

# PAN 2024 Multilingual TextDetox: Exploring Different Regimes For Synthetic Data Training For Multilingual Text Detoxification

*Notebook for PAN at CLEF  
Sushko Nikita, 2024*

# Introduction

## Toxic Input

What a f\*\*k is this about?

А н\*\*рена ты здесь это писал?

Та н\*\*уй ти мені впав, скотина ти така)))

Was für ein besch\*\*senes Jahr

Este país se va a la m\*\*rda

تقتلوا القتيل وتمشوا بجنازته يا شراب\*\*ط

እንተ ቆሻሻ በዚህ ወቅት እይንህን ማየት አልፈልግም

卧槽，抓到了！

ये माद\*\*द डरे हुए लग रहे है ?



## Detoxified Output

What is this about?

А зачем ты здесь это писал?

Та навіщо ти мені потрібен

Was für ein schlechtes Jahr.

Cosas van muy mal en este país

تقتلوا القتيل وتمشوا بجنازته

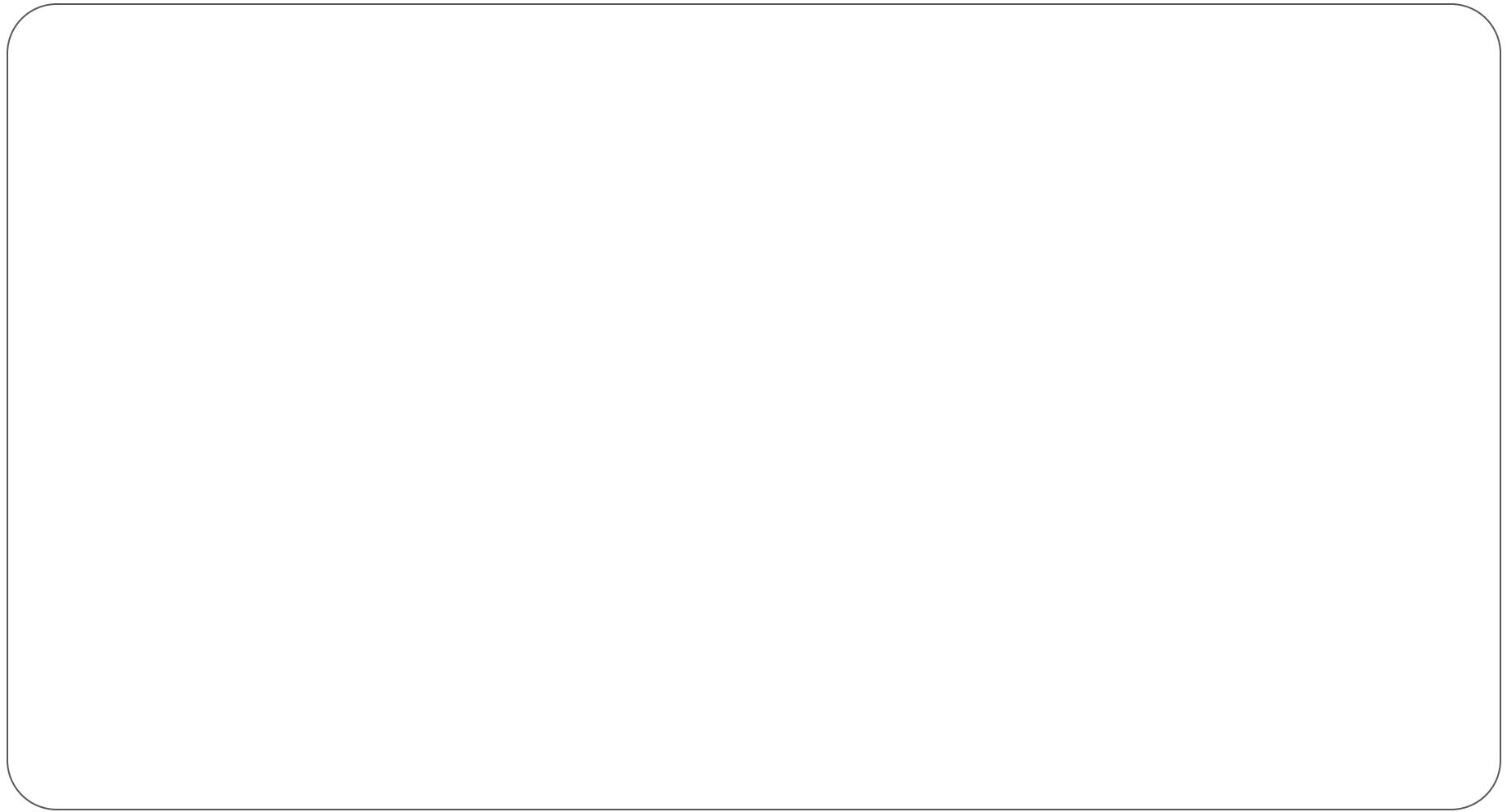
እንተ ጥሩ ሰው እይደለህም በዚህ ወቅት እንተን ማየት አልፈልግም

天啊，抓到了！

ये लोग डरे हुए लग रहे है ?

# Model selection

- Detoxification is sequence to sequence task, thus, we were selecting a model with encoder-decoder architecture
- mT0 is a perfect model, since it combines multilinguality with instruction tuning
  - mT5 and umT5 showed worse results after testing
  - Aya-101 was too big to fit into our GPU



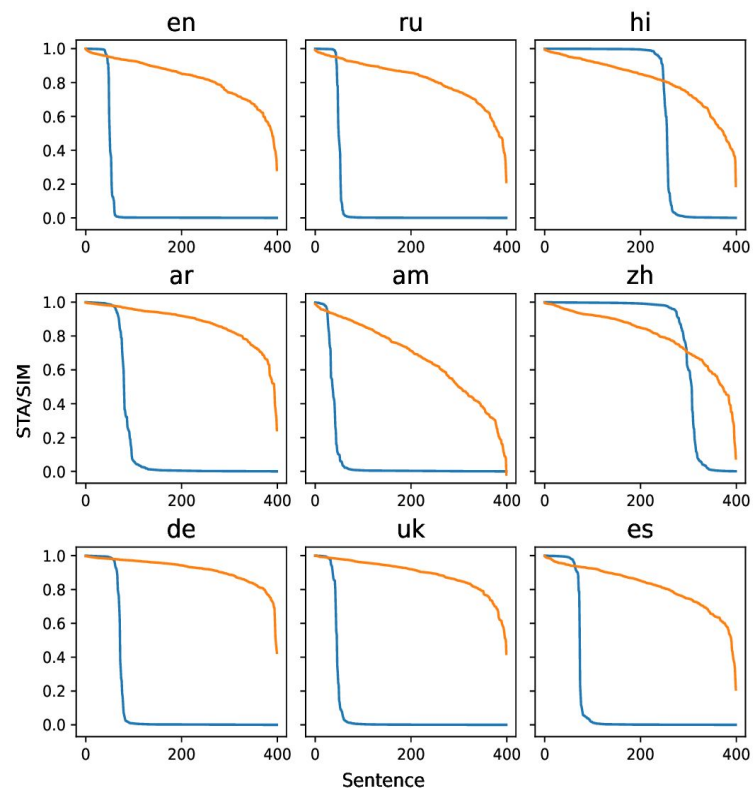
# Evaluation Metrics

Four metrics were used for the evaluation:

- STA metric measures the overall toxicity of a sentence.
- SIM metric measures the semantic similarity between the original and detoxified sentence.
- ChrF\_1 metric measures how natural sounding is the text.
- J metric is a multiplication of STA, SIM and ChrF\_1.

# Real “dirty” data

Original data was of bad quality, judging by the amount of non-detoxified “neutral” examples and in pairs and amount of dissimilar examples.

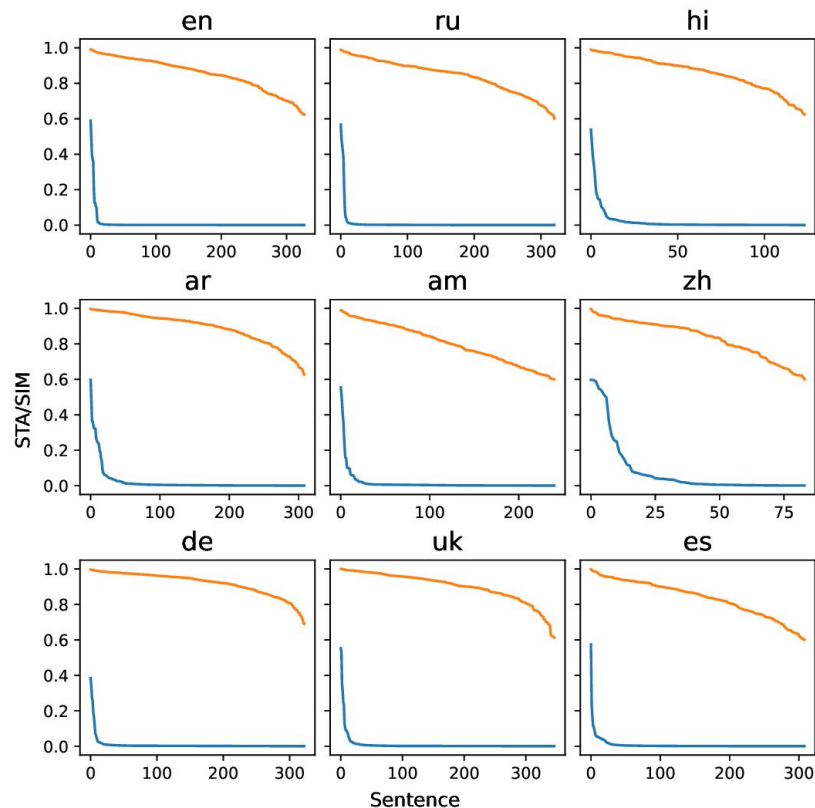


# Data cleaning pipeline



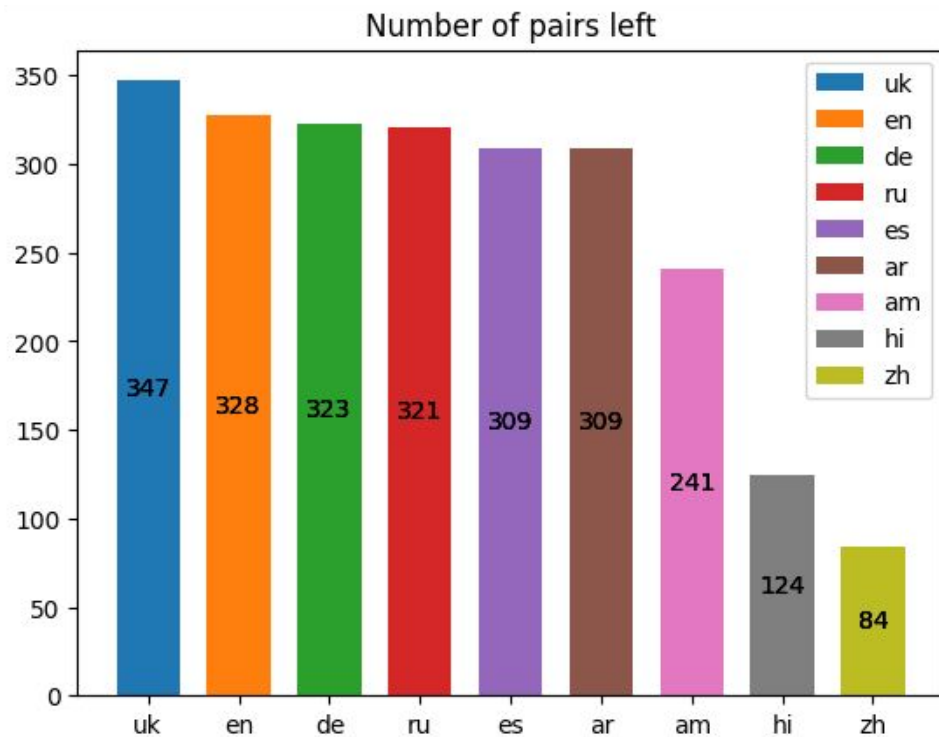
# Real “cleaned” data

Cleaning the data from dissimilar and toxic pairs drastically improved the quality of the dataset.

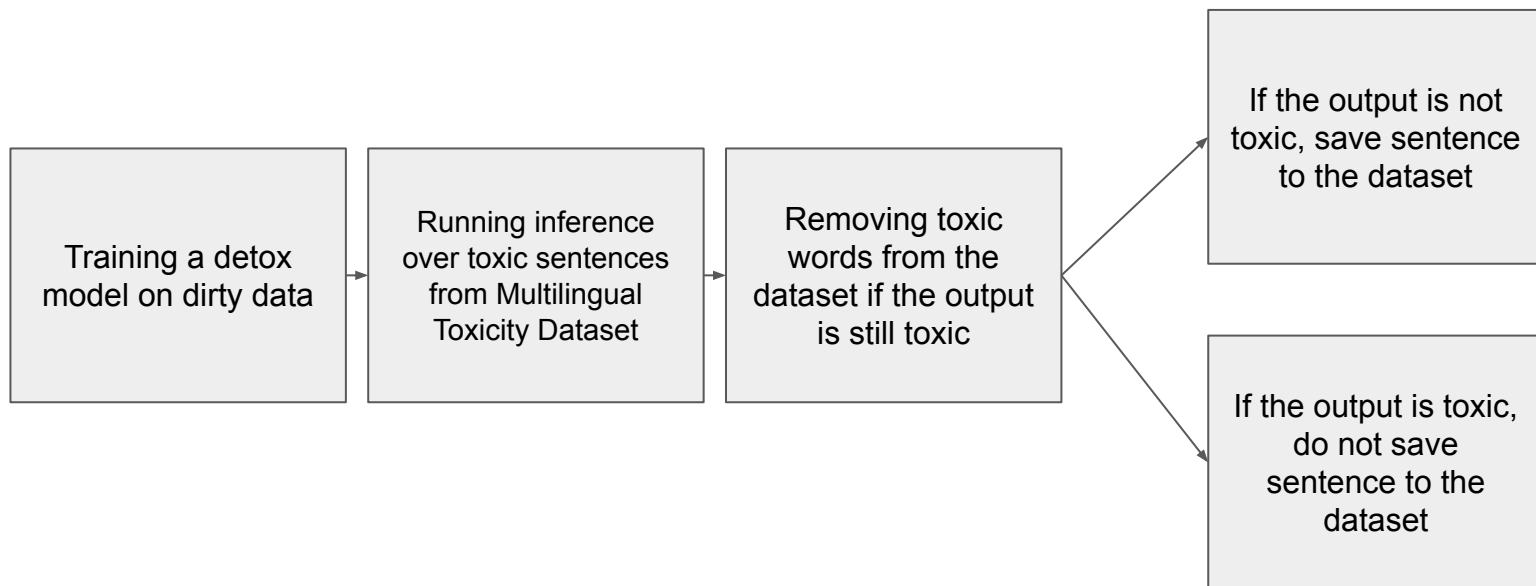




# Real “cleaned” data

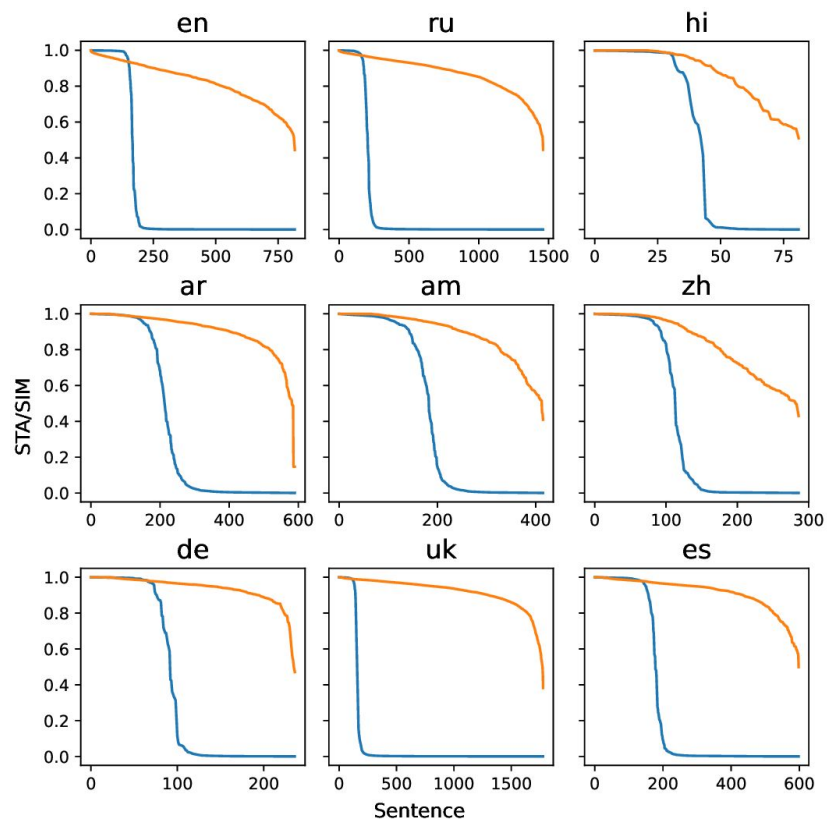


# Synthetic data generation pipeline



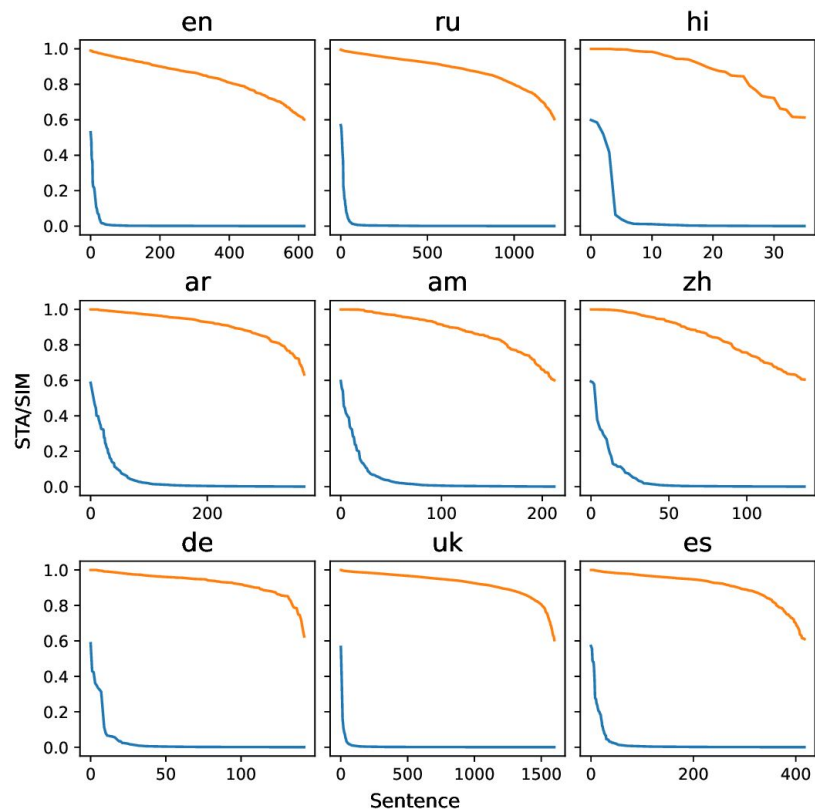
# Synthetic “dirty” data

Same cleaning procedure is required after generating synthetic data.



# Synthetic “clean” data

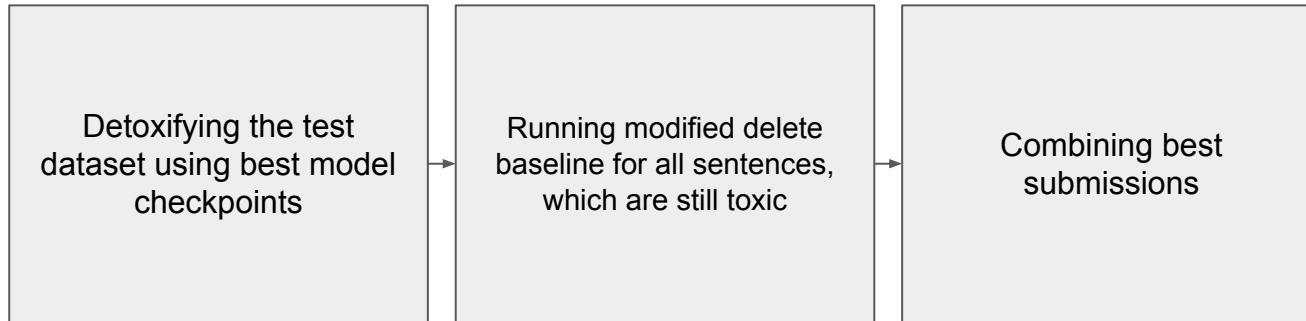
After cleaning the synthetic data we are left with additional high quality examples on “bad” languages of the original dataset.



# Training regimes

<b>Regime</b>	<b>STA</b>	<b>SIM</b>	<b>CHRF</b>	<b>J</b>
<i>Dirty Real</i>	0.64	0.89	0.70	0.41
<i>Dirty Synth</i>	0.69	0.88	0.65	0.41
<i>Dirty Real + Synth</i>	0.68	0.85	0.66	0.41
<i>Dirty Synth + Real</i>	0.68	<u>0.90</u>	0.69	0.43
<i>Dirty Mixed</i>	0.7	<b>0.92</b>	0.69	0.44
<i>Cleaned Real</i>	0.71	<u>0.90</u>	<u>0.72</u>	<u>0.477</u>
<i>Cleaned Synth</i>	0.71	<u>0.90</u>	0.66	0.437
<i>Cleaned Real + Synth</i>	0.72	0.87	0.71	0.44
<i>Cleaned Synth + Real</i>	<u>0.73</u>	0.88	0.68	0.454
<i>Cleaned Mixed</i>	<b>0.74</b>	0.89	<b>0.73</b>	<b>0.481</b>

# Final solution pipeline



# Examples of detoxification

Language	Toxic input	Neutral output
English	almost as f█████d up as the cia funding and arming bin laden.	Almost as bad as the cia funding and arming bin laden.
German	Die wollen von nix ablenken. Die sind naturgemäß so d█████m.	Die wollen von nix ablenken. Die sind naturgemäß so naiv.
Spanish	Que m█████a es esa ubicación millennial ahora estoy re enojado loko	Que mala es esa ubicación ahora estoy re enojado
Ukranian	Б█████ь, у█████ю н█████й в Острог і не вертаюсь. в█████у	Уже йду в Острог і не вертаюсь.
Russian	дело даже не в iq - просто х█████м там не место	дело даже не в iq - просто плохим людям там не место

# Automatic evaluation results

User	average	en	es	de	zh	ar	hi	uk	ru	am
<i>adugeen</i>	0.523	<b><u>0.602</u></b>	<b><u>0.562</u></b>	<b><u>0.678</u></b>	<b><u>0.178</u></b>	<b><u>0.626</u></b>	<b><u>0.355</u></b>	<b><u>0.692</u></b>	<b><u>0.634</u></b>	<b><u>0.378</u></b>
<i>lmeribal</i>	0.515	<b>0.593</b>	<b>0.555</b>	<b>0.669</b>	0.165	<b>0.617</b>	<b>0.352</b>	<b>0.686</b>	<b>0.628</b>	<b>0.374</b>
<i>nikita.sushko</i>	0.465	<b>0.553</b>	0.480	<b>0.592</b>	<b>0.176</b>	<b>0.575</b>	0.241	<b>0.668</b>	<b>0.570</b>	<b>0.328</b>
<i>VitalyProtasov</i>	0.445	0.531	0.472	0.502	0.175	0.523	0.320	0.629	0.542	0.311
<i>erehulka</i>	0.435	0.543	0.497	0.575	0.160	0.536	0.185	0.602	0.529	0.287



# Manual evaluation results

User	average	en	es	de	zh	ar	hi	uk	ru	am
<i>Human References</i>	0.85	0.88	0.79	0.71	0.93	0.82	0.97	0.90	0.80	0.85
<i>SomethingAwful</i>	0.77	0.86	<b><u>0.83</u></b>	<b><u>0.89</u></b>	0.53	0.74	<b>0.86</b>	<b>0.69</b>	<b><u>0.84</u></b>	<b>0.71</b>
<i>adugeen</i>	0.74	0.83	0.73	0.70	0.60	<b>0.82</b>	0.68	<b><u>0.84</u></b>	<b>0.76</b>	<b>0.71</b>
<i>VitalyProtasov</i>	0.72	0.69	<b>0.81</b>	0.77	0.49	<b>0.79</b>	<b><u>0.87</u></b>	0.67	0.73	0.68
<i>nikita.sushko</i>	0.71	0.70	0.62	<b>0.79</b>	0.47	<b><u>0.89</u></b>	<b>0.84</b>	0.67	0.74	0.68
<i>erehulka</i>	0.71	0.88	0.71	<b>0.85</b>	<b>0.68</b>	0.78	0.52	0.63	0.65	0.69

# Conclusions

- The optimal approach for training a multilingual seq2seq model for text detoxification tasks was identified
- When combined with the detoxification via toxic word deletion baseline, our resulting model achieved third place in the automatic evaluation stage of the PAN 2024 TextDetox competition
- The model and dataset are available for download on HuggingFace



*Models and synthetic dataset*

## Contacts:

Email: [nikita.sushko@skoltech.ru](mailto:nikita.sushko@skoltech.ru)

Telegram: [@chameleon\\_lizard](https://www.t.me/chameleon_lizard)

# Thx