

Advancing Open Knowledge

NLP for Accessibility, Fairness, and Knowledge Integrity in Collaborative Platforms

Mykola Trokhymovych

mykola.trokhymovych@upf.edu

trokhymovych.com

About me

Mykola Trokhymovych

trokhymovych.com | mykola.trokhymovych@upf.edu



- PhD student @ Universitat Pompeu Fabra
 - Co-advised by
 - Diego Saez-Trumper
 - Ricardo Baeza-Yates
- Research collaborator @ Wikimedia Foundation

Previously:

- MSc. Data Science @ Ukrainian Catholic University
- Data scientist / ML engineer @ Ciklum, Jooble, Surprise

Agenda

- 1 Wikipedia ecosystem
- 2 Readability
- 3 Vandalism detection
- 4 Discussion



Why Wikipedia ecosystem?

Largest Human-Curated Knowledge Base

- 300+ languages, 60M+ articles, 100M+ Wikidata triples;
- Maintained entirely by volunteers.

Open, Collaborative, Transparent

- Fully open-source and publicly editable;
- Rich edit history, versioning, and discussion pages.

Real-World Impact

- Search engines, voice assistants, and education;
- Data backbone for LLMs and knowledge graphs.

RICHARD COOKE BUSINESS FEB 17, 2020 6:00 AM

Wikipedia Is the Last Best Place on the Internet

People used to think the crowdsourced encyclopedia represented all that was wrong with the web. Now it's a beacon of so much that's right.

wired

Wikipedia has a solution for the deluge of AI training bots hogging its servers

Non-human traffic has strained the nonprofit's servers, but the organization may have a fix.

By [Cecily Mauran](#) on April 18, 2025

mashable

Wikipedia AI Strategy

“Not too long ago, we were asked when we're going to replace Wikipedia's human-curated knowledge with AI.

The answer? We're not.”




Our new AI strategy puts Wikipedia's humans first

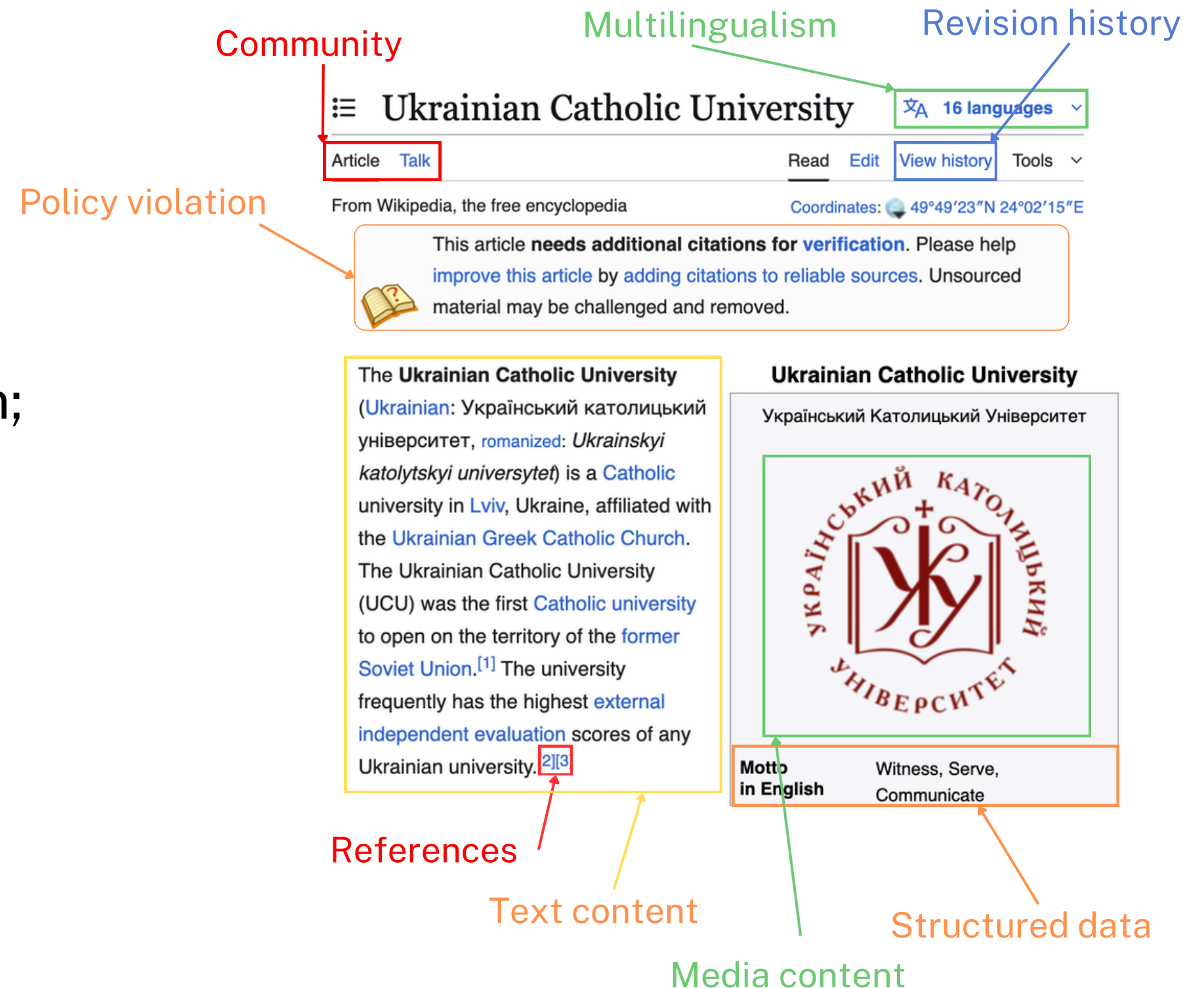
By [Chris Albon](#) and [Leila Zia](#) · 30 April 2025

[Share](#)



NLP Applications in the Wikipedia

-  **Moderation tools:**
 - Talk page tone detection;
 - Vandalism & spam detection.
-  **Editing tools:**
 - Readability evaluation;
 - Grammar/misspelling correction;
 - Detection of policy violations;
 - Source verification.
-  **Reader Tools:**
 - Article summaries;
 - Text to speech;
 - Articles search and recommendation.



The image shows a screenshot of the Wikipedia article for "Ukrainian Catholic University" with several NLP application annotations:

- Community:** A red arrow points to the "Article" and "Talk" tabs.
- Multilingualism:** A green arrow points to the "16 languages" dropdown menu.
- Revision history:** A blue arrow points to the "View history" button.
- Policy violation:** An orange arrow points to a warning box that reads: "This article needs additional citations for verification. Please help improve this article by adding citations to reliable sources. Unsourced material may be challenged and removed."
- Text content:** An orange arrow points to the main body of the article text.
- Media content:** A green arrow points to the university's logo.
- Structured data:** An orange arrow points to the "Motto in English" field, which contains the text "Witness, Serve, Communicate".
- References:** A red arrow points to a citation marker "[2][3]" at the end of a sentence in the article text.

Readability & Wikipedia

What is Readability?

The readability concept aims to capture how easy it is to read a given text. It is usually defined as the sum of all factors that affect a reader's understanding, reading speed, and level of interest.

(Dale and Chall, 1949)

Readability of Wikipedia

“[in English Wikipedia the] overall readability is poor. “

(Lucassen et al., 2012)

“Most of the health information [in English Wikipedia] remains written at a level above the reading ability of average adults.”

(Brezar and Heilman, 2019)

Measuring readability is challenging

Common approaches are
readability formulas
e.g., Flesch Reading Ease
(FRE) Score Formula

$$206.835 - 1.015 \left(\frac{\text{total words}}{\text{total sentences}} \right) - 84.6 \left(\frac{\text{total syllables}}{\text{total words}} \right)$$

Wikipedia exists in more than 300 languages:

- Adapting readability formulas to each language is not scalable/feasible
- Lack of multilingual models/tools
- Scarcity of multilingual ground-truth data
 - Standard datasets are not open (e.g., Newsela)
 - Mostly datasets are for English (no data for most languages)

In this research we aim to ...

build a multilingual model that can automatically access the readability of the text of Wikipedia articles across languages.



Тунéль Кохáння —
Die Bahnstrecke Klewan—
愛情隧道 (烏克蘭語: Тунéль
Т. Кох The Tunnel of Love (Ukrainian:
諾州 Тунéль Кохáння, *Tunel*
公司 *Kokhannya*) is a section of
且不 industrial railway located near
Klevan, Ukraine, that links it with

ukwiki	→	8.83	👩
dewiki	→	8.37	👨
zhwiki	→	12.1	🎓
enwiki	→	4.35	👶
Readability score			

An Open Multilingual System for Scoring Readability of Wikipedia

Data

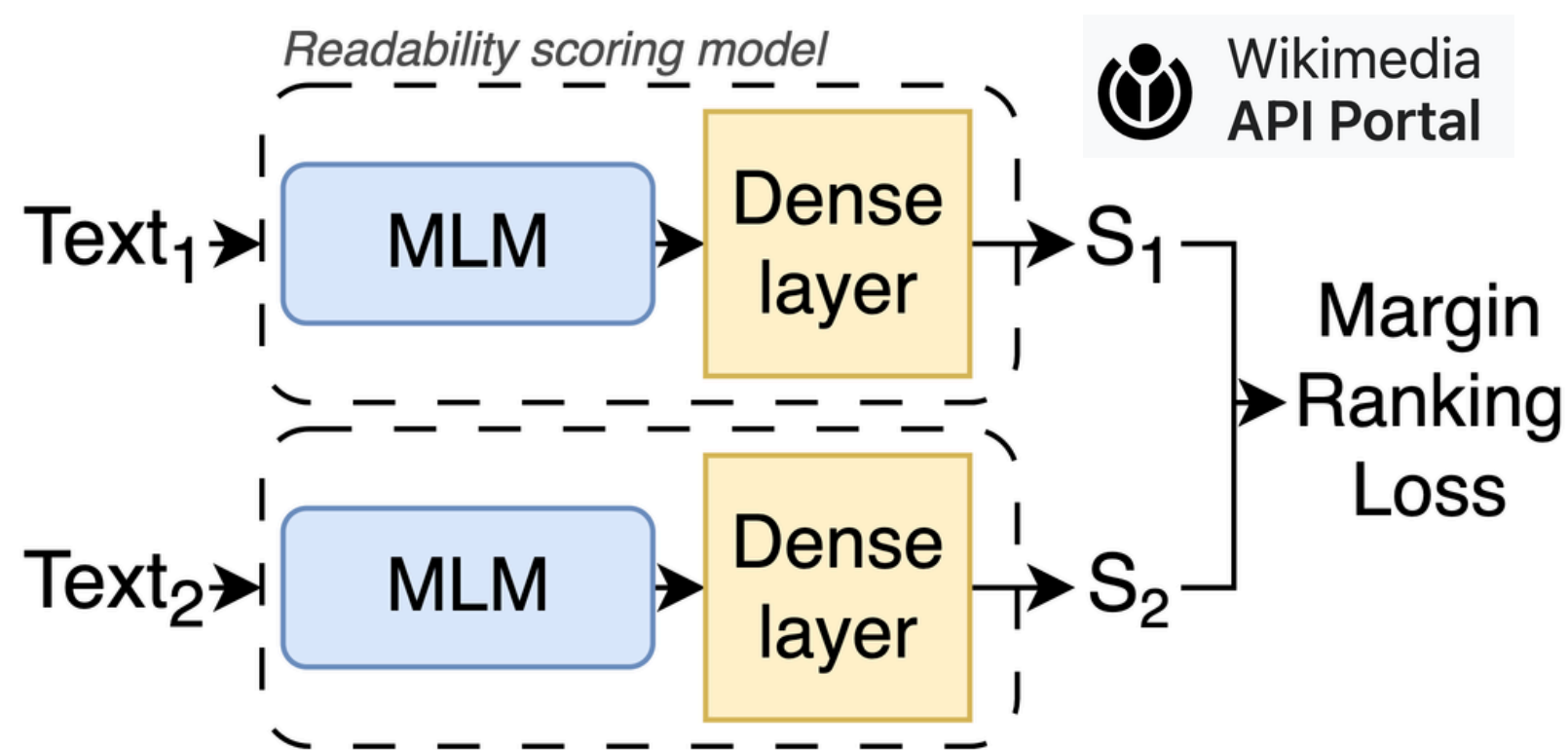
WikiReaD (Wikipedia Readability Dataset)

What? Pairs of encyclopedic articles in two readability levels (simple, hard) for 14 languages. (>100K samples)

How? Matching Wikipedia articles with a simplified version.



Model

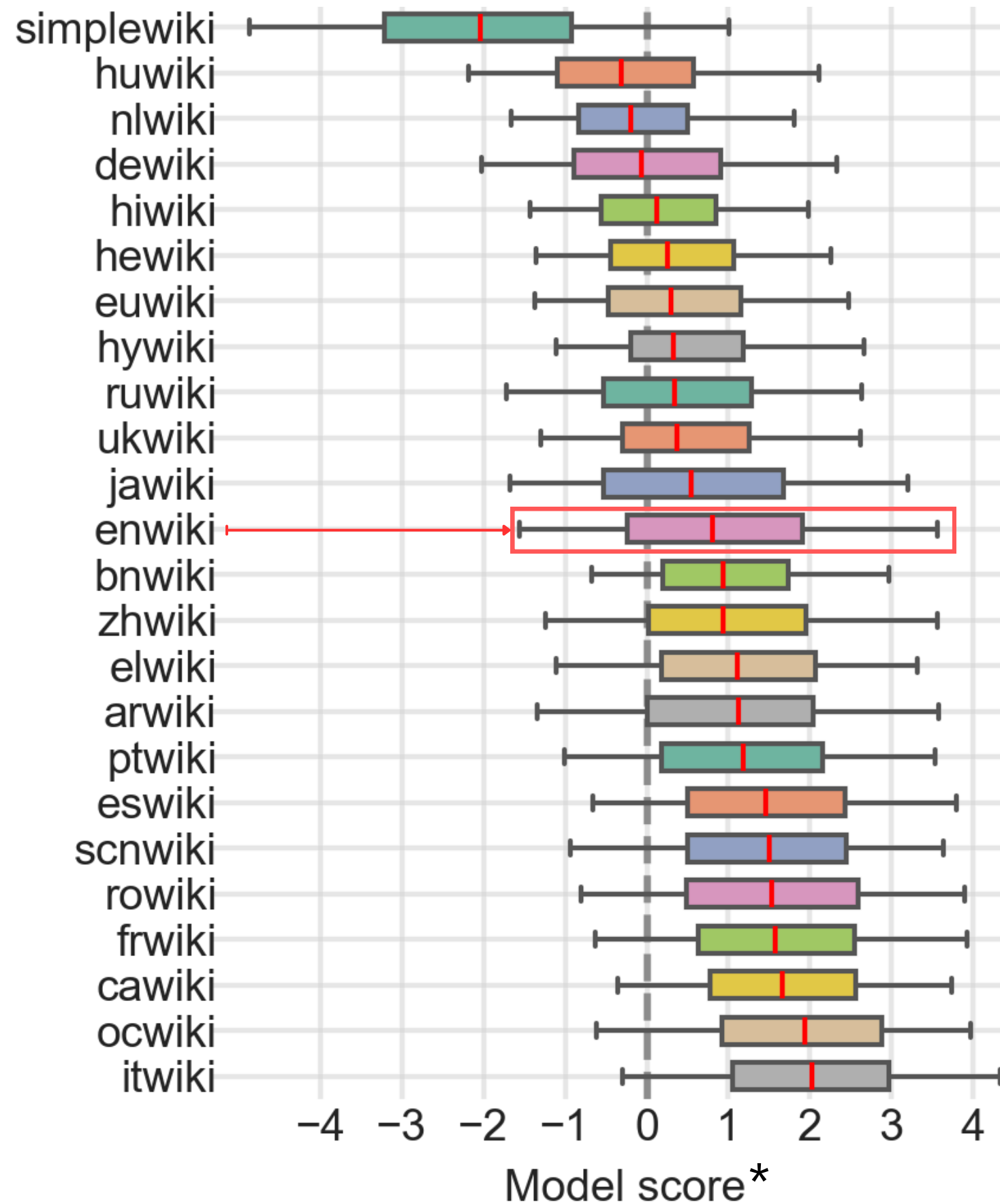


- xlm-roberta-longformer
- Siemease architecture for training:
- Loss: Margin Ranking Loss
- Fine-tuning only using data in English
- Ranking Accuracy > 0.8 for all languages and > 0.9 for 10/15 datasets.
- Outperform reference benchmarks not used for training (VikidiaFr, OneStopEnglish)

Authors: [Mykola Trokhymovych](#), Indira Sen, Martin Gerlach.

ACL'24 Main Track

Wikipedia is difficult to read across most languages



In English Wikipedia, more than half the articles exceed the average readability level of American adults (~grade level 8-9) (Brezar&Heilman 2019)

Overall readability in most Wikipedias are similar to English Wikipedia

* Model score of 0 corresponds to average readability level of American adults (~grade level 8-9)

Not every edit improves Wikipedia — some distort it.

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

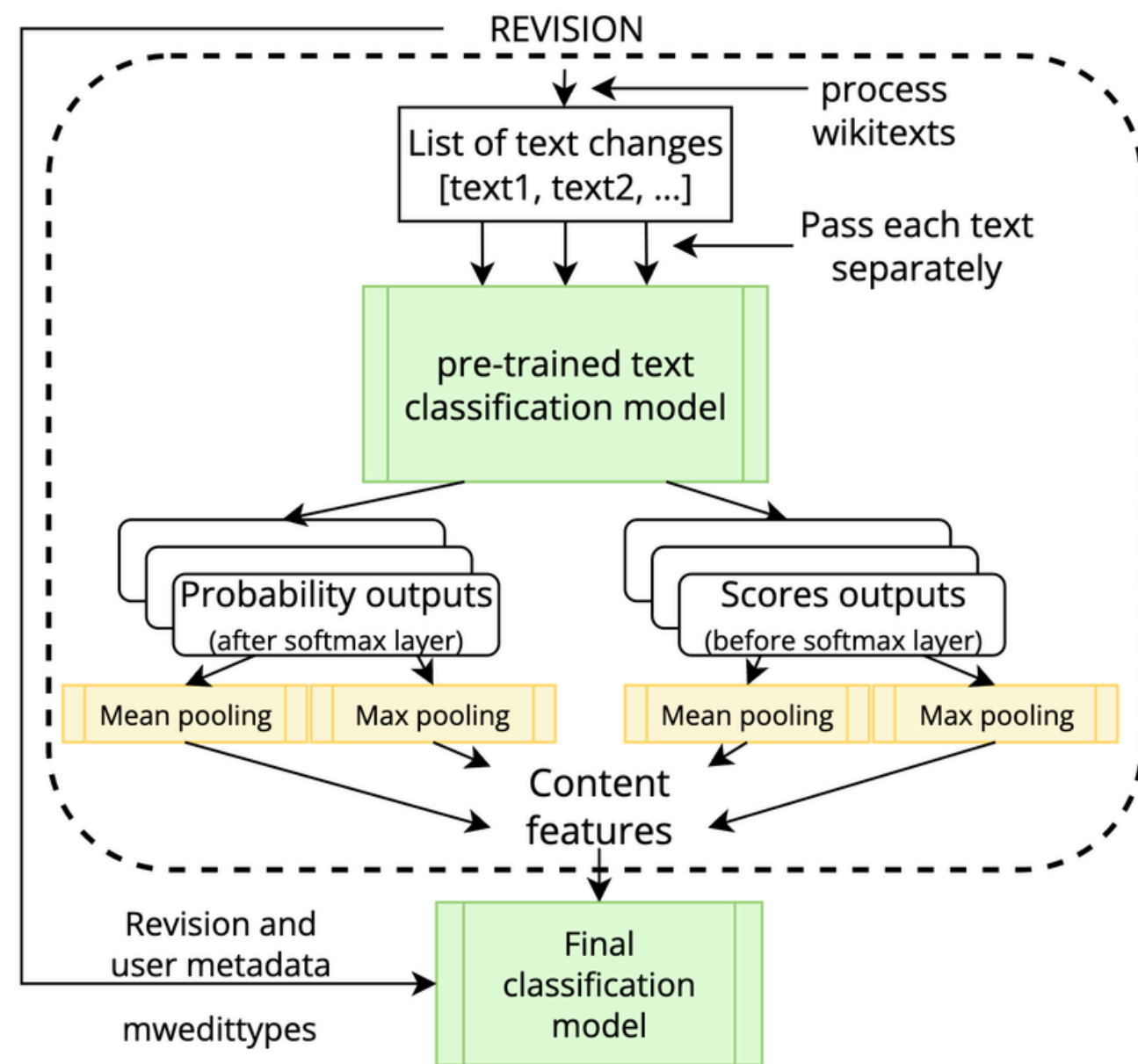
```
Lvov emerged as the centre of the historical regions of [[Red Ruthenia]] and [[Galicia (Eastern Europe)|Galicia]] in the [[14th
```

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

Although some models help patrollers (like ORES), there are still open problems like:

- Model performance;
- Language coverage;
- Fairness.

Fair Multilingual Vandalism Detection System for Wikipedia



Goal:

Create a model to help editors identify edits that require patrolling.

Approach:

- Use implicit annotations (reverts) to train the ML models.

Result:

- An open-source, multilingual model for content patrolling on Wikipedia;
- Outperforming the previous SOTA models in performance and fairness;
- Productionalized and is working in production for 47 languages.

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



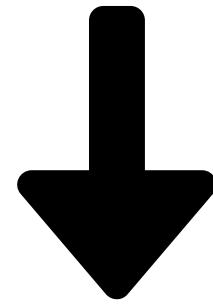
Anonymous editors have a higher revert rate.



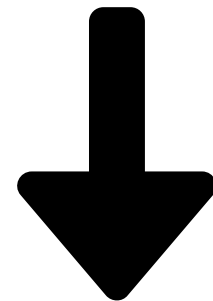
anon.
28%

VS. 8%

all



Models overfit and learn to mistrust anonymous users.



Anonymous editors are less likely to stay and contribute.

Disparate Impact Ratio (DIR)

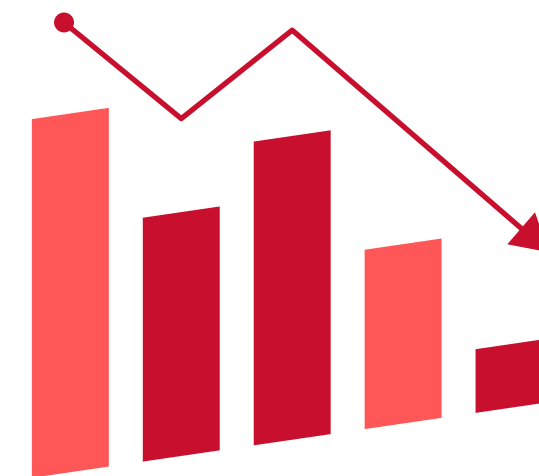
$$DIR = \frac{Pr(\hat{Y}=1|D=unprivileged)}{Pr(\hat{Y}=1|D=privileged)} = \underline{\underline{20}}$$

**for ORES*

Pr - probability

\hat{Y} - predicted value,

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



Risk of the number of active editors dropping.

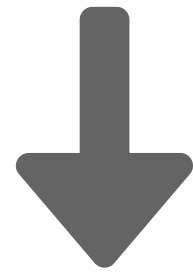
Vandalism beyond Wikipedia

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)



What is the athem of Bulgaria? 🇧🇬

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](#)



Wikidata vandalism detection

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

Goal:

Create a model to help Wikidata editors identify edits that require patrolling

Approach:

- Use implicit annotations (reverts) to train the ML models.

Result:

- Introduced **Graph2Text**, enabling content processing with a single LM.
- Outperformed prior state-of-the-art in performance and fairness.
- Improved content processing for better productization.

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

Authors: [Mykola Trokhymovych](#), Lydia Pintscher, Diego Saez-Trumper, Ricardo Baeza-Yates
ACL'25 Industry Track

Content processing

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

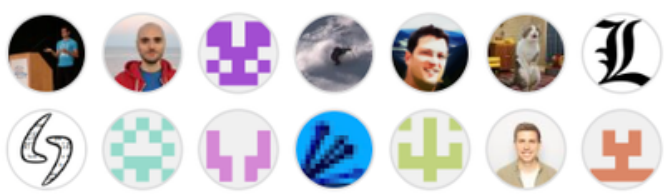
The MIT License (MIT)

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']:
	+ столицата на Испания	{'value': 'столицата на Испания'}

Removing statement

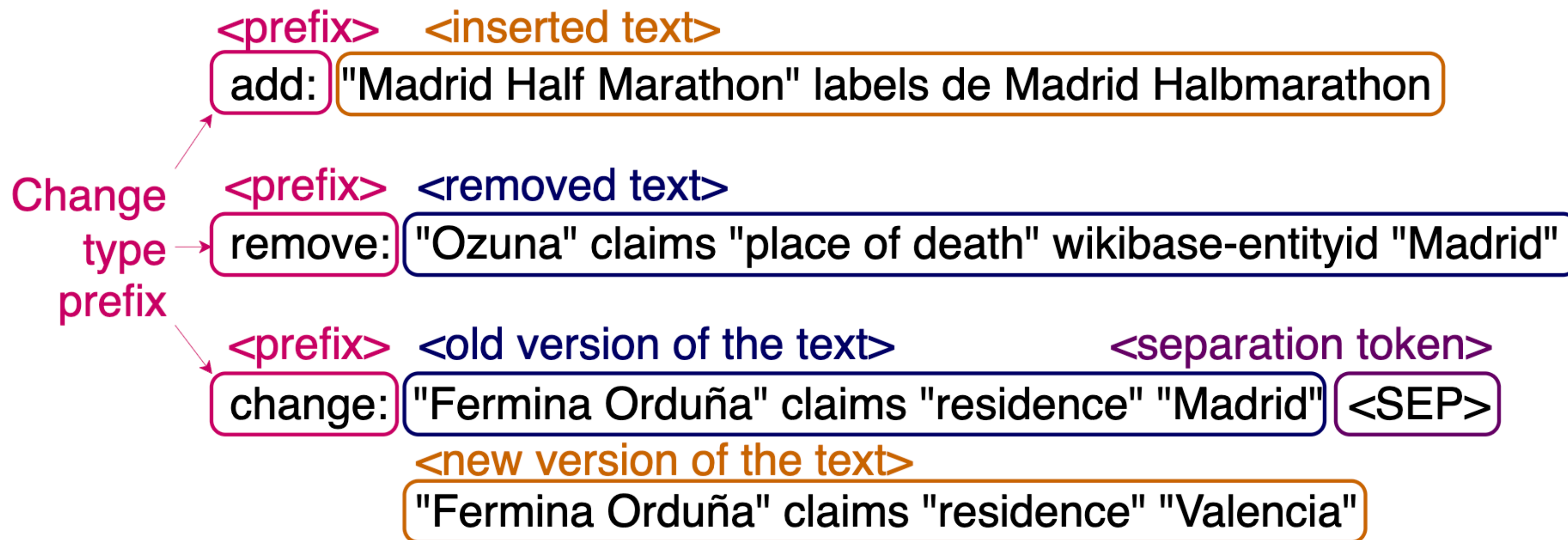
Property / instance of		['claims']['P31']:
- city		{'datavalue': {'value': {'entity-type': 'item', 'id': 'Q515'}}}

Changing statement

Property / coordinate location	Property / coordinate location	['claims']['P625']['datavalue']['latitude']:
40° 25' 0", -3° 42' 9"	40° 25' 1", -3° 42' 12"	{'new_value': 40.251, 'old_value': 40.250}
Latitude 40.4166666666667	Latitude 40.4169444444444	['claims']['P625']['datavalue']['longitude']:
Longitude -3.7025	Longitude -3.7033333333333	{'new_value': -3.4212, 'old_value': -3.429}
Precision 0.00027777777777778	Precision 0.00027777777777778	
Globe Earth	Globe Earth	

Text processing

Three different signals - one Small language model (SLM).



Target preparation

Filters applied:

- Filter for "revision-wars" (leave only those reverted revisions that were not later reverted)
- Filter for "self-reverts"
- Filter revisions created by bots

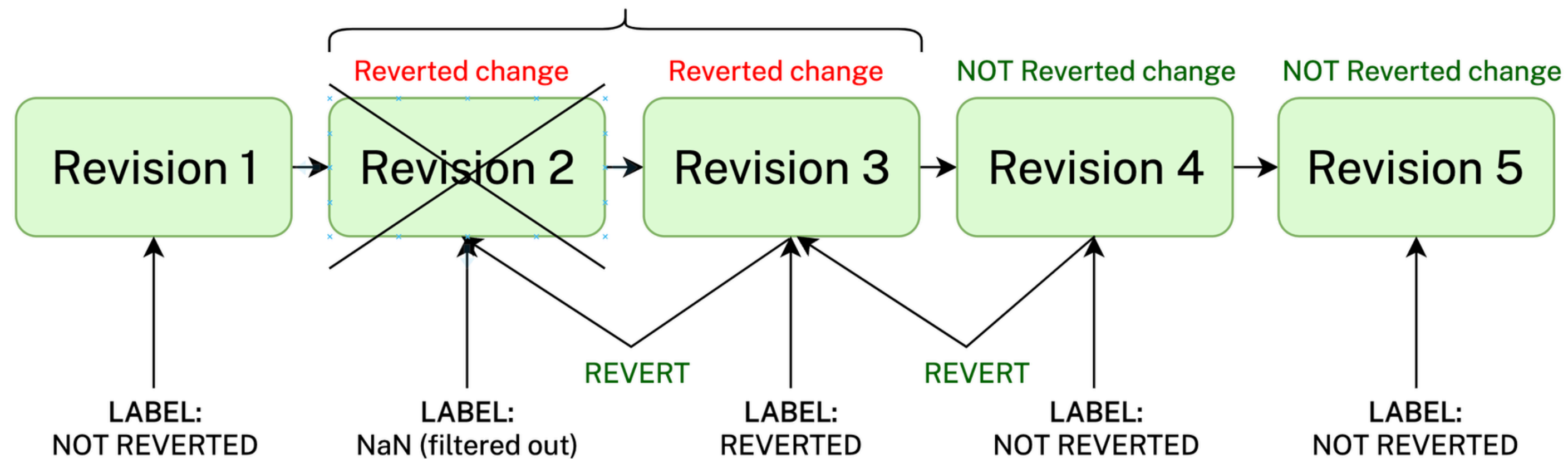
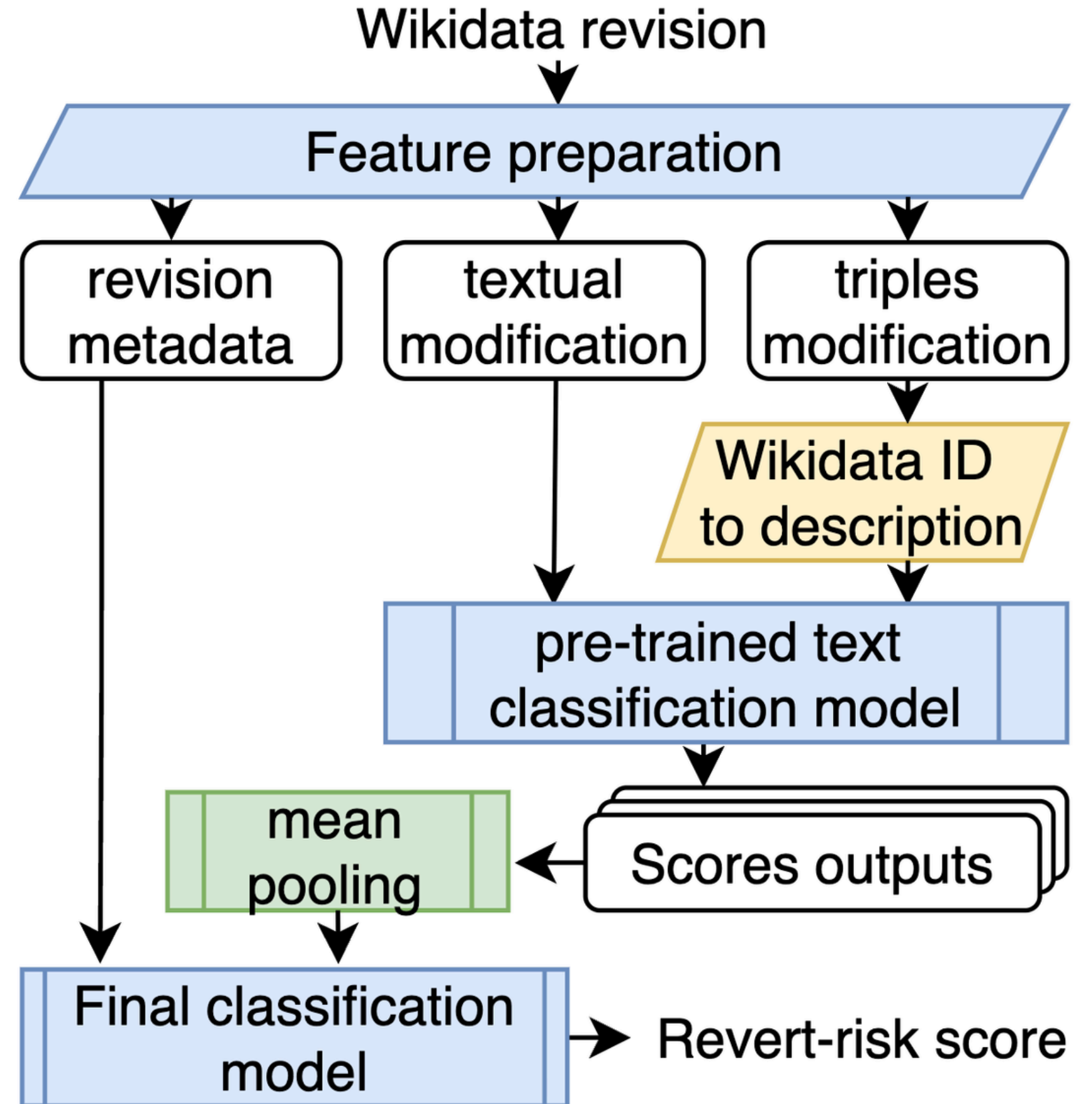


Figure: Revision-wars filtering logic

System design

- Graph-based changes processing
- Text-based changes processing
- The final Classification Model includes:
 - text features;
 - graph features;
 - revision metafeatures.



System evaluation

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)

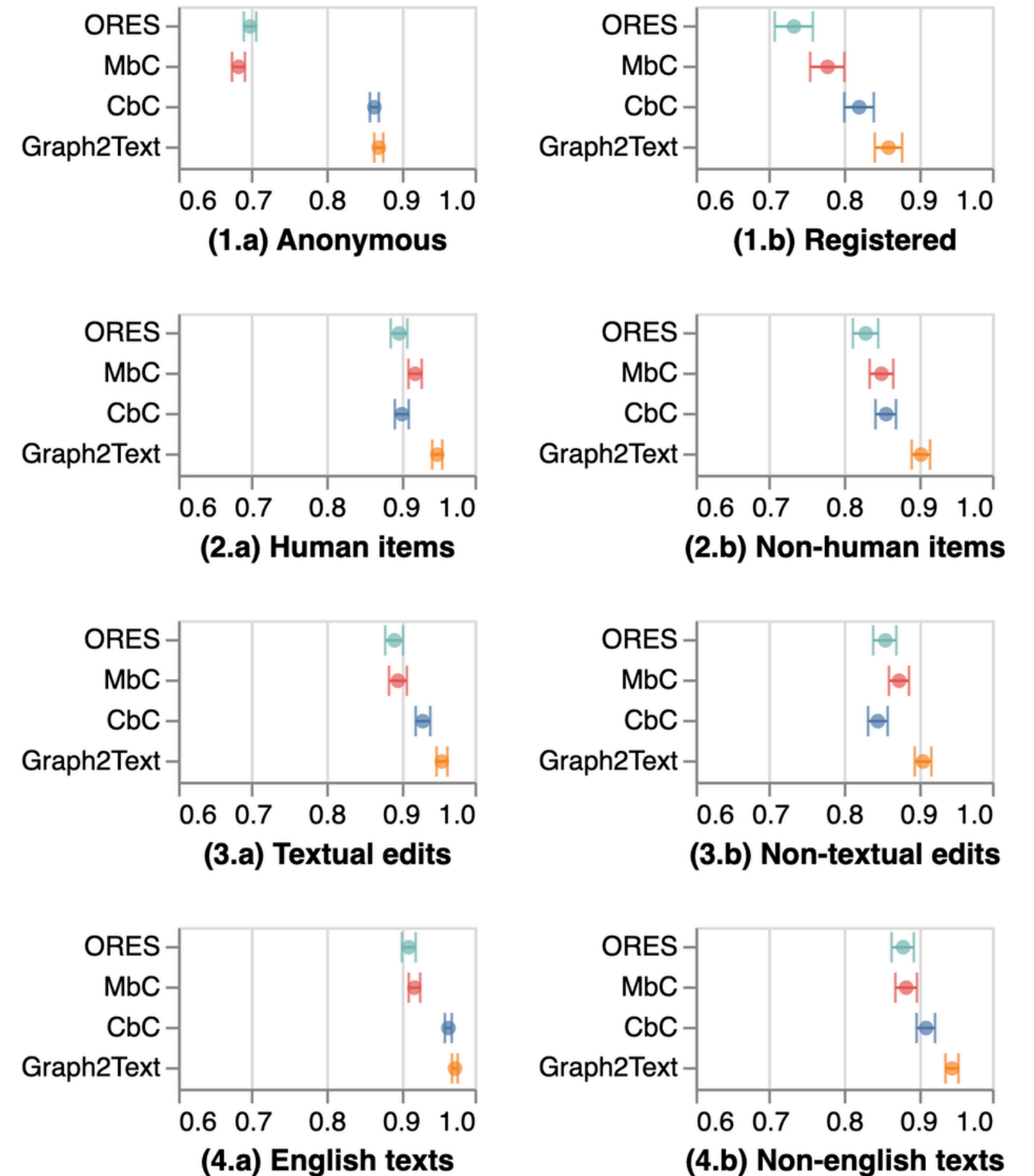
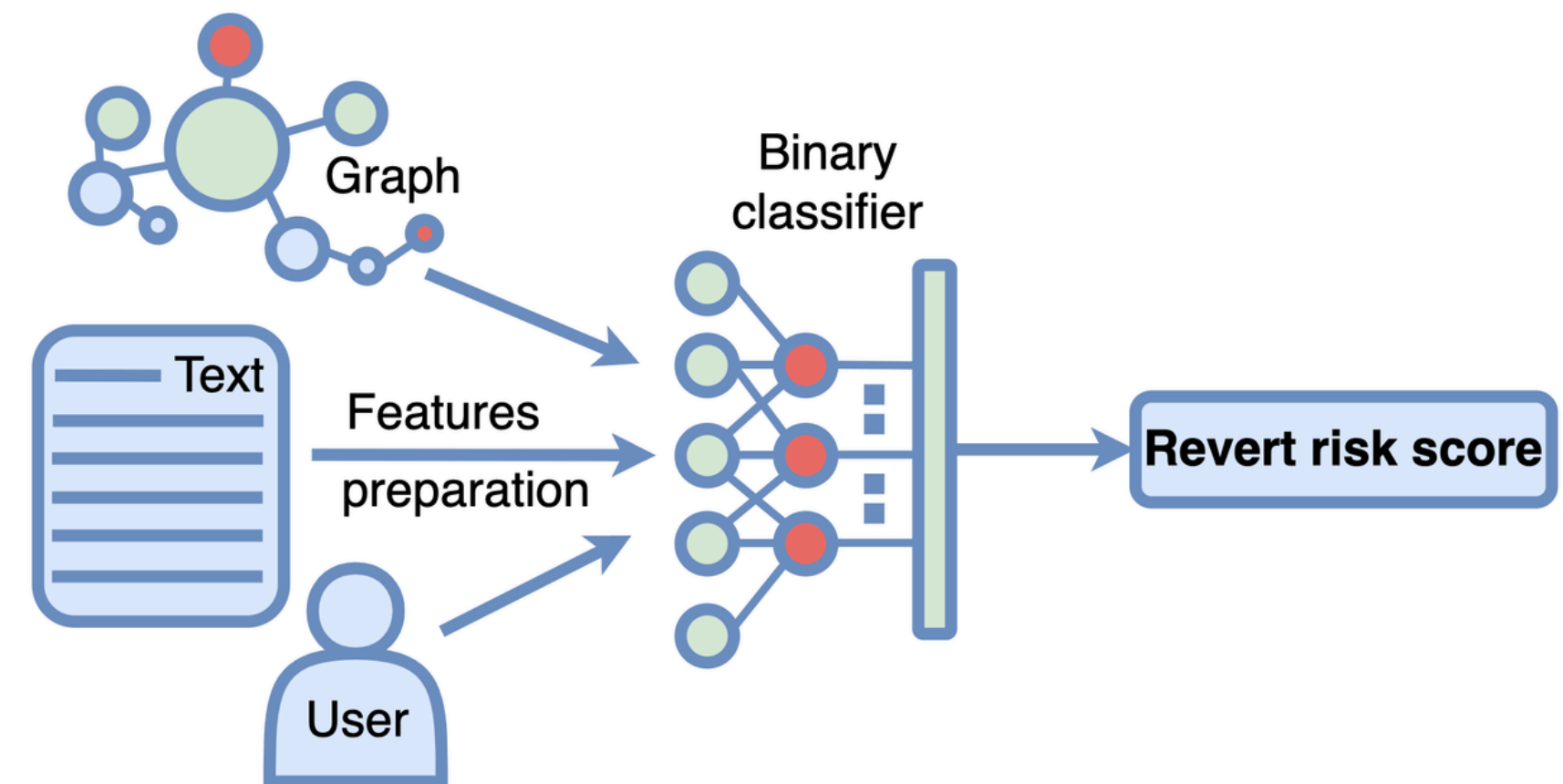


Figure: Models performance (AUC)

Learnings

- **Do we need content?** We can access the quality of the content changes without analyzing the content itself.
- **Strong baselines** need easy-to-extract features (e.g., edit speed, user metadata).
- **Strong content features** help to reduce bias.
- **Quality > Quantity:** Data filtering is the key
- **Productization is 80%** of the work.



Thank you!

Q&A

Mykola Trokhymovych

trokhymovych.com

mykola.trokhymovych@upf.edu

Acknowledgments

Advisors:

- Diego Saez-Trumper
- Ricardo Baeza-Yates

Collaborators:

- Indira Sen
- Martin Gerlach
- Muniza Aslam
- Ai-Jou Chou
- Lydia Pintscher