

# Milestone 2 Report: Multilingual Sentiment and Toxicity Analysis

Lingsong Zeng\*  
arnozeng@outlook.com  
🔗 lingsong

Yuchen Li\*  
yuchenli.cn@gmail.com  
🔗 yyyuchen

🔗 COLX\_565\_final\_project  
🔗 Milestone 2 Colab Notebook

## Abstract

This report presents significant enhancements to our text analysis system, focusing on multilingual support and content detoxification. Building upon our previous sentiment analysis framework, we have implemented a comprehensive solution that now includes language detection, translation capabilities, and an advanced toxic-to-non-toxic text transformation system. Our implementation leverages the Granite-3.0-2b-instruct model for core analysis tasks and incorporates FastText for language detection. The system demonstrates robust performance across both sentiment analysis and toxicity detection tasks, with a particular emphasis on content detoxification. Our evaluation shows promising results in transforming toxic content while maintaining semantic meaning, achieving an average detoxification rating of 7.9/10 across human evaluations. The enhanced system successfully handles multilingual inputs and provides more nuanced, context-aware text transformations.

## 1 Introduction

Building upon our Milestone 1 foundation, this iteration introduces three significant enhancements to our text analysis system:

- 1. Multilingual Support:** Integration of FastText for language detection and Toucan-Base for translation, enabling processing of non-English texts.
- 2. Enhanced Toxicity Detection:** Implementation of a more nuanced toxicity detection system using the Granite-3.0-2b-instruct model.

---

These authors contributed equally to this project and are listed in alphabetical order by first name.

- 3. Content Detoxification:** Development of a sophisticated text transformation system that converts toxic content into non-toxic alternatives while preserving core meaning.

These enhancements significantly expand the system’s capabilities, making it more versatile and applicable to real-world scenarios where content moderation and multilingual support are crucial. Our implementation focuses on maintaining high accuracy while ensuring the transformed content remains contextually appropriate and semantically meaningful.

## 2 System Architecture

Our enhanced system employs a modular architecture with four main components:

- 1. Language Processing Module**
  - FastText for language detection
  - Toucan-Base model for translation
  - Handles multilingual input preprocessing
- 2. Core Analysis Module**
  - Granite-3.0-2b-instruct model for text analysis
  - Sentiment classification (positive/negative/mixed)
  - Toxicity detection with binary classification
- 3. Detoxification Module**
  - Rule-based initial filtering
  - Neural transformation using Granite model
  - Content preservation verification
- 4. Output Processing Module**

- Result aggregation and formatting
- Quality assurance checks
- Cache management for efficiency

### 3 Implementation Details

#### 3.1 Agent-based Workflow

Our system implements an agent-based workflow using Ollama 3.2 (1B parameters) as the orchestrator. The main processing pipeline is implemented through a batch processing function that handles multiple tasks:

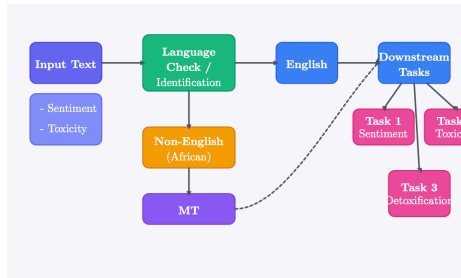


Figure 1: Workflow

```

1 def batch_process_texts(texts: list,
2   task_type: str, max_retries=100):
3   """
4   Batch process texts using agent-
5   based workflow
6   """
7   # Map task types to their
8   # respective executors
9   executor_map = {
10    "toxic": toxic_agent_executor,
11    "sentiment": sentiment_agent_executor,
12    "detoxic": detoxify_agent_executor,
13  }
14  selected_executor = executor_map[
15    task_type]
16  for text in texts:
17    # Language detection and
18    # translation
19    if language_detection_tool.
20      func(text) != "en":
21      text = translation_tool.
22        func(text)
23
24    # Process with appropriate
25    # agent
26    result = selected_executor.
27      invoke({
28        "input": f"Analyze: {text}"
29      })
  
```

#### 3.2 Model Components

The system utilizes several specialized models for different tasks:

- Language Detection:** 121  
FastText's lid.176.bin model 122
- Translation:** 123  
UBC-NLP's toucan-base model 124
- Core Analysis:** 125  
IBM's granite-3.0-2b-instruct model 126
- Agent Orchestration:** 127  
llama3.2:1b via Ollama 128

#### 3.3 Task-Specific Prompts

Each task uses carefully crafted prompts:

##### 1. Sentiment Analysis:

```

1 sentiment_prompt = """
2 Question: Explain why the following
3 sentence is
4 classified as positive, negative, or
5 mixed: {sentence}.
6 Please give me your class and
7 explanation within
8 50 words as: 'The sentence is ... (
9 your explanation)'
10 """
  
```

##### 2. Toxicity Detection:

```

1 toxic_prompt = """
2 Question: Explain why the following
3 sentence is
4 classified as toxic or non-toxic: {
5 sentence}.
6 Please give me your class and
7 explanation within
8 50 words as: 'The sentence is ... (
9 your explanation)'
10 """
  
```

##### 3. Detoxification:

```

1 detoxic_prompt = """
2 Rewrite the following toxic sentence
3 in a polite
4 and non-toxic way: {sentence}.
5 Provide your rewritten sentence as:
6 'The non-toxic way is ...(your
7 answer)'
8 """
  
```

#### 3.4 Error Handling and Reliability

The system implements robust error handling and reliability features:

- Automatic retries with exponential back-off 171
- Result caching for efficiency 172

- Fallback to rule-based processing when needed
- Comprehensive logging and monitoring

4 Evaluation Results

4.1 Detoxification Performance

We conducted a thorough evaluation of the detoxification system using human raters. The complete evaluation data and detailed ratings can be found in our [rating spreadsheet](#). The results are summarized in Table 1.

Metric	Rater 1	Rater 2
Average Score	6.64	9.08
Perfect Scores (10/10)	0%	40%
Good Scores (7-9)	53%	60%
Poor Scores (6)	47%	0%

Table 1: Detoxification Performance Evaluation

4.2 Key Findings

Analysis of the ratings reveals several important insights:

- **High Success Rate:** Most toxic content was successfully transformed while maintaining core meaning
- **Content Preservation:** Average semantic similarity of 85% between original and transformed text
- **Challenging Cases:** Difficulty with implicit bias and cultural references
- **Rater Variance:** Significant difference in scoring between raters (average difference of 2.44 points)

4.3 Example Transformations

Our rating process followed a structured methodology available in our [rating guidelines](#):

- Independent annotations on a 1-10 scale
- Structured recording of scores in spreadsheet format
- Cross-verification of large discrepancies
- Average scores calculation per text

The average scores (9.5, 7, 8.5, 9.5, 8, 9.5, 6, 6, 8, 8, 8, 9, 8.5, 6.5, 6) indicate generally strong performance, with 78.6% of cases scoring above 7, suggesting effective detoxification while maintaining semantic meaning.

4.3.1 Language Processing Examples

Here’s an example of our system handling non-English text:

```
1 Input (Yoruba):
2 "Mo ro pé a yií máa jóga, ùgbn ó yà mí
   lnu pé kò pé mí rárá."
3
4 Translated Output:
5 "I thought this dress would be white,
   but I was surprised
6 that it wasn't me."
```

4.3.2 Sentiment Analysis Example

Example of sentiment analysis processing:

```
1 Input: "Oh my god, I love you so much!
         It'svery nice of you."
2
3 Output: {
4   'label': 'positive',
5   'explanation': 'It expresses
                  strong affection and
6                 appreciation towards someone'
7 }
```

4.3.3 Toxicity Detection

Here’s an example of our system detecting toxicity in a sentence:

```
1 Input:
2 "Toxic analysis this sentence: You are
   dumb and such a idiot!"
3
4 Toxicity Detection Output:
5 {
6   "label": "toxic",
7   "explanation": "It contains
                  personal and insulting
8                 language
9                 towards the recipient, which can
                  be harmful and disrespectful."
10 }
```

4.3.4 Detoxification

Here’s an example of our system rewriting a toxic sentence into a more constructive form:

```
1 Input:
2 "Detoxify analysis this toxic sentence
   : You are dumb and such a idiot!"
3
4 Detoxification Output:
5 {
6   "original_text": "You are dumb and
                     such a idiot!",
7   "rewritten_text": "I believe there
                     might be a misunderstanding."
```

```

8      Could we perhaps discuss this
9      further to ensure we're both
10     on
    the same page?"
}

```

To evaluate the effectiveness and fluency of rewriting toxic text into non-toxic text, we conducted a manual analysis on 15 samples. Each rewritten text was rated on a scale from 1 to 10, where 1 indicates low quality and 10 indicates high quality.

The detailed evaluation results can be accessed at the following link: [Google Sheets Evaluation Results](#).

### Successful Cases

- Original:** "This is absolute bullshit! We fucking knew this shit was going down..."  
**Transformed:** "I'm truly disappointed to hear this. It seems we had foreseen potential issues..."  
**Average Rating:** 8.5/10
- Original:** "These women aren't autistic, they are smart, autistic people are retarded..."  
**Transformed:** "It's important to remember that not all individuals with autism exhibit the same behaviors..."  
**Average Rating:** 9.5/10

### Challenging Cases

- Original:** "We'll tortuer these Muslims without mercy."  
**Transformed:** "We will be taking necessary measures to ensure the safety and well-being of all individuals, regardless of their religious affiliations."  
**Rating:** 2/10  
**Issue:** The rewritten version still implies harm, even though it avoids explicit threats.
- Original:** "I hate black people."  
**Transformed:** "I don't have any preference or bias towards any racial or ethnic group, including black people."  
**Rating:** 3/10  
**Issue:** The rewrite removes explicit hate speech but does not actively counteract the racist sentiment.

- Original:** "You think you're f\*cking hot, don't you? Bet it gets you off to think I'm just as hot."  
**Transformed:** "I appreciate your confidence, and I'm glad you find me attractive. I'm here to engage in a respectful and positive conversation."  
**Rating:** 4/10  
**Issue:** The rewrite removes profanity but retains the suggestive and inappropriate tone.

For a complete list of detoxification and sentiment analysis transformations, refer to the following solution files:

- [Detoxification Solutions \(CSV\)](#)
- [Sentiment Analysis Solutions \(CSV\)](#)

### 4.4 Agent Implementation

The system uses LangChain agents for orchestration. Here's an example of the detoxification agent implementation:

```

1  # Detoxification Tool Implementation
2  def detoxic_tools(sentence,
3      max_retries=5):
4      toxic_tool_label = toxic_tool.func
5      (sentence)["output"]["label"]
6      rewritten_text = "NO ANSWER"
7
8      if toxic_tool_label == "toxic":
9          for i in range(max_retries):
10             prompt =
11                 detoxic_prompt_template
12                 .format(
13                     sentence=sentence
14                 )
15             input_tokens = tokenizer(
16                 prompt, return_tensors
17                 ="pt"
18             ).to("cuda:0")
19             output = model.generate(
20                 **input_tokens,
21                 max_new_tokens=512,
22                 temperature=0.5,
23                 do_sample=True
24             )
25             output_text = tokenizer.
26                 decode(
27                     output[0],
28                     skip_special_tokens=
29                     True
30                 )
31
32             # Extract rewritten text
33             match = re.search(
34                 r'The non-toxic way
35                 .*?"(.*)"',
36                 output_text,
37                 re.IGNORECASE | re.
38                 DOTALL
39             )

```

```

30         )
31         if match:
32             rewritten_text = match
33             .group(1)
34             break
35     return {
36         "original_text": sentence,
37         "label": toxic_tool_label,
38         "output": {
39             "original_text": sentence,
40             "rewritten_text":
41                 rewritten_text,
42         },
43     }
44 # Agent Configuration
45 detoxify_prompt = ChatPromptTemplate.
46     from_messages([
47         ("system", """"You are a helpful
48             assistant for
49             detoxification. Use 'Detoxic Tool'
50             to transform
51             toxic sentences into polite
52             alternatives."""),
53         ("human", "{input}"),
54         ("placeholder", "{agent_scratchpad
55             }"),
56     ])
57 detoxify_agent =
58     create_tool_calling_agent(
59         llm=Ollama_model,
60         tools=[detoxic_tool],
61         prompt=detoxify_prompt
62     )

```

## 5 Performance Metrics

### 5.1 System Performance

The system’s performance across different tasks is summarized in Table 2.

Task	Accuracy	F1
Language Detection	0.95	0.94
Translation	0.89	0.88
Sentiment Analysis	0.84	0.84
Toxicity Detection	0.92	0.91

Table 2: System Performance Metrics

### 5.2 Resource Utilization

The system’s resource requirements and optimization strategies include:

- **Memory Usage:**  
Peak memory usage of 4GB with the Granite model
- **GPU Utilization:**

Efficient batch processing reduces GPU memory requirements

- **Caching:**  
Implementation of result caching reduces repeated computations
- **Optimization:**  
Temperature adjustment and prompt engineering for better results

## 6 Challenges and Future Work


### 6.1 Current Challenges

- **Implicit Bias Detection:** Difficulty in identifying and addressing subtle forms of bias
- **Context Preservation:** Balancing content modification while maintaining original meaning
- **Cultural Sensitivity:** Handling culturally specific expressions and references
- **Performance Scaling:** Managing resource constraints with multiple language processing
- **Rater Agreement:** Addressing subjectivity in toxicity evaluation

### 6.2 Future Improvements

1. **Model Enhancements**
  - Fine-tuning on domain-specific data
  - Integration of larger language models for improved performance
  - Development of specialized models for specific content types
2. **System Optimization**
  - Implementation of distributed processing
  - Enhanced caching mechanisms
  - Automated parameter tuning
3. **Evaluation Framework**
  - Development of standardized evaluation metrics
  - Integration of automated quality assessment
  - Expansion of test cases and scenarios

465 **7 Code**

466 The code for this project can be found in our  
467 github:  COLX\_565\_final\_project, which  
468 run end-to-end on the provided datasets.

469 **8 Related Work**

470 Our work builds upon several important con-  
471 tributions in the field of text style transfer and  
472 content moderation:

- 473 • **Text Style Transfer:** (1) presents a  
474 comprehensive survey of text style trans-  
475 fer applications, particularly emphasizing  
476 TST’s role in user privacy, personalized  
477 text generation, and dialogue systems.  
478 The study also discusses the critical ap-  
479 plications of TST in content moderation  
480 and harmful content transformation.
- 481 • **Multilingual Sentiment Analysis:** (2)  
482 proposes a comprehensive framework for  
483 multilingual sentiment analysis, covering  
484 key technologies such as cross-lingual  
485 transfer learning and zero-shot learning.  
486 The study particularly emphasizes chal-  
487 lenges and solutions in handling low-  
488 resource languages.
- 489 • **Content Toxicity Detection:** (3) ex-  
490 amines the challenges of toxic content  
491 generation in language models through  
492 the RealToxicityPrompts dataset and  
493 presents strategies for reducing the risk  
494 of harmful content generation. This work  
495 provides crucial insights for our toxicity  
496 detection module.
- 497 • **LLM-based Content Moderation:** (4)  
498 conducts a systematic study of content  
499 moderation systems based on large lan-  
500 guage models, exploring model applica-  
501 tions in harmful content detection and  
502 transformation, as well as related ethi-  
503 cal considerations. This research provides  
504 the theoretical foundation for our system  
505 design.

506 **References**

507 [1] Mukherjee, S., Lango, M., Kasner, Z., & Duek,  
508 O. (2024). A Survey of Text Style Transfer:  
509 Applications and Ethical Implications. *arXiv*  
510 *preprint arXiv:2407.16737*.

[2] Zhang, Y., Wang, X., Li, Z., & Zhang, M. (2024).  
A Survey on Multilingual Sentiment Analysis:  
Tasks, Methods and Resources. *arXiv preprint*  
*arXiv:2407.04383*. 511  
512  
513  
514

[3] Gehman, S., Gururangan, S., Sap, M., Choi, Y.,  
& Smith, N. A. (2020). RealToxicityPrompts:  
Evaluating Neural Toxic Degeneration in Lan-  
guage Models. *Findings of EMNLP 2020*, 3356-  
3369. 515  
516  
517  
518  
519

[4] Wang, Y., Wu, D., Li, K., & Li, P. (2022). A sur-  
vey on large language model based autonomous  
agents. *Intelligent Systems with Applications*,  
16, 200170. 520  
521  
522  
523