

# Polarization Detection at POLAR @ SemEval-2026: A MARBERT-based Approach for Arabic Polarization Detection

*Milestone Report*

Guanghui Ma

guanghui.ma@sabanciuniv.edu

Ali Khaled A. Ishtay Altamimi

ali.altamimi@sabanciuniv.edu

Nuh Al Sharafi

nuh.sharafi@sabanciuniv.edu

Hussein MH Nasser

h.nasser@sabanciuniv.edu

## Abstract

This milestone report presents our current progress on developing a system for POLAR @ SemEval-2026 Task 9: subtask 1, focusing on polarization detection in Arabic social media text. We employ MARBERT, a pretrained Arabic language model, fine-tuned on the official challenge dataset supplemented with translated English data. Our approach addresses the unique challenges of Arabic text processing, including dialectal variation, morphological complexity, and code-switching. We are developing two model variants: Smol-MARBERT trained on Arabic data only, and BigData-BERT leveraging multilingual translated datasets. This report describes our current methodology, preprocessing pipeline, initial results, and planned optimizations for improving polarization detection across diverse Arabic dialects.

## 1 Introduction

Polarization refers to the divergence of political attitudes away from the center and towards more extreme views to the far-left or far-right. The internet and communication through it has allowed polarization to spread in ways previously unseen, the veil of anonymity encouraging more extreme behaviors and thoughts. This poses a threat to social cohesion as well as the proliferation of harmful thoughts.

This project's purpose is to detect online polarization as part of POLAR @ SemEval-2026 (Task 9): Detecting Multilingual, Multicultural, and Multievent Online Polarization. The task aims to develop models capable of identifying and classifying polarized content in text form. This project focuses on the Arabic language as part of Subtask 1: Polarization Detection. In this subtask, the goal is to determine whether a given text exhibits polarized opinion or not. Texts are classified into one of two categories: polarized or non-polarized.

**Current Progress:** Up to this stage, our work has included an in-depth literature review covering related papers in Arabic natural language processing and polarization detection. In parallel, we have carried out preliminary data exploration and analysis on the dataset provided by the POLAR @ SemEval-2026 challenge organizers, examining aspects such as label balance, text length distribution, and dialectal variation. We have also conducted preprocessing and exploratory data analysis on the datasets used for training as well as beginning to train and evaluate prospective models. This milestone report documents our progress to date and outlines our planned next steps toward the final submission.

## 2 Dataset Selection

For this project, we use the official Arabic dataset provided by the organizers of the POLAR @ SemEval-2026 (Task 9) challenge. This dataset has been specifically curated for the Polarization Detection subtask and contains social media texts labeled as either polarized or non-polarized. Each instance is annotated

based on whether it expresses polarized attitudes, opinions, or language within its contextual setting, thereby removing the need for any additional manual labeling.

This dataset is particularly well suited to our objectives for several reasons. First, it was designed explicitly for polarization detection, ensuring that all included samples are directly relevant to our research goals. Second, the Arabic subset covers a broad range of dialects and regional varieties, offering rich linguistic and cultural diversity that reflects real-world online discourse. This diversity enhances our model’s ability to generalize across different forms of Arabic expression. Finally, since all participating teams in the challenge utilize the same dataset, our results can be directly benchmarked against other approaches, allowing for fair comparison and meaningful evaluation of model performance across methodologies.

In addition, in order to have a wider range of data, we used the official English dataset provided by the organizers, and translated it into a variety of Arabic dialects using OpenL translation API to serve as supplementary training data. This translated corpus expands our training pool, supporting the development of models with improved robustness and cross-linguistic understanding.

## 3 Approach

Our approach for this project is guided by established conventions in Arabic Natural Language Processing and insights from polarization-related research. The main objective is to develop models capable of accurately classifying Arabic social media texts as either polarized or non-polarized while maintaining robust generalization across different dialects and contexts. This section outlines our planned methodology and the work completed to date.

### 3.1 Literature Foundations and Methodology

Polarization detection goes beyond standard sentiment or stance classification, as it focuses on identifying division, hostility, and opinion extremity rather than mere agreement or disagreement POLAR @ SemEval (2025). Previous research highlights that polarization is highly context-dependent and non-linear, necessitating deep, high-capacity models that can capture subtle linguistic, affective, and discourse cues Vasist et al. (2023). Early approaches relied on feature-engineering-based baselines, including lexical affect scoring and discourse structure modeling. While these methods offered interpretability, Transformer-based architectures have consistently outperformed them in capturing nuanced contextual relationships Wolf et al. (2020).

### 3.2 Selected Computational Framework

Based on comparative studies in Arabic stance detection and fake news classification, MARBERT has been identified as the most suitable pretrained language model for this task AlShenaifi et al. (2024); Al Hariri and Abu Farha (2024). Fine-tuning MARBERT is a practical choice given the limited computational resources available. It was trained on a large mixture of Modern Standard Arabic (MSA) and dialectal user-generated content, making it well suited to the noisy and informal language of online discussions. This alignment between training data and target domain reduces the need for extensive retraining while supporting the context-heavy analysis required for Arabic.

The first model, referred to as Smol-MARBERT, is being fine-tuned exclusively on the Arabic portion of the SemEval dataset. In contrast, a second model, BigData-BERT, is under development to leverage translated datasets from English into Arabic. Since the polarization rubric for Arabic is not publicly available, this approach expands the available training data while approximating polarization labels. Several challenges accompany this strategy:

- **Relevance filtering** is necessary since many English datasets include political content (e.g., Biden/Trump) that is irrelevant in Arabic contexts.

Dialect	Code	Est. %
Standard Arabic	ar	12.5
Egyptian Arabic	arz	20.0
North Levantine Arabic	apc	20.0
Sudanese Arabic	apd	2.5
Gulf Arabic	afb	10.0
Iraqi Arabic	acm	7.5
Hejazi Arabic	acw	10.0
Najdi Arabic	ars	2.5
Yemeni Arabic	ayn	5.0
Maghrebi Arabic	ary	10.0

Table 1: Arabic dialect distribution estimates in social media data.

- **Translation quality** must be verified, as the accuracy of machine translation varies; validation by speakers of each source language is required to ensure semantic fidelity.
- **Mixed-language comments** are common in Arabic social media, where users frequently code-switch with English or French, requiring specialized preprocessing to handle multilingual text effectively.

Finally, we will replicate the same training process using MARBERT V2. By comparing the performance and accuracy of both versions, we will determine which model is better suited for this task. The more efficient version will then be optimized and adopted for the final system.

### 3.3 Handling Arabic Linguistic Challenges

Arabic presents unique linguistic challenges due to diglossia, morphological complexity, and orthographic variation Farghaly and Shaalan (2009). To mitigate these issues, we employ character normalization to standardize inconsistent letter forms (e.g., various forms of alif) and subword tokenization by retaining SentencePiece from MARBERT’s pretrained tokenizer, which is effective for morphologically rich and agglutinative languages Qarah and Alsanoosy (2024).

### 3.4 Dialect Coverage

The model training and evaluation will explicitly account for major Arabic dialects present on social media based on distribution estimates. Details are shown in Table 1.

## 4 Next Steps

With Subtask 1 of the competition already in progress and with a trained model showing promising initial results, our next steps are aimed at moving on to Subtask 2 of the project, which is polarization type classification. Additionally, we aim to improve our model’s accuracy and performance for Subtask 1 by exploring other techniques and methods, including new approaches to data preprocessing and hyperparameter tuning.

The distribution of the planned next steps and the work we have completed so far are as follows:

**Data Exploration and Analysis** Perform exploratory data analysis (EDA) to examine label distribution, text characteristics, and class imbalance issues in the dataset, providing insights that will inform model training and evaluation. *Status: In Progress. Responsible: Ali Khaled A. Ishtay Altamimi*

**Preprocessing and Data Preparation** Implement and validate the Arabic-specific preprocessing pipeline, including character normalization, tokenization, and handling of dialectal text variations to ensure data consistency and model readiness. *Status: Partially Complete. Responsible: Guanghui Ma*

**Model Training** Fine-tune the BDB model on the official Arabic and English dataset using both weighted and contrastive loss configurations to handle class imbalance and improve classification performance. *Status: In Progress. Responsible: Nuh Al Sharafi*

**Testing and Optimization** Evaluate the trained models on the development and test sets, perform hyperparameter tuning, and apply optimization strategies to enhance the Macro F1 score and overall model robustness. *Status: Planned. Responsible: Hussein MH Nasser*

## 5 Conclusion

This milestone report has outlined our current progress in developing a MARBERT-based polarization detection system for Arabic social media text. We have established a solid foundation through literature review, dataset preparation, and initial model development. Our preliminary work demonstrates the feasibility of our approach, though significant optimization and evaluation work remains. The next phase of our project will focus on refining our models, conducting comprehensive experiments, and extending our work to Subtask 2. A final report will document our complete methodology, experimental results, and findings upon completion of the shared task.

## Acknowledgments

We thank the organizers of POLAR @ SemEval-2026 for providing the datasets and evaluation framework for this shared task.

## References

- N. AlShenaifi, N. Alangari, and H. Al-Negheimish. 2024. Rasid at StanceEval: Fine-tuning MARBERT for Arabic Stance Detection. In *Proceedings of the Second Arabic Natural Language Processing Conference*, pages 828–831, Bangkok, Thailand.
- Y. Al Hariri and I. Abu Farha. 2024. SMASH at StanceEval 2024: Prompt Engineering LLMs for Arabic Stance Detection. In *Proceedings of the Second Arabic Natural Language Processing Conference*, pages 800–806, Bangkok, Thailand.
- J. Gatto, O. Sharif, and S. M. Preum. 2023. Chain-of-Thought Embeddings for Stance Detection on Social Media. In *Findings of the Association for Computational Linguistics: EMNLP*, pages 4154–4161.
- Z. Zhang, J. Zhang, H. Xu, J. Guo, and X. Cheng. 2025. MPRF: Interpretable Stance Detection through Multi-Path Reasoning Framework. In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 454–470, Suzhou, China.
- POLAR @ SemEval. 2025. Task Description: Detecting Multilingual, Multicultural, and Multievent Online Polarization.
- H. Najadat, M. Tawalbeh, and R. Awawdeh. 2022. Fake news detection for Arabic headlines-articles news data using deep learning. *International Journal of Electrical and Computer Engineering (IJECE)*, 12(4):3951–3959.

- P. N. Vasist, D. Chatterjee, and S. Krishnan. 2023. The Polarizing Impact of Political Disinformation and Hate Speech: A Cross-country Configural Narrative. *Information Systems Frontiers*, pages 1–26.
- T. Wolf, L. Debut, V. Sanh, and J. Chaumond. 2020. Transformers: State-of-the-Art Natural Language Processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations (EMNLP Demos)*, pages 38–45.
- F. Qarah and T. Alsanoosy. 2024. A Comprehensive Analysis of Various Tokenizers for Arabic Large Language Models. *Applied Sciences*, 14(13):5696.
- A. Farghaly and K. Shaalan. 2009. Arabic Natural Language Processing: Challenges and Solutions. In *Proceedings of the 2009 Workshop on Arabic Natural Language Processing*.