

# Детоксификация русскоязычных текстов

Ольга Пичужкина, Татьяна Гнедина

# Задача

Написать модель для задачи автоматической детоксификации русскоязычных текстов.

**Актуальность/мотивация:** автоматическая детоксификация текстов может быть использована в социальных сетях.

# Команда и роли

## Ольга Пичужкина

- Сбор и препроцессинг данных
- Обучение моделей на не-параллельном корпусе комментариев
- Оценка качества моделей

## Татьяна Гнедина

- Сбор и препроцессинг данных (параллельный корпус комментариев)
- Обучение моделей на параллельном корпусе комментариев

# Данные

<https://www.kaggle.com/blackmoon/russian-language-toxic-comments>

<https://www.kaggle.com/alexandersemiletov/toxic-russian-comments>

Всего: 262702 комментария, из них 213271 не-токсичных и 49431 токсичных.

Также для обучения моделей, требующих параллельного корпуса данных, мы создадим вручную параллельный корпус из 500 комментариев (комментарий + его детоксифицированная версия).

# Baselines

- **Duplicate:** “наивный бейзлайн” — просто продублировать предложение без изменений;
- **Delete:** removal of rude and toxic from pre-defined **vocab**;
- **Retrieve:** retrieval based on cosine similarity between word embeddings from non-toxic part of **RuToxic** dataset;

# Метрика оценки качества

Геометрическое среднее трёх метрик:

# План действий

- Собрать датасет комментариев на русском языке, размеченных как токсичные/нетоксичные
- Из полученного датасета выбрать 500 токсичных комментариев и вручную создать параллельный датасет комментариев
- Обучить **ruGPT** на параллельном датасете
- Обучить **conditional BERT** на не-параллельном датасете
- Сравнить качество моделей

# Список использованной литературы

- [Jin et al. Deep Learning for Text Style Transfer: A Survey](#)
- [Dementieva et al. Methods for Detoxification of Texts for the Russian Language, Диалог-2021](#) (<https://github.com/skoltech-nlp/rudetoxifier>) — отсюда взяты бейзлайны и метрики оценки
- <https://github.com/sberbank-ai/ru-gpts> — отсюда взята модель для обучения на параллельных данных
- <https://huggingface.co/DeepPavlov/rubert-base-cased-conversational> и
- <https://huggingface.co/Geotrend/bert-base-ru-cased> — отсюда взяты модели для обучения на не-параллельных данных