

---

# Adapting Safe-for-Work Classifier for Malaysian Language Text: Enhancing Alignment in LLM-Ops Framework

---

Aisyah Razak\*

Ariff Nazhan†

## Abstract

## 1 Introduction

## 2 Data Source

The data for this study was collected from various platforms, including social media, public forums, and publicly available datasets. The majority of the data is in the Malay language and relevant to the Malay context. By utilizing these comprehensive datasets from multiple sources, we have strengthened the robustness and accuracy of our classification model, enabling it to effectively tackle the challenges of identifying self-harm and sexism in online content.

### 2.1 Social Media

Data was collected from popular social media platforms such as Twitter and Facebook. Two main approaches were employed to gather relevant data:

1. Keyword-based scraping: a list of keywords associated with explicit content was compiled. These keywords were used to extract tweets from the platform.
2. Profile-based scraping: a list of profiles known for regularly posting NSFW content was curated. Posts from these profiles were then scraped to obtain a more targeted dataset.

The combination of these two scraping methods resulted in a comprehensive and diverse dataset from social media, capturing both keyword-specific content and data from profiles that frequently share explicit material.

### 2.2 Public Articles

For public articles, we have collected data from various articles and blogs which are [b.cari](#), which hosts a wide range of user-generated content in Malay.

The collected dataset consists of human dialogues extracted from these articles and blog posts. It notably includes some dialogues that contain nsfw content. The inclusion of such dialogues, while potentially controversial, is important to allow the trained classifier to effectively detect explicit content that may realistically occur in open-ended dialogue systems.

---

\*aisyahrazak171@gmail.com

†ariffnazhn@gmail.com

## 2.3 Public Datasets

We have also collected data from publicly available dataset on kaggle such as [Kaggle: Suicide and Depression Detection](#) datasets. The dataset is a collection of posts from the "SuicideWatch" and "depression" subreddits of the Reddit platform. This dataset contains a wide range of suicide ideation contexts, providing valuable insights for our research.

We also leverage [Explainable Detection of Online Sexism \(EDOS\)](#) from github and [EXIST: sEXism Identification in Social neTwork \(EXIST\)](#) from web that contain a diverse collection of sexism statements, which have significantly contributed to the success of our classifier in identifying and categorizing such content.

## 3 Methodology

### 3.1 Overall

### 3.2 Knowledge Distillation

### 3.3 Centroid Filtering

### 3.4 Pseudolabeling

## 4 Result

### 4.1 Evaluation

## 5 Acknowledgement

Special thanks to Malaysia-AI volunteers especially [Ammar Azman](#), [M. Amzar](#), [Muhammad Farhan](#), [Syafie Nizam](#), [Alif Aiman](#), [Azwan Zuharimi](#) and [Haziq Zikry](#) for contributing dataset to train Malaysian Reranker models.

We would like to express our gratitude to NVIDIA Inception for generously providing us with the opportunity to train our model on the Azure cloud. Their support has played a crucial role in the success of our research, enabling us to leverage advanced technologies and computational resources.

We extend our thanks to the wider research community for their valuable insights and collaborative discussions, which have greatly influenced our work. This paper reflects the collective efforts and contributions from both NVIDIA Inception and the broader research community.

## 6 Conclusion

## References