# ModGuard App

**Team 6**

Ali Alper Sakar
Evren Can

AIE Lab - Week 3

1

# Content

- Task
- Motivation
- Workflow Chart
- Metrics
- Swagger UI
- POST / GET
- Youtube Toxic Comments
- Grafana Graphs
- Toxic Comment Model
- Prompts for LLM
- Some Cases LLM outperforms
- Scores for LLM1, LLM2, BERT
- Zero-shot & Few-shot prompts
- Evaluation
- Improvements

# Task

- Creating a API for a content moderation system
  - uses an LLM and BERT
  - **Goals**:
    - to classify text on whether it includes undesirable language
    - providing some metrics endpoints with relevant statistics via API
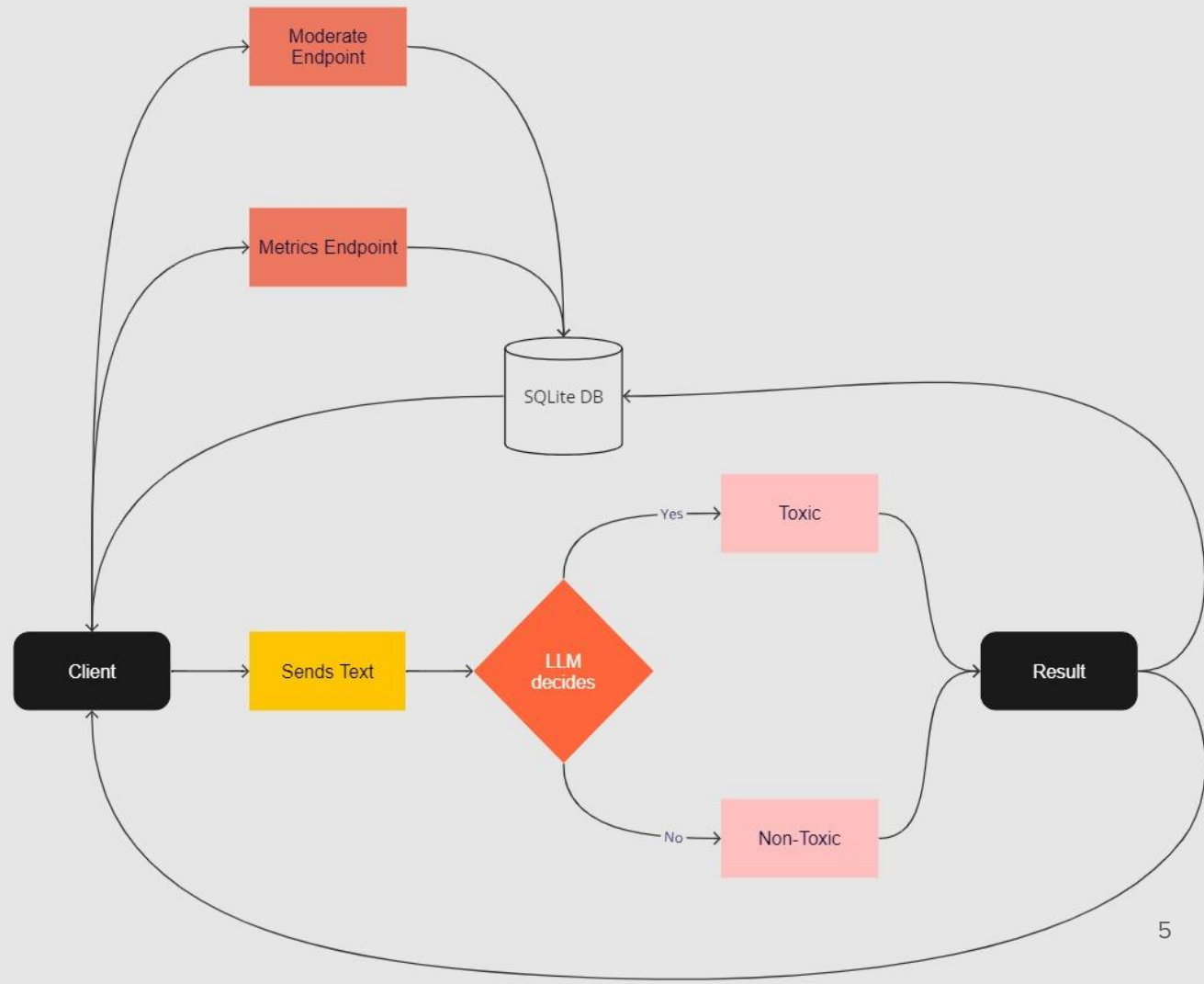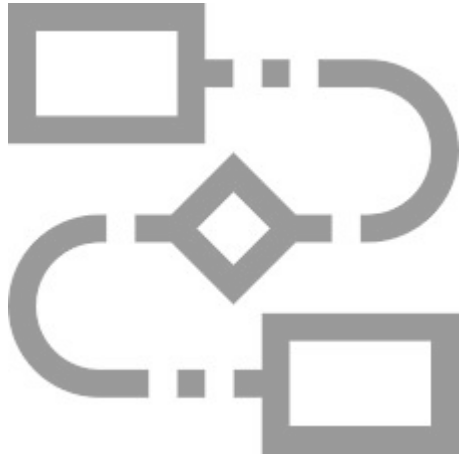    - Including documentation for API

# Motivation

- Promote Safe and Respectful Communities
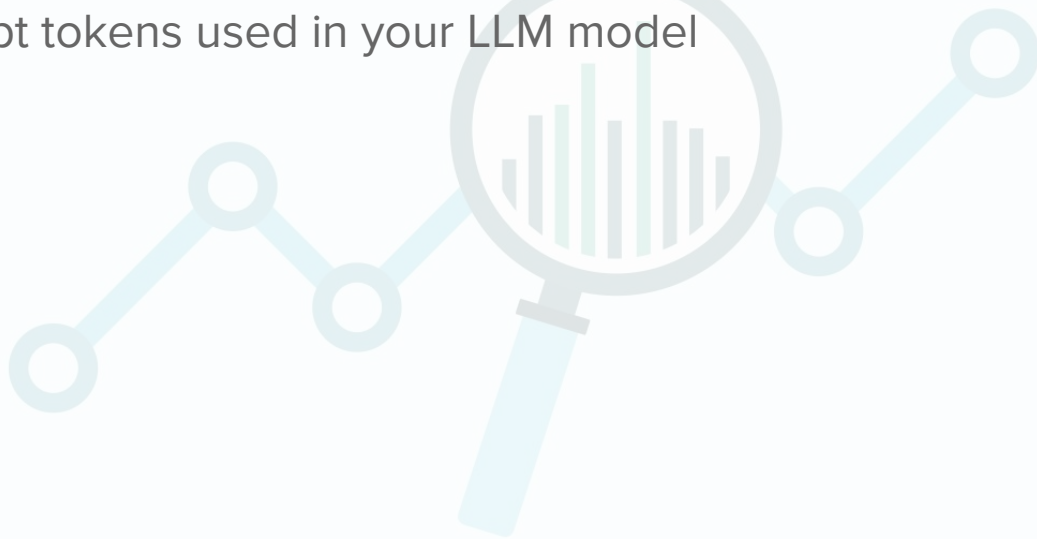- Quality Content Assurance
- User Empowerment

# Workflow Chart

# Metrics

- The sum of response times for all requests
- The count of requests
- The count of error rates
- The total number of requests received
- The total number of prompt tokens used in your LLM model

Swagger UI

# Modguard API 1.0
/static/swagger.json

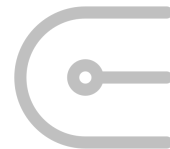## default

POST /moderate Moderate text

GET /metrics Get metrics

# POST/ moderate endpoint

# GET/ metrics endpoint

# Youtube Toxic Comments

- 1000 comments
  - 538 non-toxic
  - 462 toxic
- 11 subclassification of toxicity
  - non is considered
- Some comment examples;

```
['How did they get the Time off from work?  JK',
 'Those dumb fucks, screwing fellow citizens over instead of protesting against the police. 🧑\u200dd',
 'I would've pressed my m button and shoot the tsar bombar outta my car into them',
 "Everyone wants to be the next MLK but the issue is that there aren't problems in the US anymore that require another MLK.",
 "It's hilarious seeing these morons get hit by cars"]

[False, True, True, False, True]
```

Non-Toxic
46.2%

Toxic
53.8%

https://www.kaggle.com/datasets/reihanenamdari/youtube-toxicity-data

10

# GRAFANA Graphs

# Toxic Comment Model

- Fine-tuned version of **DistilBERT** to classify toxic comments
- Tensor size of 512
  - Index of 174 was too long, so shorten is applied
- DistilBERT
  - Smaller       -       40% less parameters
  - Runs 60% faster
  - Cheaper
  - Lighter
  - Preserving over 95% of BERT's performances
- Poor performance for the muslim and jewish group
  - 'Muslims are people who follow or practice Islam, an Abrahamic monotheistic religion.' returns toxic

# Prompts for LLM

```
prompt1 = "You are tasked with the role of a content moderation system, \
responsible for analyzing user-submitted texts to identify the presence \
of undesirable or inappropriate language. Your task is to review each text \
and categorize it accordingly: if the text contains any form of toxic language, \
please respond with 'toxic'. If the text is free of inappropriate content, \
respond with 'non-toxic'. Answer only with lowercase and without punctuation."


prompt2 = "You are tasked with the role of a content moderation system, \
responsible for analyzing user-submitted texts to identify the presence \
of undesirable or inappropriate language. Your task is to review each text \
and categorize it accordingly: if the text contains any form of toxic language, \
please respond with 'true'. If the text is free of inappropriate content, \
respond with 'false'. Answer only with lowercase and without punctuation."
```

# Some Cases LLM outperforms

- <u>Sentence 1</u>
  - "This interviewer is right. This guy should show some respect to the paramilitary police that are killing and getting away with it. Ha!"
    - **llm1** -> non toxic
    - **llm2** -> non toxic
    - **bert** -> toxic
    - **ground_truth** -> non toxic
- <u>Sentence 2</u>
  - "I don't think women are suited for leadership positions."
    - **llm1** -> toxic
    - **llm2** -> toxic
    - **bert** -> non toxic
    - **ground_truth** -> toxic
- But doesn't mean that it performs always better!

# Scores of LLM1, LLM2, and BERT

# Zero-shot prompt

```
innap = "You are a worthless piece of trash and nobody likes you."
poten_inapp = "I don't think women are suited for leadership positions."
app = "I think you did a great job on this project, well done!"
ambiguous = "I don't mean to be rude, but that's a terrible idea."
```

|  | **inapp** | **poten_inapp** | **app** | **ambiguous** |
|---|---|---|---|---|
| **LLM1** | Toxic | Toxic | non-toxic | non-toxic |
| **LLM1** | Toxic | Toxic | non-toxic | non-toxic |
| **BERT** | Toxic | non-toxic | non-toxic | non-toxic |

# Few-shot prompt

```
innap = "You are a worthless piece of trash and nobody likes you."
poten_inapp = "I don't think women are suited for leadership positions."
app = "I think you did a great job on this project, well done!"
ambiguous = "I don't mean to be rude, but that's a terrible idea."
```

```
system_content =  "'You are worthless.' - Result: Toxic, 'I love the way you think.' - Result: Not toxic, 'TEXT_TO_BE_ANALYZED' Result:?"
```

|  | **inapp** | **poten_inapp** | **app** | **ambiguous** |
|---|---|---|---|---|
| **LLM1** | Toxic | Toxic | non-toxic | non-toxic |
| **LLM1** | Toxic | Toxic | non-toxic | non-toxic |
| **BERT** | Toxic | Toxic | Toxic | Toxic |

# Evaluation

- Overview evaluation between LLM and BERT
  - Similar performance
- Individual Model Analysis
  - **LLM**
    - **Strength**: better understanding delicate language
    - **Weakness**: oversensitivity to certain phrases
  - **BERT**
    - **Strength**: efficient handling of explicit toxic phrases
      - (identifying sarcasm and reference to violence)
    - **Weakness**: potential inability to detect underlying biases

# Areas for Improvements

- Feedback Loop
  - Establishing a feedback loop for continuous improvements
- Custom Training
  - Exploring opportunities for custom training of the models for better performance

# Thank you !

Ali Alper Sakar

Evren Can