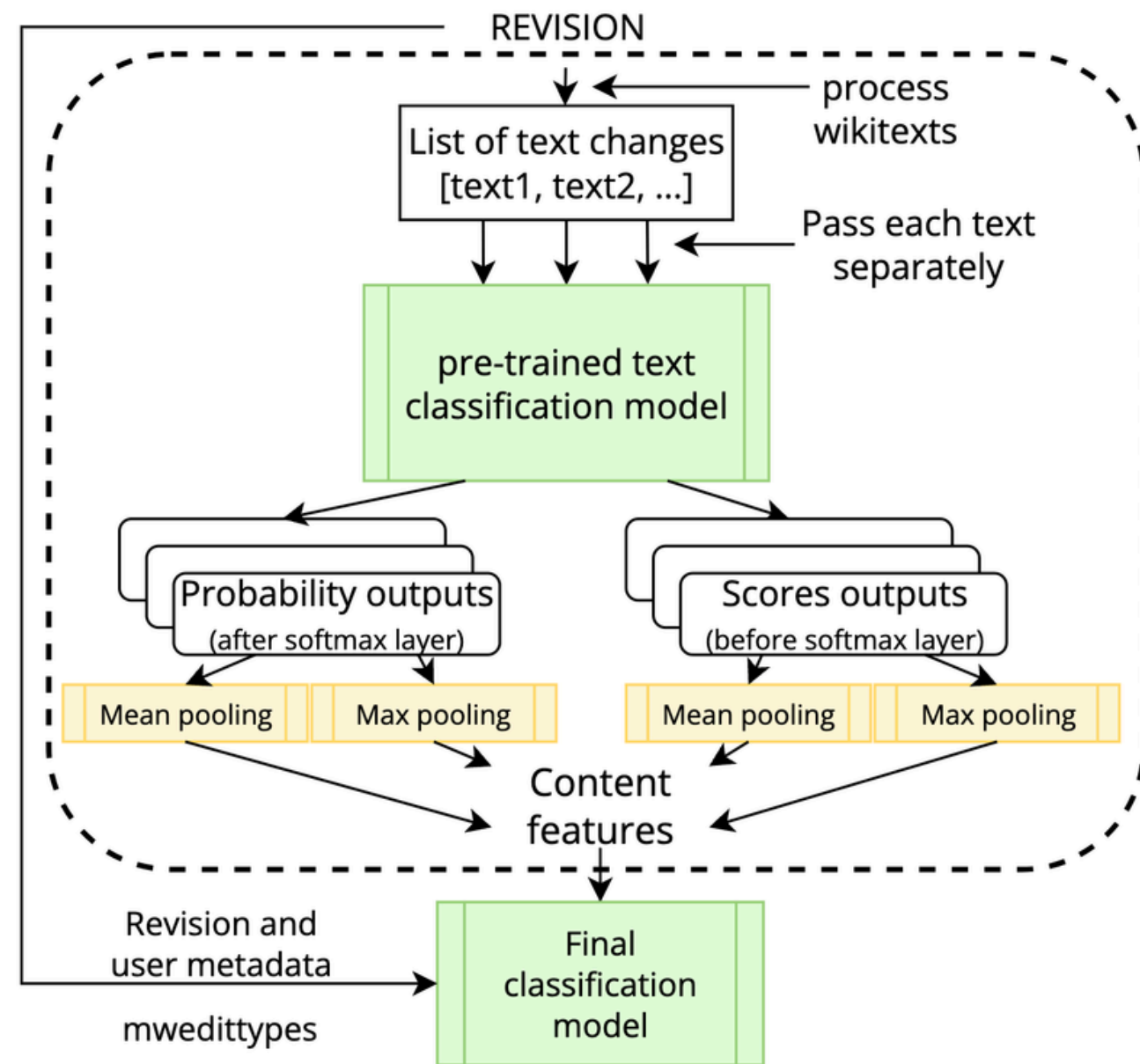
Authors:  [Mykola Trokhymovych](#),  [Muniza Aslam](#),  [Ai-Jou Chou](#),  [Ricardo Baeza-Yates](#),  [Diego Saez-Trumper](#) | [Authors](#)
[Info & Claims](#)

[KDD '23: Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining](#) • Pages 4981 - 4990
<https://doi.org/10.1145/3580305.3599823>

Published: 04 August 2023 [Publication History](#)



Fair Multilingual Vandalism Detection System for Wikipedia



Мета:

Створити модель, яка допомагає редакторам визначати правки, що потребують патрулювання.

Підхід:

- Використовувати історичні дані про скасування змін для навчання моделей машинного навчання.

Результат:

- Відкрита, багатомовна модель для патрулювання контенту на Вікіпедії;
- Перевершує попередні SOTA-моделі за продуктивністю та неупередженістю;
- Впроваджена і працює для 47 мов.

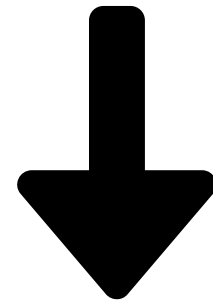
Authors: [Mykola Trokhymovych](#), Muniza Aslam, Ai-Jou Chou, Ricardo Baeza-Yates, Diego Saez-Trumper
KDD'23 Industry Track

Анонімні редактори мають вищий відсоток скасованих змін.

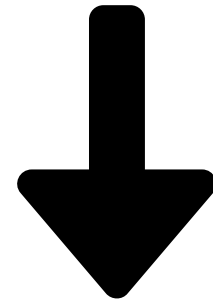


anon.
28%

all
8%



Моделі перенавчаються і починають дискримінувати анонімних користувачів.



Анонімні редактори рідше залишаються та продовжують робити внесок.

Disparate Impact Ratio (DIR)

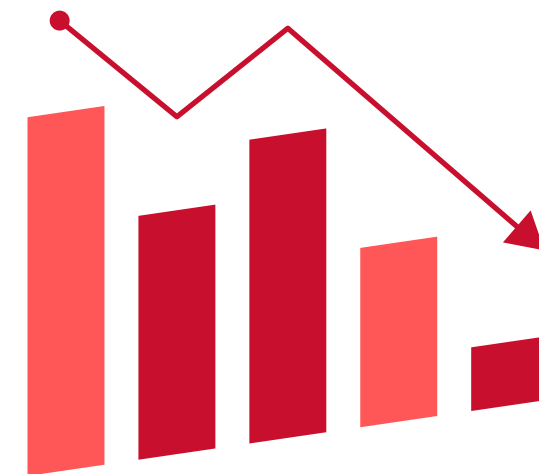
$$DIR = \frac{Pr(\hat{Y}=1|D=unprivileged)}{Pr(\hat{Y}=1|D=privileged)} = \underline{\underline{20}}$$

**for ORES*

Pr - probability

\hat{Y} - predicted value,

D - a group of users (anon. or registered)



Ризик зменшення кількості активних редакторів.

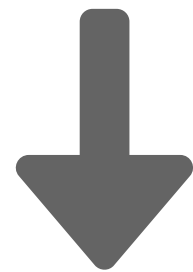
Вандалізм поза Вікіпедією

Malicious changes to Wikidata:

Example of a revision (ID: 593195479) vandalizing the Wikidata entry for Bulgaria. Original triple IDs are mapped to their corresponding English labels.

before Q219 (Bulgaria) P85 (anthem) Q182115 (Mila Rodino)

after Q219 (Bulgaria) P85 (anthem) Q28572509 (Despacito)



Який гімн Болгарії? 🇧🇬

TL ; DR

Why does Siri think the national anthem of Bulgaria is 'Despacito?' / A cursory investigation with no clear answer in sight

by [Nick Statt](#)

Oct 5, 2017, 1:15 AM GMT+2

[theverge](#)



Graph-Linguistic Fusion: Using Language Models for Wikidata Vandalism Detection

Mykola Trokhymovych, Lydia Pintscher, Ricardo Baeza-Yates, Diego Sáez Trumper

Abstract

We introduce a next-generation vandalism detection system for Wikidata, one of the largest open-source structured knowledge bases on the Web. Wikidata is highly complex: its items incorporate an ever-expanding universe of factual triples and multilingual texts. While edits can alter both structured and textual content, our approach converts all edits into a single space using a method we call Graph2Text. This allows for evaluating all content changes for potential vandalism using a single multilingual language model. This unified approach improves coverage and simplifies maintenance. Experiments demonstrate that our solution outperforms the current production system. Additionally, we are releasing the code under an open license along with a large dataset of various human-generated knowledge alterations, enabling further research.

[PDF](#)[Cite](#)[Search](#)[Fix data](#)

Anthology ID: 2025.acl-industry.21

Volume: [Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics \(Volume 6: Industry Track\)](#)

Month: July

Year: 2025

Address: Vienna, Austria

Editors: [Georg Rehm](#), [Yunyao Li](#)

Venue: [ACL](#)

Що ми знаємо про Wikidata?

Label Madrid (Q2807) **Item Identifier (QID)**

Description municipality and capital of Spain

▼ In more languages
Configure

Language	Label	Description
English	Madrid	municipality and capital of Spain
Russian	Мадрид	столица и крупнейший город Испании
Spanish	Madrid	capital y municipio más poblado de España
Catalan	Madrid	capital d'Espanya

Statements

Property	Value	Qualifier
country	Spain	
	start time	23 January 1516
	▶ 1 reference	Collapsed references

- Структуроване сховище знань
 - ~100 млн+ триплетів у Wikidata
 - ~8 правок на секунду
- Різноманіття типів даних:
 - Тексти більш ніж 300 мовами
 - Мітки часу
 - Числові значення
 - Географічні координати
 - І багато іншого...
- Живить численні сервіси в Інтернеті

Парсинг ознак

DeepDiff v 8.2.0

downloads 21M/month python 3.8 | 3.9 | 3.10 | 3.11 | 3.12 | 3.13 license MIT Unit Tests passing codecov 97%

Modules

- [DeepDiff](#): Deep Difference of dictionaries, iterables, strings, and ANY other object.
- [DeepSearch](#): Search for objects within other objects.
- [DeepHash](#): Hash any object based on their content.
- [Delta](#): Store the difference of objects and apply them to other objects.
- [Extract](#): Extract an item from a nested Python object using its path.
- [commandline](#): Use DeepDiff from commandline.

🕒 1,220 Commits

last week


last week

last week

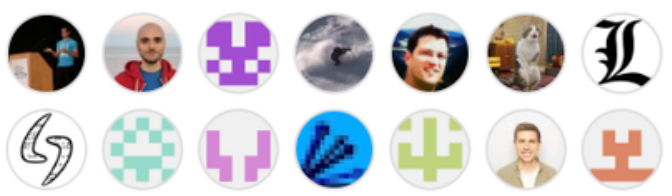
The MIT License (MIT)

Copyright (c) 2014 – 2021
www.zepworks.com

Used by 15.3k

 + 15,257

Contributors 72



+ 58 contributors

DeepDiff

Inserting description

description / bg	description / bg	['descriptions']['bg']: {'value': 'столицата на Испания'}
	+ столицата на Испания	

Removing statement

Property / instance of		['claims']['P31']: [{'datavalue': {'value': {'entity-type': 'item', 'id': 'Q515'}}}]
- city		

Changing statement

Property / coordinate location	Property / coordinate location	['claims']['P625']['datavalue']['latitude']: {'new_value': 40.251, 'old_value': 40.250} ['claims']['P625']['datavalue']['longitude']: {'new_value': -3.4212, 'old_value': -3.429}
40° 25' 0", -3° 42' 9" Latitude 40.4166666666667 Longitude -3.7025 Precision 0.00027777777777778 Globe Earth	40° 25' 1", -3° 42' 12" Latitude 40.4169444444444 Longitude -3.703333333333333 Precision 0.00027777777777778 Globe Earth	

Деталі обробки тексту

Три різні сигнали – одна мала мовна модель (SLM).

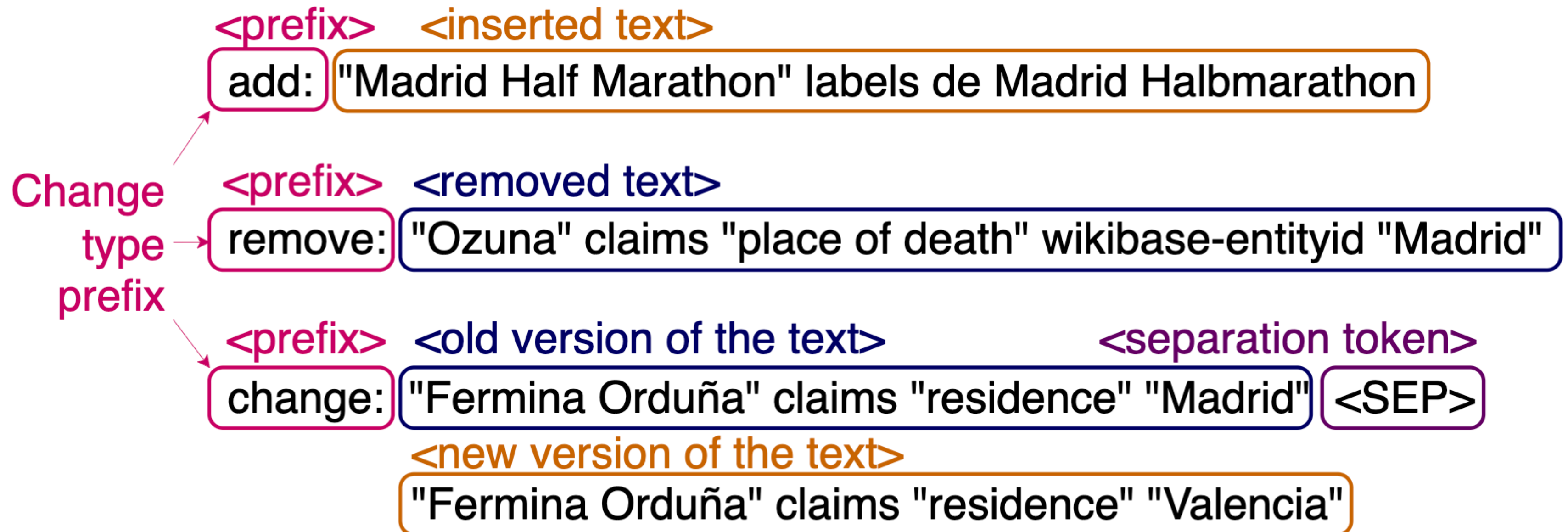
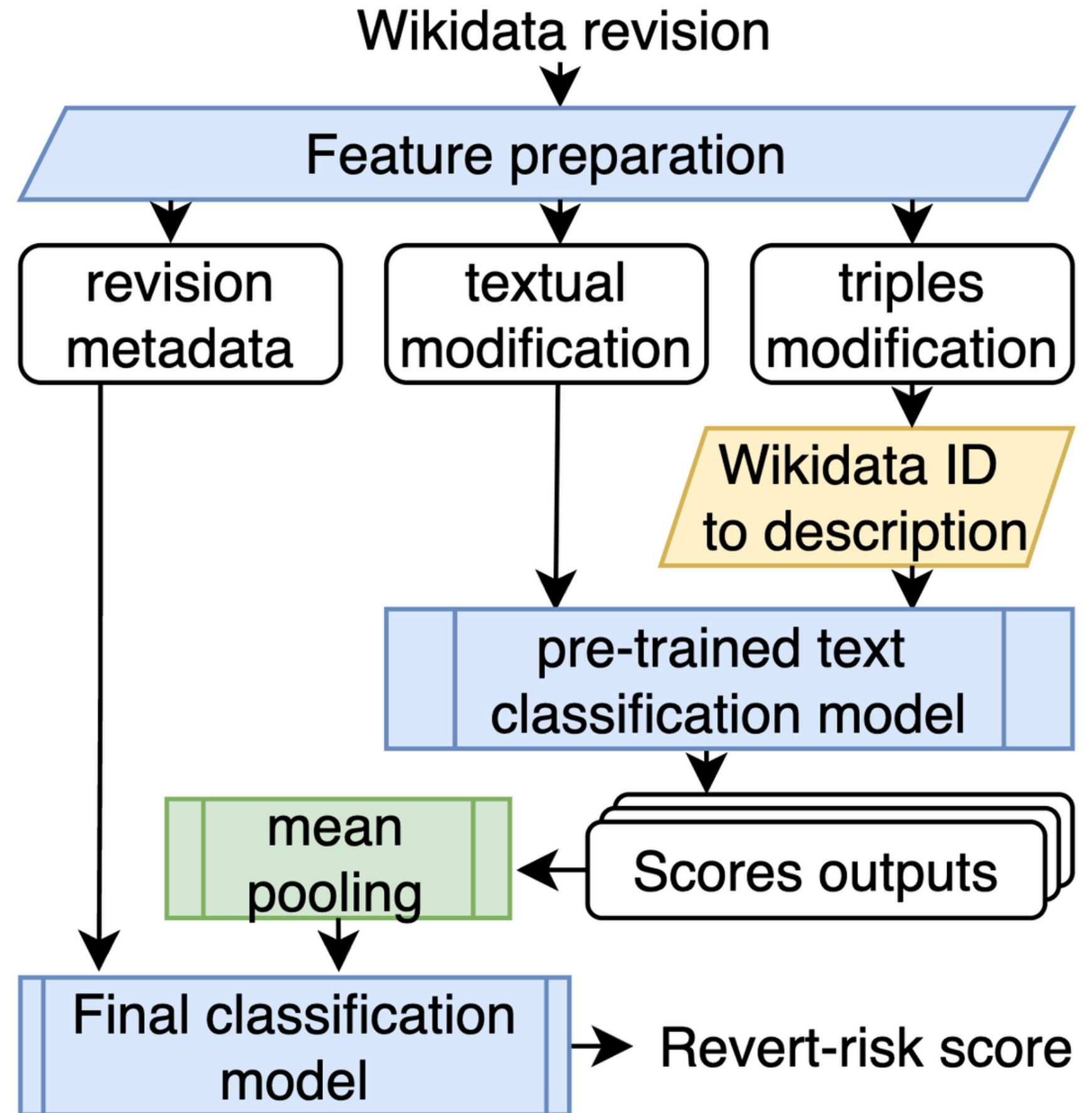


Схема системы



Оцінка якості системи

Table 3: System performance on holdout testing set.

Model	AUC	CI	FR@0.99	FR@0.9	FR@0.7
Rule-based	0.760	[0.74, 0.78]	0.0	0.0	0.92
ORES	0.859	[0.84, 0.87]	0.45	0.88	0.94
MbC	0.880	[0.87, 0.89]	0.55	0.89	0.94
CbC	0.876	[0.86, 0.89]	0.60	0.82	0.93
Graph2Text	0.924	[0.91, 0.93]	0.71	0.91	0.96

- **Rule-based:** all revisions by anonymous users are reverted;
- **ORES:** previous SOTA model;
- **MbC:** Metadata-based Classifier (MbC)
- **CbC:** Content-based Classifier (CbC)
- **Graph2Text:** proposed model (CbC + MbC)

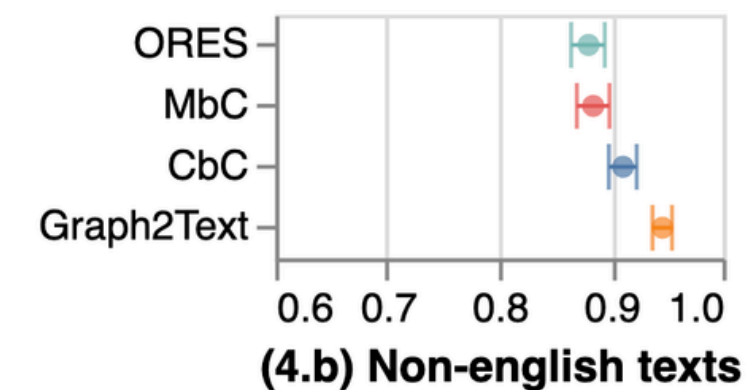
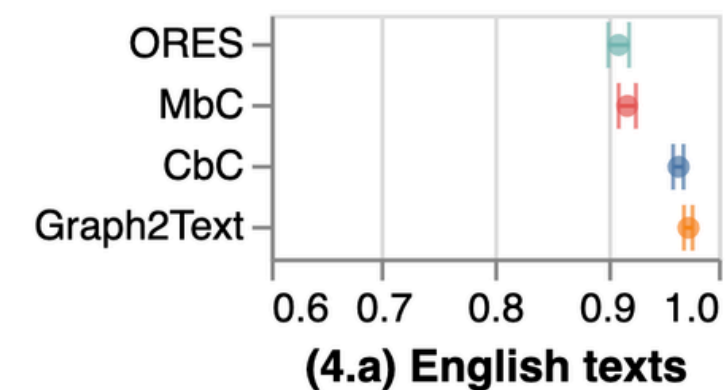
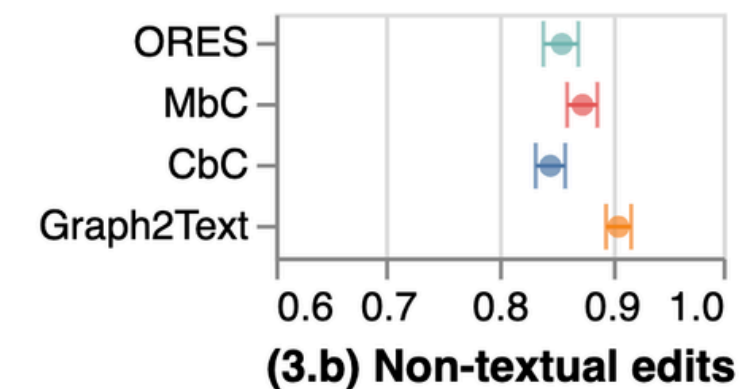
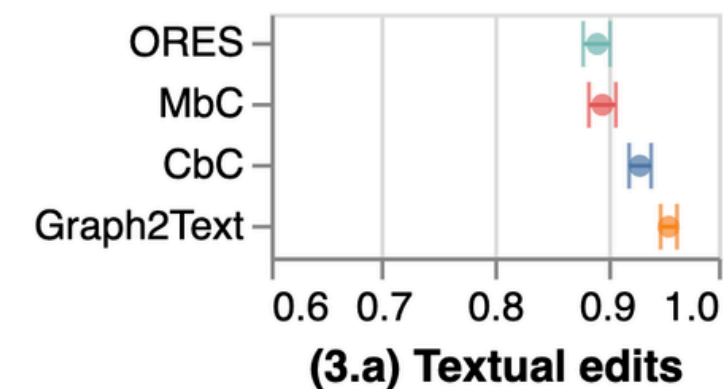
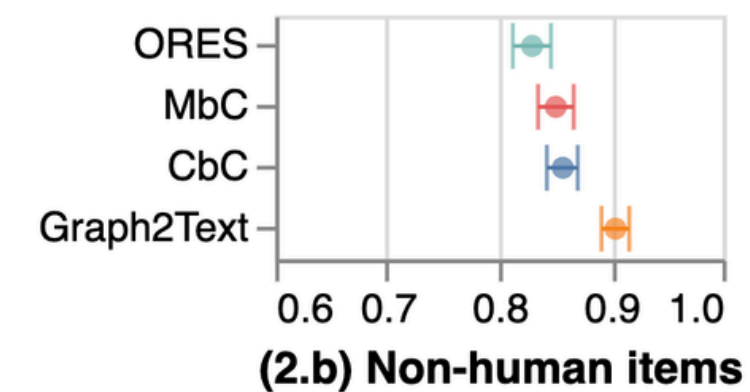
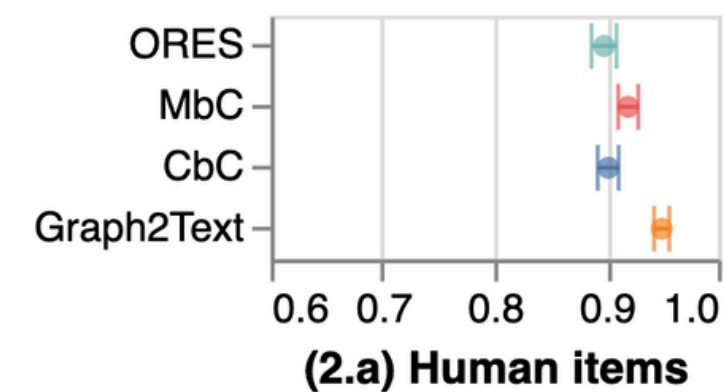
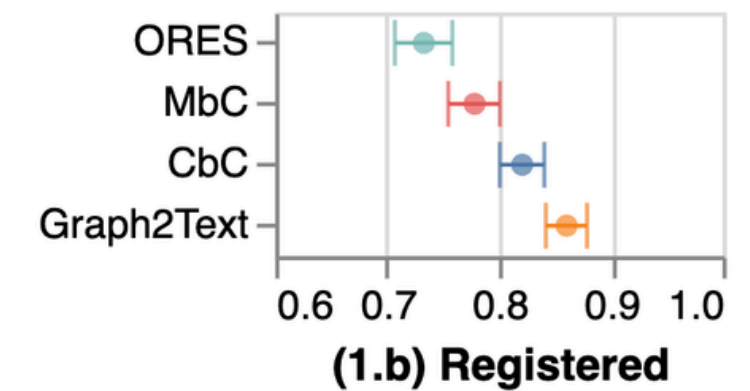
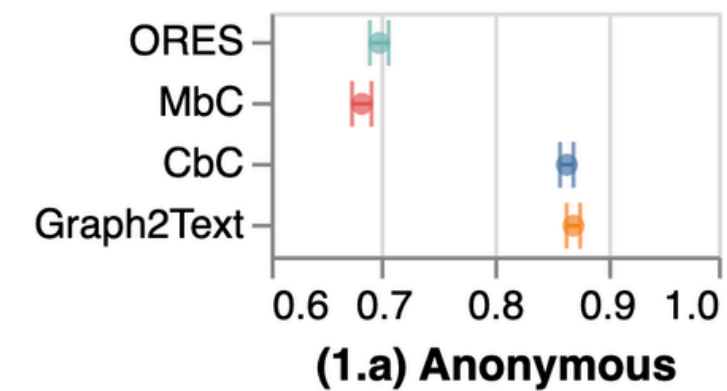
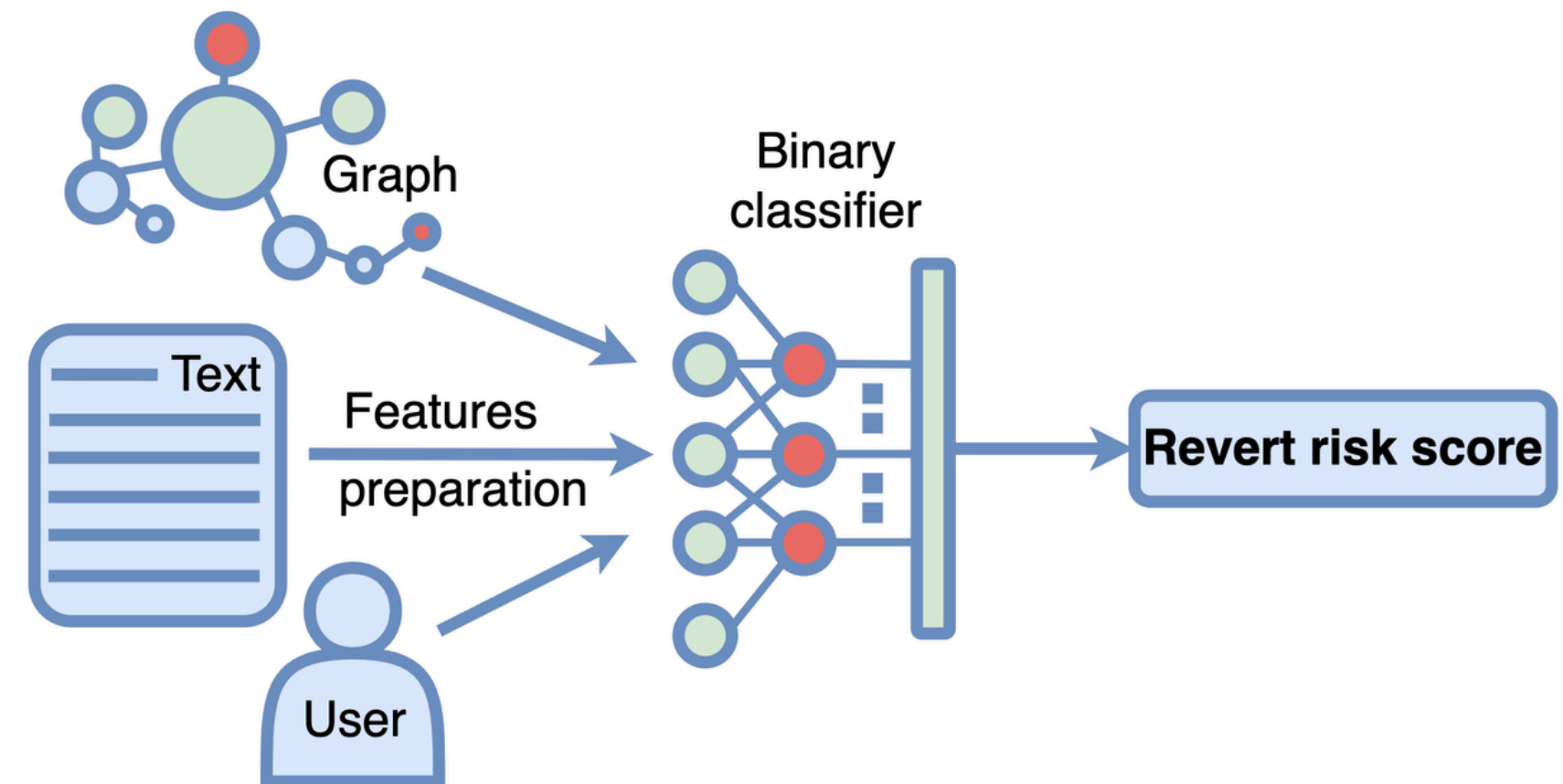


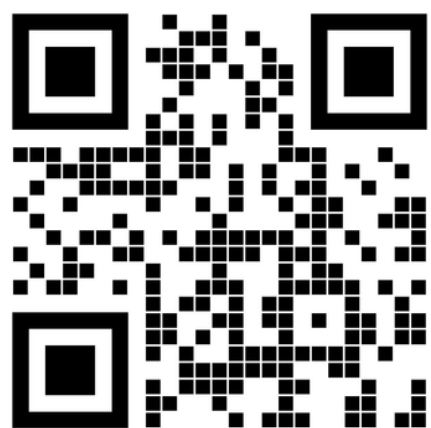
Figure: Models performance (AUC)

Набуті знання

- Ми можемо оцінювати якість змін контенту, не аналізуючи сам контент.
- **Сильні базові моделі** потребують швидкого процесингу ознак (наприклад, метадані користувача).
- **Сильні ознаки контенту допомагають зменшити упередженість.**
- Впровадження становить **80% роботи.**



Дякую!



Mykola Trokhymovych

trokhymovych.com

mykola.trokhymovych@upf.edu