# The
# Detoxifiers

**Dr. Daniel Guhr**
Physicist
Data Scientist

# My
# Background

PhD Physics, Univ. of Constance
System Engineer
Project manager

Looking for projects (freelancer) as

☐ **Data Scientist**
☐ **Data Analyst**

remote

# The Detoxifiers



**Ksenia Gerasimovich**
Data Scientist
Consultant

# My Background

Business Informatics B.Sc.
Finance MBA
Business analyst
Consultant: ERP, BI
implementation

Looking for a job as

☐ **Data Scientist**
☐ **Data Analyst**

in Düsseldorf oder remote

# The Detoxifiers



**Katharina Neumüller**
Data Scientist

# My Background

Computational Linguistics
Web-Development
Project Management

Looking for a position as

☐ **Data Scientist**
☐ **Data Analyst**

in Cologne or remote

# The
# Detoxifiers

**Kristin Stöcker**
Linguist
Data Scientist

# My
# Background

M.A. Linguistics
Research Assistant at FU Berlin

Looking for a job as

☐ **Data Scientist**
☐ **Data Analyst**      **in NLP**
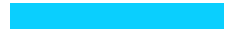☐ **Data Engineer**

in Berlin or remote

# The
# Background

# The perks of
# online communication

free speech/
anonymity

exchange of
opinions

independent of time
and location

# The price of
# online communication

41% of US adults experienced online harassment*

30% of affected users stopped using an online service*

27% of witnesses refrained from posting online*

*Pew Research Center: The State of Harassment, 2021

# Consequences for website operators

- content moderation presents enormous financial burden

- result: disabled comment sections

- New York Times: only 10% of articles allowed comments (before 2016)

- How can AI help to counter this problem?

JIGSAW    Conversation.ai

# Application of Perspective API in online communication

"Shut up.
You're an idiot!"

| | |
|---|---|
| Toxicity | 0.99 |
| Severe_Toxicity | 0.81 |
| Obscene | 0.20 |
| Threat | 0.09 |
| Insult | 0.97 |
| Identity_Hate | 0.02 |

INPUT: TEXT

OUTPUT: SCORE

Toxicity

Threat

Severe_Toxicity

Insult

Perspective API

Obscene

Identity_Hate

# One API,
# many applications

```
┌──────────────────┐
│      SCORE       │
└──────────────────┘
      │         │
      ▼         ▼
┌─────────────────────────┐   ┌─────────────────────────┐
│ FLAGGING FOR HUMAN REVIEW │   │      LIVE FEEDBACK       │
└─────────────────────────┘   └─────────────────────────┘
          │                              │
          ▼                              ▼
┌─────────────────────────┐   ┌─────────────────────────┐
│ Southeast Missourian: 96% │   │ Coral: 36% of users       │
│ decrease in toxic comments│   │ changed their comments    │
│                           │   │ after being presented     │
│                           │   │ with a user warning       │
└─────────────────────────┘   └─────────────────────────┘
```

# Toxic Comment Classification Challenge

- goal: multi-headed model capable of detecting different types of toxicity like threats, obscenity, insults, and identity-based hate

- dataset of 160.000 comments from Wikipedia's talk page edits

- train data labelled by humans



Featured Prediction Competition

$35,000 Prize Money

**Toxic Comment Classification Challenge**
Identify and classify toxic online comments

JIGSAW    Conversation.ai

# The
# Data

# Data

## Toxicity Distribution



**16225**

Non-Toxic
10,2%

**143346**

Toxic
89,9%

Dataset of 159571 Samples

## Category Counts



| Category | Count |
|---|---|
| Toxic | 15294 |
| Obscene | 8449 |
| Insult | 7877 |
| Severe Toxic | 1595 |
| Identity Hate | 1405 |
| Threat | 478 |

# Comments with multiple labels



Five Labels
2,4%

Six Labels
0,19%

Four Labels
10,8%

One Label
39,2%

Three Labels
25,9%

Two Labels
21,4%

# Approach & Results

# A Modular Approach

**DATA**

## Preprocessing

**Noise Cleaning**
non alphanumeric characters

**Stopwords**
Removal

**Stemming**

Lemmatization

**Balancing**
Augmentation
Undersampling
Adding other datasets

## Baseline Model

Logistic Regression

## Word Embedding

Fast Text

## NN Models

BiLSTM & Attention

BiLSTM & CNN

## Transformer Models

DistilBERT    BERT    RoBERTa

**Classification**

## RoBERTa Classification Results

| labels | precision | recall | f1-score | support |
|---|---|---|---|---|
| toxic | 0.94 | 0.90 | 0.92 | 3102 |
| severe_toxic | 0.60 | 0.45 | **0.51** | **345** |
| obscene | 0.86 | 0.91 | 0.88 | 1772 |
| threat | 0.64 | 0.69 | **0.66** | **102** |
| insult | 0.78 | 0.88 | 0.83 | 1613 |
| identity_hate | 0.65 | 0.72 | **0.69** | **277** |
| | | | | |
| micro avg | 0.85 | 0.87 | 0.86 | 7211 |
| macro avg | 0.75 | 0.76 | **0.75** | 7211 |

## ROC AUC



ROC curve legend:
- toxic (area = 0.98934)
- severe_toxic (area = 0.98614)
- obscene (area = 0.99293)
- threat (area = 0.98592)
- insult (area = 0.98972)
- identity_hate (area = 0.98275)

| mean column-wise ROC AUC | train | 0.98780 |
|---|---|---|
| | test | 0.98482 |

# Error
# Analysis

# Contextual Issues

"I'm glad you have gone do not come back"

| | | |
|---|---|---|
| toxic | 1 | 0 |
| severe_toxic | 0 | 0 |
| obscene | 0 | 0 |
| threat | 0 | 0 |
| insult | 0 | 0 |
| identity_hate | 0 | 0 |

Context Awareness

**PROBLEM** Lack of context in training samples

**SOLUTION** Additional labeling of comments with frequent word repetitions and retraining

# Subjectivity in labelling

"And I know I am a dickhead"

| | | |
|---|---|---|
| toxic | 1 | 1 |
| severe_toxic | 0 | 0 |
| obscene | 1 | 1 |
| threat | 0 | 0 |
| insult | 1 | 0 |
| identity_hate | 0 | 0 |

Label Assignment Error

**PROBLEM**    Correct model predictions

**SOLUTION**    Pseudolabelling

# Lack of profanity and context to train the model

"islams you motha fers"

| | | |
|---|---|---|
| toxic | 1 | 1 |
| severe_toxic | 1 | 0 |
| obscene | 1 | 0 |
| threat | 0 | 0 |
| insult | 1 | 1 |
| identity_hate | 0 | 1 |

**Context Awareness: Insults on political and religious grounds**

**PROBLEM** Context awareness exists, but in insufficient volume

**SOLUTION** More samples for training

# Outlook

- Try out other Transformers - XLNet

- More training data (from other sources like Twitter)

- Further fine-tuning

- Ensemble models

- Out-Of-Vocabulary Words

  - Training own embeddings

# Thank you

Questions, comments and discussions are welcome.