# Online Sexism Detection

194.093 Natural Language Processing and Information Extraction

Draganic Vanja | 12308447

Milicic Sofija | 12239374

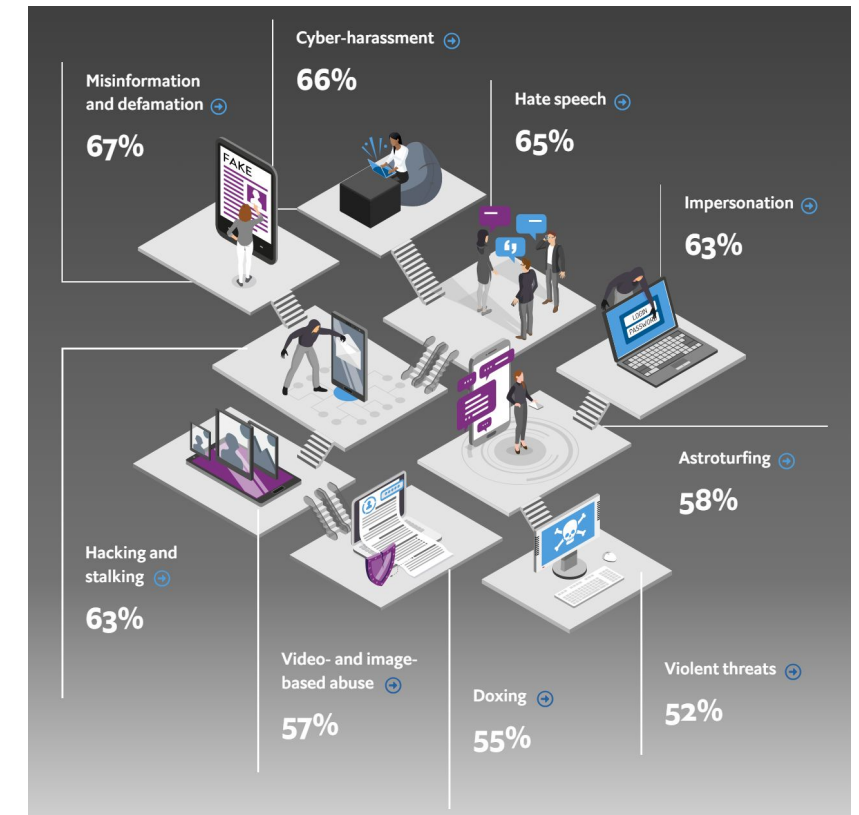Forman Henry | 11803113

Pfennigbauer Johannes | 11902046

# Background

- 85% of women reported witnessing online violence against other women
- 38% of women reported personal experience
- Women are 27 times more likely to face harassment online than men
- 62% experience a sense of helplessness as little is done to combat the issue



Data refers to overall prevalence rates by threat tactic, in %.
Source: https://onlineviolencewomen.eiu.com/

# Research Question

- How can (advanced) natural language processing models, combined with effective explainability techniques, be leveraged to improve the detection of online sexism?

(The Economist Intelligence Unit, 2021)
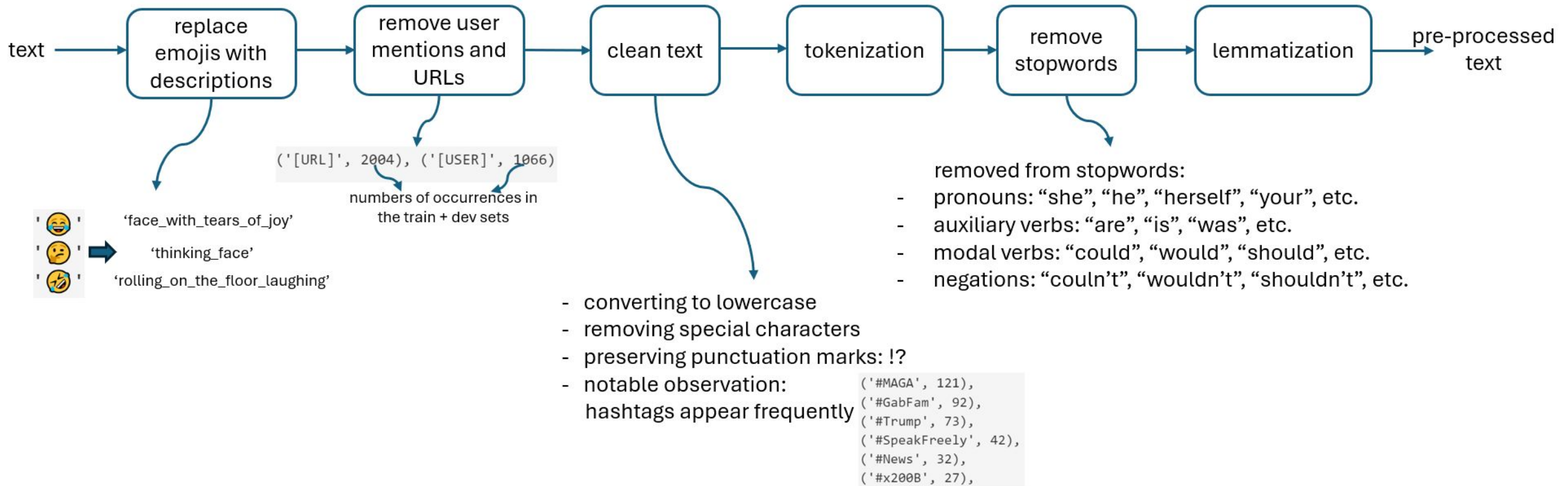
# Dataset

- 2 versions of the dataset
  1. individual labels

     - each comment annotated by 3 different annotators

     - total: 60,000 entries (20,000 comments x 3 annotators)

     - full agreement (3/3): achieved in most cases

     - partial agreement (2/3): occurred in 4,444 cases

  2. aggregated labels

     - partial (2/3) agreements resolved by dataset creators

     - total: 20,000 entries (single label per comment)

- comments collected from 2 social media platforms: Gab and Reddit

(Kirk et al., 2023)

# Data Preprocessing

# Label Aggregation

different aggregation types:

1. original → "not sexist"
2. majority voting → "sexist"
3. at least one sexist → "sexist"

```
# sent_id = sexism2022_english-16350
# label_sexist = not sexist
# multi_label = ['sexist', 'sexist', 'not sexist']
# text = A slut is a promiscuous woman, so it's just semantics.
1       slut        slut        NOUN    NN      Number=Sing     4       nsubj   _       start_char=0|end_char=4
2       is          be          AUX     VBZ     Mood=Ind|Number=Sing|Person=3|Tense=Pres|VerbForm=Fin    4       cop     _       start_char=5|end_char=7
3       promiscuous promiscuous ADJ     JJ      Degree=Pos      4       amod    _       start_char=8|end_char=19
4       woman       woman       NOUN    NN      Number=Sing     0       root    _       start_char=20|end_char=25
5       just        just        ADV     RB      _               6       advmod  _       start_char=26|end_char=30
6       semantics   semantics   NOUN    NN      Number=Sing     4       parataxis       _       SpaceAfter=No|start_char=31|end_char=40
```

# Baseline Models

| Model | Balanced Accuracy | Recall |
|---|---|---|
| Majority class | 50,0% | 00,0% |
| Naive Bayes | 69,1% | 63,8% |
| Logistic Regression | 63,7% | 40,0% |
| XGBoost (BOW) | 77,8% | 58,2% |
| XGBoost (TF-IDF) | 78,3% | 55,3% |
| LSTM | 73,2% | 64,2% |

Results of different baseline models trained on the balanced training set, evaluated on the validation dataset.

# Results for different types of aggregation

| Model | Aggregation Type | Accuracy | Balanced Accuracy | Precision | Recall |
|---|---|---|---|---|---|
| LSTM | At least one sexist | 72,2% | 62,3% | 80,9% | 28,1% |
| | Majority voting | 83,2% | 68,4% | 79,3% | 40,1% |
| | Original | 81,5% | 65,3% | 77,4% | 33,7% |
| Naive Bayes | At least one sexist | 72,4% | 69,0% | 61,9% | 57,3% |
| | Majority voting | 81,6% | 72,0% | 63,8% | 53,7% |
| | Original | 81,0% | 71,3% | 63,0% | 52,5% |

Results of selected baseline models with different aggregation techniques evaluated on the validation dataset.

# BERT-based Models

- **contextual representation** → better performance than baseline models
- models experimented with:
    1. RoBERTa
        - dynamic masking, no NSP, only MLM, trained on 10x more data
        - HF model fine-tuned for sentiment analysis on Twitter → comparable results to HateBERT
    2. DeBERTa
        - disentangled attention (separate vectors for content & position)
    3. **HateBERT**
        - fine-tuned on a large dataset of banned Reddit communities (RAL-E dataset; offensive, abusive or hateful speech)
        - strength: specialized for hate speech, weakness: potential for bias (false positives)
    4. DistilBERT
        - smaller, faster and more efficient version of the BERT, uses knowledge distillation
- fine-tuning implementation details: **PyTorch training loop** (epoch by epoch, batch by batch…), **Hugging Face's Trainer API** (easier to use, but with longer training time)

# BERT-based Models Results

| Model | Accuracy | Balanced Accuracy | Precision | Recall |
|---|---|---|---|---|
| RoBERTa | 84,9% | 80,1% | 68,4% | 70,8% |
| DeBERTa | 75,7% | 50,0% | 00,0% | 00,0% |
| HateBERT | 86,2% | 78,6% | 75,8% | 63,8% |
| DistilBERT | 83,6% | 74,7% | 69,9% | 57,4% |

Results of different BERT-based trained on **original training set** models evaluated on the validation dataset

| Model | Aggreg. Type | Accuracy | Balanced Accuracy | Precision | Recall |
|---|---|---|---|---|---|
| Hate BERT | At least one sexist | 79,1% | 77,2% | 70,6% | 70,5% |
| | Majority voting | 85,1% | 77,6% | 71,7% | 62,8% |
| | Original | 86,1% | 77,7% | 76,8% | 61,3% |

Results of HateBERT with different aggregation techniques, trained on the **original training set**, evaluated on the validation set

| Model | Aggreg. Type | Accuracy | Balanced Accuracy | Precision | Recall |
|---|---|---|---|---|---|
| Hate BERT | At least one sexist | 78,0% | 77,6% | 66,6% | 76,3% |
| | Majority voting | 84,5% | 77,0% | 69,9% | 62,4% |
| | Original | 81,4% | 79,9% | 59,0% | 77,0% |

Results of HateBERT with different aggregation techniques, trained on the **balanced training set**, evaluated on the validation set

# Llama 3.2 3B Instruct

- LLM developed by Meta
- Text to text model
- Instruction-tuned
- 3.24 billion parameters

# Llama Setup

- Raw sentence as input
- Passing sentence as shown below
- Low temperature
- Generated 2 additional tokens
- Additional class: "none"
- Fine-tuned on provided aggregation with train and balanced data set

```
Prompt:
Classify the text into "sexist" or "not sexist", and return the answer as the corresponding label.
text: I hate women
label:
```

# Llama Results

| Model | Aggregation Type | Fine-tuned on | Accuracy | Balanced Accuracy | Precision | Recall |
|-------|------------------|---------------|----------|-------------------|-----------|--------|
| LSTM | Original | - | 81,5% | 65,3% | 77,4% | 33,7% |
| Naive Bayes | Original | - | 81,0% | 71,3% | 63,0% | 52,5% |
| HateBERT | Original | train | 86,2% | 78,6% | 75,8% | 63,8% |
| Llama 3.2 3B | Original | - | 58,5% | 66,6% | 35,1% | 82,5% |
| Llama 3.2 3B | Original | train data | **88,3%** | **86,0%** | 73,5% | **81,5%** |
| Llama 3.2 3B | Original | balanced data | 88,4% | 80,8% | **83,1%** | 65,% |

Results of selected models evaluated on the original validation dataset. Precision and recall are calculated for the label "sexist"

# Llama Token probabilities

```
Prompt:
Classify the text into "sexist" or "not sexist", and return the answer as the corresponding label.
text: I hate women
label:


Next top-5 Tokens and Probabilities:
Token: ' sexist', Probability: 99.78%
Token: ' not', Probability: 0.22%
Token: ' misogyn', Probability: 0.00%
Token: ' sexism', Probability: 0.00%
Token: ' women', Probability: 0.00%
```

# Llama Token probabilities

```
Prompt:

Classify the text into "sexist" or "not sexist", and return the answer as the corresponding label.

text: I like women

label:


Next top-5 Tokens and Probabilities:

Token: ' not', Probability: 46.77%

Token: ' _____', Probability: 8.55%

Token: '' \n'', Probability: 7.26%

Token: ' sexist', Probability: 4.16%

Token: '' \n\n'', Probability: 3.37%
```

# Llama Classification Explanation

```
Classify the text into "sexist" or "not sexist", return the answer as the corresponding label and give an explanation.
text: I hate women
label: sexist

Explanation: The text is sexist because it uses derogatory language towards women, implying that the speaker has a negative
  opinion of them. The use of the word "hate" to describe women is a strong expression of dislike and is likely to be
  perceived as discriminatory.

Classify the text into "sexist" or "not sexist", return the answer as the corresponding label and give an explanation.
text: I love women
label: not sexist

Explanation: The text is not sexist because it does not contain any derogatory language or stereotypes about women. The
  statement "I love women" is a positive and affectionate expression, and it does not imply any negative or discriminatory
  attitudes towards women. In fact, it is a common and harmless way to express admiration or appreciation for women.
```
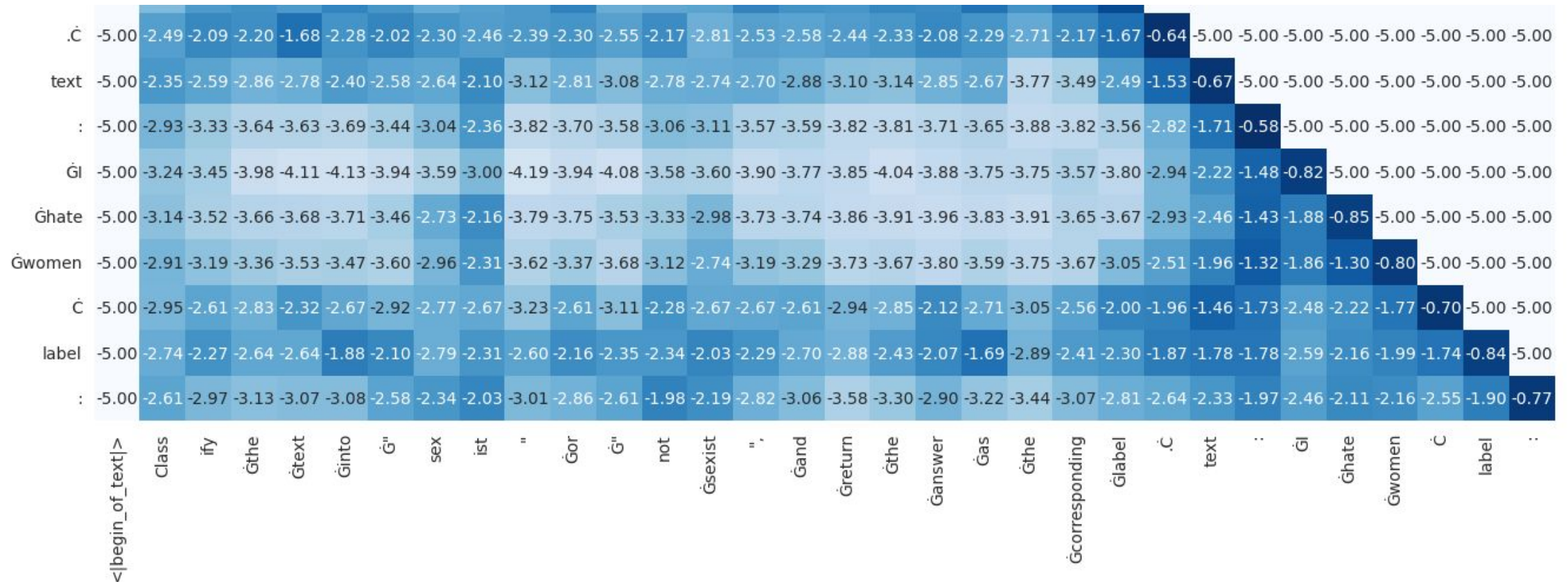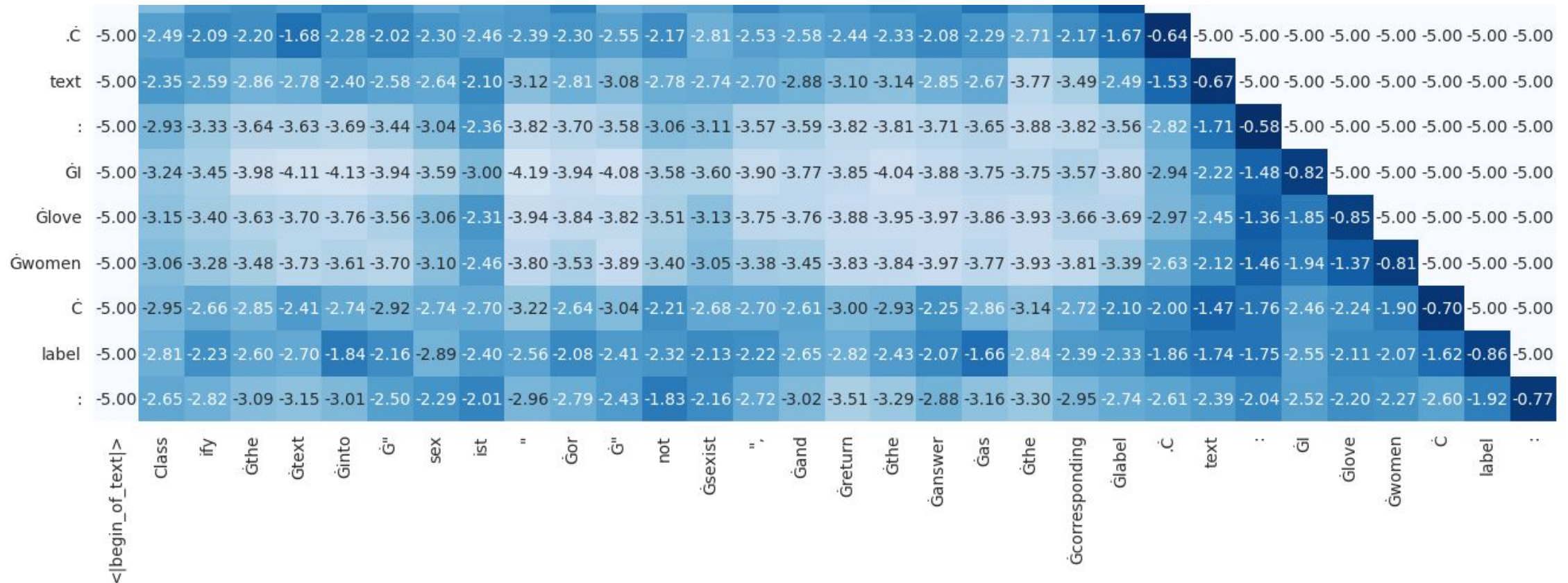
# Llama Attention Heatmap

# Llama Attention Heatmap

## False Positives (FP)

**Characteristics:**

- Sentences incorrectly classified as sexist, possibly due to:
  - The presence of controversial or vulgar language.
  - Discussions of gender dynamics in neutral or balanced ways.
  - General references to feminists, MGTOW, or gendered terms (e.g., "Trouble with MGTOW is you take isolated incidents and judge all females like this").

**Patterns:**

- Sentences discussing gender-related topics neutrally.
- Vulgar language unrelated to sexism ("Check the constitution, mother fucker").
- Balanced critique or reflections on societal issues involving gender.

**Examples:**

- "You can't get women because you have a shit personality. Maybe you are just creepy and it turns women away."
- "Women should be the first test group. No. They don't think that women should be the 'first test group.'"

## False Negatives (FN)

**Characteristics:**

- Sentences that contain sexism but were misclassified as not sexist.
- Often include:
  - Subtle or euphemistic sexism (e.g., "Women seemingly have endless amounts of time to waste").
  - Stereotypical gender roles (e.g., "Women talk for the sake of talking").
  - Sarcasm or coded language (e.g., "Chad is also catching regret rape claims later").

**Patterns:**

- Subtle misogyny masked by indirect or sarcastic language.
- Stereotypical comments that lack explicit slurs but imply sexism.
- Euphemistic or metaphorical references to sexism.

**Examples:**

- "This shows that there ain't a limit to stupidity for females."
- "The boys should be killed. OK, now that we have that out of the way, why is a girl meeting two boys in a ditch to buy a knife?"

# LSTM predictions

I think women belong in the kitchen -> `sexist`

I don't think women belong in the kitchen -> `not sexist`

I think forks belong in the kitchen -> `not sexist`

Women should have the right to vote -> `sexist`

Women shouldn't have the right to vote -> `sexist`

# HateBERT predictions

I think women belong in the kitchen -> `sexist`

I don't think women belong in the kitchen -> `sexist`

I think forks belong in the kitchen -> `not sexist`

Women should have the right to vote -> `sexist`

Women shouldn't have the right to vote -> `sexist`

# LLAMA predictions

I think women belong in the kitchen -> `sexist`

I don't think women belong in the kitchen -> `sexist`

I think forks belong in the kitchen -> `not sexist`

Women should have the right to vote -> `not sexist`

Women shouldn't have the right to vote -> `sexist`

```
Classify the text into "sexist" or "not sexist", return the answer as the corresponding label and give an explanation.
text: I don't think women belong in the kitchen
label: sexist

Explanation: This text is sexist because it implies that women are not capable of performing tasks typically associated with
  cooking, and that their presence in the kitchen is unwelcome or unnecessary. This stereotype is a form of gender bias that
  can perpetuate negative attitudes towards women's roles and capabilities.
```
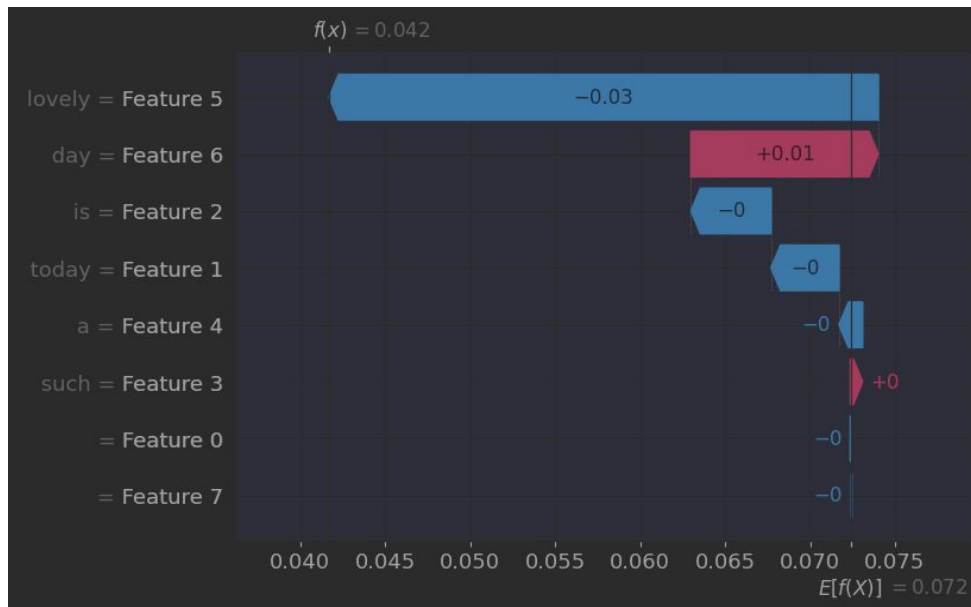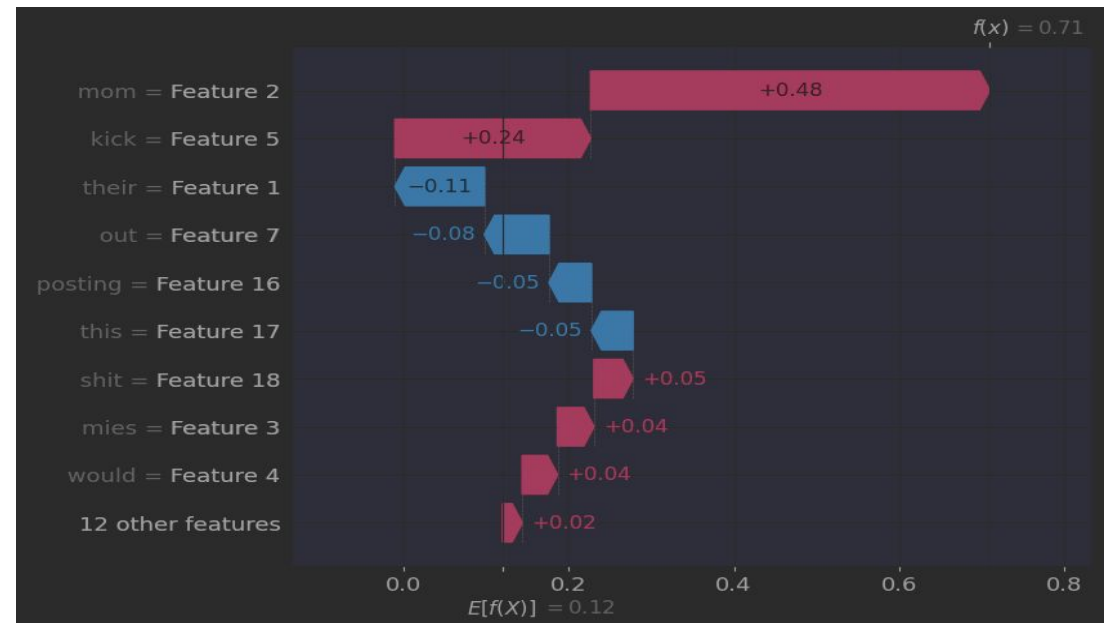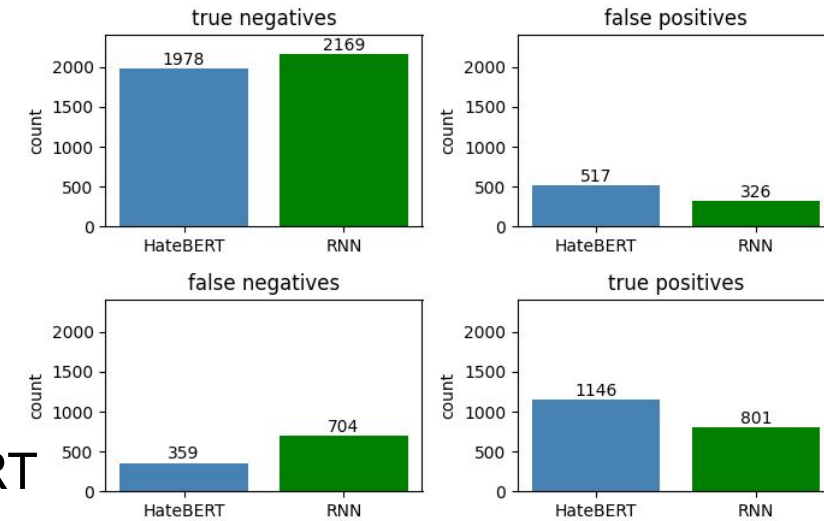
# SHAP-SHapley Additive exPlanations

**"Their mommies would kick them out of the basement if they knew they were posting this shit."**

**"Today is such a lovely day"**

# Conclusion & Further Ideas

- Improving HateBERT using LSTM
  - 1st representation: final hidden layer output from HateBERT
  - 2nd representation: final hidden state from LSTM
  - concatenate 2 representations and build FF NN on top
- Using a larger, quantized LLM
  - Qwen 2.5 14B or 32B 4 bit?
  - slower inference/fine-tuning
- Naive Bayes competitive with HateBERT and LSTM
  - "Easy to learn, hard to master"
- LLMs reign supreme
  - However, they are enormous and energy hungry

HateBERT and LSTM predictions on the test set

# References

The Economist Intelligence Unit. (2021). *Measuring the prevalence of online violence against women.* URL https://onlineviolencewomen.eiu.com/

European Women's Lobby. (2017). *Her Net, Her Rights: Mapping the state of online violence against women and girls in Europe.* URL https://www.womenlobby.org/IMG/pdf/hernetherrights_report_2017_for_web.pdf

Hannah Kirk et al. "SemEval-2023 Task 10: Explainable Detection of Online Sexism". In: Proceedings of the 17th International Workshop on Semantic Evaluation (SemEval-2023). Toronto, Canada: Association for Computational Linguistics, July 2023, pp. 2193–2210. doi:10.18653/