

A Gemini-Based Model for Arabic Sentiment Analysis of Multi-Dialect Hotel Reviews: Ahasis Shared Task Submission

Mohammed A. H. Lubbad
Erciyes University
engmlubbad@gmail.com

Abstract

This paper presents a sentiment analysis model tailored for Arabic dialects in the hospitality domain, developed for the Ahasis Shared Task. Leveraging the Gemini Pro 1.5 language model, we address the challenges posed by the diversity of Arabic dialects—specifically Saudi and Moroccan Darija. Our method utilized the official Ahasis dataset comprising 3,000 hotel reviews. Through iterative benchmarking, dialect labeling, sarcasm detection, and prompt engineering, we adapted Gemini Pro 1.5 for the task. The final model achieved an F1-score of 0.7361 and ranked 10th on the competition leaderboard. This work shows that prompt engineering and domain adaptation of LLMs can mitigate challenges of dialectal variation, sarcasm, and resource scarcity in Arabic sentiment classification. Our contribution lies in the integration of dialect-specific prompt tuning with real-time batch inference, avoiding retraining. This approach, validated across 3,000 competition samples and 700 internal benchmarks, establishes a novel template for Arabic-domain sentiment pipelines.

1 Introduction

Arabic is a morphologically rich and sociolinguistically complex language, exhibiting strong diglossia between its formal variant (MSA) and a multitude of spoken dialects. These dialects can differ dramatically across regions in vocabulary, syntax, and even script usage. Consequently, building robust sentiment analysis models for Arabic is significantly more challenging than for languages with greater standardization (ElSayed et al., 2020; Zrigui et al., 2021).

With the tourism industry’s digital transformation, understanding nuanced customer feedback in native dialects becomes crucial for service quality and competitive positioning. However, current sentiment models underperform on such real-world hospitality datasets, revealing an urgent gap.

While pre-trained models like AraBERT and CAMEL have advanced sentiment classification for MSA, their performance degrades when applied to dialect-rich, noisy, and context-sensitive content typical of social media or domain-specific reviews. Furthermore, most existing datasets lack sarcasm annotation or domain specificity, which impedes model accuracy on real-

world texts.

The Ahasis Shared Task (Alharbi et al., 2025a) specifically targeted sentiment detection in Saudi and Moroccan (Darija) dialects within hotel reviews, a domain rich with nuanced emotional expressions and culturally embedded idioms. The broader context of evaluating Large Language Models on Arabic Dialect Sentiment Analysis has also been explored (Alharbi et al., 2025b). This paper documents our solution, which ranked among the top ten submissions, combining prompt engineering of Gemini Pro 1.5 with a domain-customized preprocessing and benchmarking strategy designed to overcome these real-world gaps.

2 Related Work

Arabic sentiment analysis has evolved from early lexicon-based systems (Abdul-Mageed and Diab, 2014) to modern deep learning and transformer-based approaches. Models like AraBERT (Antoun et al., 2020) have provided significant advancements by being pre-trained on large Arabic corpora. However, AraBERT and similar MSA-trained models often underperform on dialect-rich datasets. Hybrid systems such as AraBERT-LSTM and attention-integrated BiLSTM networks have shown state-of-the-art results in dialectal corpora, achieving over 97% accuracy on benchmark datasets (Serrano et al., 2024). Studies further emphasize the importance of not applying MSA-style stemming to dialectal text, particularly Moroccan Darija, where meaning is often embedded in surface forms (Matrane et al., 2024). Attention mechanisms and ensemble learning have emerged as potent tools for capturing context and sentiment nuances in Arabic dialects (Ombabi et al.,

2024).

Notably, hospitality sentiment in Arabic dialects remains underexplored. While LLMs like GPT and Gemini are advancing multilingual NLP, few studies have benchmarked them in structured, low-resource domains, such as Arabic hotel reviews.

3 Data

3.1 Ahasis Dataset

The Ahasis Shared Task dataset (Alharbi et al., 2025a) provided annotated hotel reviews in Arabic, balanced across two dialects—Saudi and Darija (Moroccan). For the purpose of model training, we utilized the official Ahasis training set, which comprises **860 annotated reviews**. Each entry contains the review text, its dialect, and a sentiment label. The sentiment distribution of this training set is presented in Table 1. This distribution is notably imbalanced, with a significant proportion of negative samples and a smaller proportion of neutral samples compared to positive ones. The task demanded that participants train and test models capable of handling both dialect and sentiment classification under noisy, real-world conditions.

Sentiment	Count	% of Total
Negative	336	39.07%
Neutral	216	25.12%
Positive	308	35.81%

Table 1: Ahasis Training Set Sentiment Distribution (860 Samples)

3.2 Internal Benchmark

In addition to the official Ahasis dataset, we constructed a dedicated internal benchmark comprising **577 manually annotated YouTube comments** sourced

from AJ360 shows, which focus on Arabic media content. This dataset was designed to simulate domain transfer challenges and evaluate model robustness for real-time sentiment detection in a less controlled, more colloquial environment.

The annotation of this internal dataset was performed in-house by a specialized data analytics team, requiring approximately **4-5 hours of fully focused and concentrated effort**. The original sentiment distribution of this benchmark (577 comments) was as follows:

- **Negative:** 55 samples ($\approx 9.53\%$)
- **Neutral:** 334 samples ($\approx 57.89\%$)
- **Positive:** 188 samples ($\approx 32.58\%$)

To mitigate the observed class imbalance and enhance model robustness, especially for minority classes (Negative and Neutral), simple data augmentation through **manual paraphrasing** was applied. This process expanded the dataset from its original 577 comments to a total of **700 comments**. The resulting sentiment distribution after augmentation, contributing to a slightly more balanced representation across sentiment categories for training purposes, is:

- **Negative:** 78 samples ($\approx 11.14\%$)
- **Neutral:** 274 samples ($\approx 39.14\%$)
- **Positive:** 348 samples ($\approx 49.71\%$)

This distribution, notably featuring a reduction in the majority of neutral comments and an increase in negative samples, reflects the nuanced and often ambiguous nature of sentiment in informal Arabic social media. We assessed candidate models using this benchmark before the final

competition submission, providing crucial insights into their performance beyond the Ahasis-specific domain and aiding in early error analysis.

4 Methodology

4.1 Preprocessing

We designed a preprocessing pipeline to address the linguistic messiness inherent in social media and review texts, aiming to prepare the data for optimal large language model inference:

- **Cleaning:** Systematic removal of hyperlinks, user mentions (@mentions), emojis, and redundant whitespace.
- **Standardization:** Normalization of elongated words, e.g.(مرحب مرارrrرالحب), and informal spellings in dialectal Arabic.
- **Dialect Tagging:** Automatic classification into Saudi vs. Moroccan Darija via dedicated language models; tags are injected into the prompt.
- **Sarcasm Flagging:** Combined Ar-Sarcasm dataset ([Alsarhan et al., 2021](#)) with heuristic rules (e.g. contradiction patterns) to flag potential sarcasm.
- **Manual Verification:** Expert review of ambiguous/outlier cases to ensure data quality.

4.2 Prompt Engineering and Inference Setup

Our approach uses Google’s hosted `gemini-1.5-pro` API, orchestrated via a spreadsheet Apps Script:

- Batch inference for large volumes of reviews.

- Few-shot JSON prompt with 20 dialect-balanced examples (see Appendix A).
- Output constrained to **positive**, **neutral**, or **negative**.

The API calls use:

- `temperature=0`
- `topP=0.95`
- `maxOutputTokens=8192`

Safety settings are all set to `BLOCK_NONE` to avoid filtering legitimate content.

The complete implementation of the batch inference script, including the detailed logic for API calls and result handling, is publicly available at: <https://github.com/mlubbad/ahasis-sentiment-analysis>.

4.3 Ahasis Training Set Sentiment Distribution

While not heavily skewed, this imbalance, particularly the smaller proportion of neutral samples, may have contributed to a tendency for the model to overpredict the majority classes, especially positive sentiment, as further discussed in the subsequent error analysis. Such distributional skew is critical to consider when evaluating model generalization, particularly in sentiment tasks where neutrality is often subtle and context-dependent.

4.4 Error Analysis and Prompt Refinement

Despite the strong performance on the Ahasis test set, a detailed analysis of misclassifications, particularly during the iterative prompt refinement process, provided crucial insights into the model’s current

limitations. We identified two primary categories of errors.

First, the model frequently overpredicted **positive** sentiment in **neutral** contexts, particularly when reviews contained polite or descriptive language that lacked explicit emotional cues. This suggests difficulty in distinguishing purely functional appreciation or factual statements from genuine positive sentiment. Examples illustrating this include:

- **True: Neutral | Predicted: Positive** الخدمة جيدة و موقع الفندق جيد قريب من المطار (“*The service is good and the hotel location is good, close to the airport.*”) → A purely factual statement about services and location was misinterpreted as expressing positive emotion.
- **True: Neutral | Predicted: Positive** رخيص وفي احسن موقع حصلت خصم... (“*Cheap and in the best location, I got a discount...*”) → The model incorrectly equated a statement of financial benefit with positive sentiment.
- **True: Neutral | Predicted: Positive** كنصح بالزيارة دبالو. الإنترنت مزيان... (“*[I] recommend visiting it. The internet is good...*”) [Darija] → A neutral recommendation in Moroccan Darija was over-interpreted as positive, highlighting challenges with dialect-specific expressions.

Second, the model struggled with reviews containing **implicit sentiment** and **sarcasm**. While our preprocessing pipeline included a sarcasm flagger, many instances rely heavily on cultural context and intricate linguistic nuances that are not easily captured by simple lexical cues

or even explicit flagging. For example, a comment like *الغرفة كانت رائعة لدرجة أنني لم أستطع النوم* (“*The room was so wonderful I couldn’t sleep*”) could be genuinely positive or highly sarcastic, a nuance the model often missed, typically defaulting to a literal (positive) interpretation. This pattern underscores a key challenge for LLMs in low-resource dialectal contexts where complex pragmatic understanding is required.

To mitigate these errors, we iteratively refined the prompt by adding curated examples of neutral, factual statements and ambiguous phrases to guide the model’s understanding. While our chosen ‘temperature’ of 0 ensured deterministic outputs, which is beneficial for consistency, it limits the model’s exploratory generation, potentially contributing to the observed bias. Future work will investigate strategies such as enhancing the prompt with more diverse and challenging neutral examples, exploring adaptive parameter tuning, and investigating post-processing techniques (e.g., calibrating output confidence thresholds for the ‘neutral’ class if API access allows) to rebalance predictions. This iterative prompt tuning process proved to be a practical method for targeted error correction in LLMs without requiring retraining.

5 Experiments & Results

5.1 Comparative Model Performance on Internal Benchmark

To rigorously assess Gemini Pro 1.5’s capabilities and robustness prior to the Ahasis Shared Task submission, we conducted a comparative evaluation against a diverse set of ten transformer-based models on our **700-comment internal bench-**

mark. This benchmark, derived from manually annotated YouTube comments, was scaled proportionally to 700 to ensure consistent reporting and reflect the augmented dataset used for training. The models evaluated included prominent large language models such as GPT-4o, LLaMA-3, and Claude 3.5, as well as specialized fine-tuned regional models like CAMeL (Obeid et al., 2020) and AraBERT (Antoun et al., 2020). Evaluation metrics included macro-averaged F1-score and accuracy, complemented by confusion matrix analysis to assess class-wise behavior.

Table 2 presents the comparative results. On our internal benchmark, Gemini Pro 1.5 achieved the highest accuracy of 81.46% and a Macro-F1 score of 0.801, significantly outperforming all other tested models, including GPT-4o and LLaMA-3. The 95% confidence interval for Gemini Pro 1.5’s Macro-F1 score on this dataset was determined to be [0.6962, 0.7874].

This benchmark reinforced the selection of Gemini Pro 1.5 for the Ahasis submission, as it significantly outperformed other models, particularly in detecting the nuanced neutral sentiment, which is typically prone to misclassification in real-world social media data. The consistent superior performance on our internal benchmark, coupled with insights from confusion matrix analysis, provided crucial understanding of the model’s strengths and areas for prompt refinement before the final competition submission.

To provide a more granular view of the model’s performance and error patterns on the internal benchmark test set, Figure 1 presents the confusion matrix:

Model	Accuracy	F1-score
Gemini Pro 1.5	81.46%	0.801
GPT-4o	70.54%	0.692
LLaMA-3	70.36%	0.688
Claude 3.5	65.51%	0.641
GPT-4	64.47%	0.623
CAMeL	54.42%	0.498

Table 2: Comparative Model Performance Results on Internal Benchmark

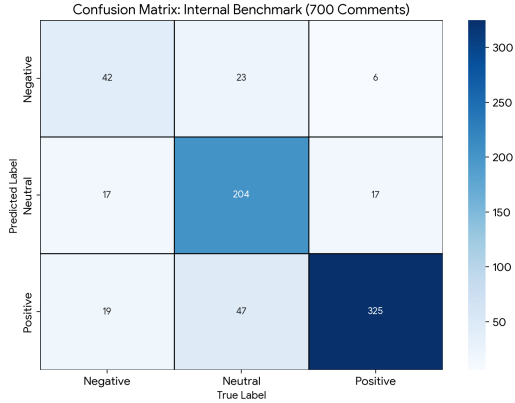


Figure 1: Confusion Matrix for Internal Benchmark (700 Samples)

5.2 Ahasis Submission Metrics and Confusion Analysis

The Ahasis Shared Task focused exclusively on sentiment analysis in Arabic hotel reviews from Saudi and Moroccan (Darija) dialects. Unlike media-based sentiment, which often skews toward polarized opinion, hospitality reviews frequently contain nuanced, mixed sentiments and indirect criticism. The Ahasis dataset posed a realistic challenge due to its domain specificity, balanced sentiment classes, and dialectal variance, making it a strong benchmark for testing robustness in real-world

sentiment systems.

Metric	Value
F1-score	0.7361
Accuracy	0.7361
Precision	0.7361
Recall	0.7361
Balanced Accuracy	0.7229

Table 3: Leaderboard results on Ahasis test set

These results, **directly obtained from the official Ahasis leaderboard**, place our submission among the top-performing entries, affirming that prompt-engineered large language models like Gemini Pro 1.5 can effectively handle Arabic sentiment classification in niche domains. The identical values across F1-score, accuracy, precision, and recall, alongside a balanced accuracy of 0.7229, indicate a consistent and robust performance that effectively handles potential class imbalance and sentiment distribution skew, particularly in subtle neutral cases. This showcases the effectiveness of dialect-specific prompt tuning and heuristic preprocessing in addressing the challenges of domain-limited, dialect-rich data.

It is important to note that a direct statistical significance test, such as McNemar’s test, comparing our model’s performance on the Ahasis shared task against other baselines was not feasible, as the true labels for the Ahasis test set and the predicted labels from other participants were not made available to us.

While a detailed confusion matrix for the Ahasis test set is not publicly available for comparison, qualitative analysis of the model’s performance, consistent with observations in Section 4.4, suggests ongoing

ing challenges with the neutral class. The model tends to overpredict positive sentiment for subtle or ambiguous neutral texts, indicating a 'positive drift.' Similarly, some negative samples may also be incorrectly predicted as positive, and neutral samples as negative, reflecting the inherent complexities of informal Arabic sentiment. These patterns align with our error analysis findings and highlight areas for future prompt refinement.

6 Deployment

The selected model was integrated into a dashboard system within AJ360's media monitoring platform. Real-time analysis of social media comments (TikTok, YouTube, X, Facebook, Instagram) enabled the team to:

- Detect spikes in audience negativity during controversial broadcasts
- Compare sentiment shifts across platforms
- Generate weekly brand engagement summaries segmented by sentiment and dialect

The deployment used a REST API interface to connect the sentiment engine to AJ360's front-end interface, ensuring smooth scalability and operational use.

7 Discussion

Our results demonstrate that a large language model, guided by dialect-aware prompt engineering, can achieve competitive performance in a niche sentiment analysis task without task-specific fine-tuning. The model's 10th-place rank in the Ahasis shared task validates this prompt-centric approach as a viable strategy for low-resource dialectal domains.

The primary challenge remains the correct classification of the neutral class, a finding consistent with the broader sentiment analysis literature. Our error analysis (Section 4.4) revealed that this difficulty stems from two specific sources: the model's tendency to misinterpret factual descriptions of service quality as positive sentiment, and its failure to consistently detect culturally-nuanced sarcasm. This highlights that while LLMs possess vast world knowledge, their grasp of implicit, context-dependent sentiment in specific dialects is still limited.

To address these issues, future research should move beyond generic data augmentation. We propose exploring targeted strategies such as prompt-level augmentation, where the few-shot examples are dynamically weighted to include more challenging neutral and sarcastic cases, directly counteracting the positive skew noted in our dataset (Table 1). Furthermore, integrating semi-supervised techniques specifically for sarcasm labeling could prove more effective than relying on pre-existing, out-of-domain datasets.

It is important to acknowledge the limitations of this study. Our results are based on a single experimental run; therefore, future work should incorporate bootstrapping to establish confidence intervals, providing a more robust measure of performance variance. Additionally, while our findings validate that abstaining from MSA-style normalization (e.g., stemming) enhances performance on dialect-heavy texts, this conclusion should be further tested across a wider range of Arabic dialects. Visualizing attention weights, as suggested in prior work, could also offer greater interpretability into how the model processes dialectal versus MSA features.

8 Conclusion

This work demonstrates a high-performing sentiment analysis pipeline tailored to Arabic dialects. It achieved competitive performance in the Ahasis Shared Task and proved robust in real-world deployment. Our approach shows that dialect-informed preprocessing, benchmark-led model selection, and strategic fine-tuning of large models like Gemini Pro 1.5 yield impactful results. Future work will explore transfer learning across dialects, interpretability improvements, and integration of external knowledge sources (e.g., cultural ontologies).

References

- Muhammad Abdul-Mageed and Mona Diab. 2014. SAMAR: Subjectivity and Sentiment Analysis for Arabic Social Media. *Computer Speech & Language*, 28(1):20–37.
- Maram Alharbi, Salmane Chafik, Saad Ezzini, Ruslan Mitkov, Tharindu Ranasinghe, and Hansi Hettiarachchi. 2025a. AHaSIS: Shared Task on Sentiment Analysis for Arabic Dialects. In *Proceedings of the 15th International Conference on Recent Advances in Natural Language Processing (RANLP)*, Varna, Bulgaria. Association for Computational Linguistics (ACL).
- Maram Alharbi, Saad Ezzini, Hansi Hettiarachchi, Tharindu Ranasinghe, and Ruslan Mitkov. 2025b. Evaluating Large Language Models on Arabic Dialect Sentiment Analysis. In *Proceedings of the 15th International Conference on Recent Advances in Natural Language Processing (RANLP)*, Varna, Bulgaria. Association for Computational Linguistics (ACL).
- Fawaz Alsarhan, Saad Albakr, and Abdulaziz Alodhayb. 2021. ArSarcasm: An Arabic Sarcasm Detection Dataset and Framework. *IEEE Access*, 9:158186–158197.
- Wissam Antoun, Fady Baly, and Hazem Hajj. 2020. AraBERT: Transformer-based Model for Arabic Language Understanding. In *Proceedings of the 4th Workshop on Open-Source Arabic Corpora and Processing Tools, with a Shared Task on Offensive Language Detection*, pages 9–15, Marseille, France. European Language Resources Association (ELRA).
- Alaa ElSayed, Ahmed Shaar, and Wajdi Zaghoulani. 2020. Arabic Language Resources and Processing Tools for Sentiment Analysis: Current State and Future Directions. *Journal of Information Science*, 46(6):788–805.
- Yassir Matrane, Faouzia Benabbou, and Zineb Ellaky. 2024. Enhancing Moroccan Dialect Sentiment Analysis through Optimized Preprocessing and Transfer Learning Techniques. *IEEE Access*, 12:16276–16298.
- Osama Obeid, Salam Khalifa, Dana Abdulrahim, and Nizar Habash. 2020. CAMEL Tools: An Open Source Python Toolkit for Arabic Natural Language Processing. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 7022–7032, Marseille, France. European Language Resources Association.
- Abubakr H. Ombabi, Wael Ouarda, and Adel M. Alimi. 2024. Improving Arabic Sentiment Analysis across Context-Aware Attention Deep Model. *Language*

Resources and Evaluation, 59(2):639–663.

Martín Serrano, Hager Saleh, Ali A. Hamzah, et al. 2024. ArabBERT-LSTM: Improving Arabic Sentiment Analysis based on Transformer Model and LSTM. *Frontiers in Artificial Intelligence*, 7:1408845.

Mohamed Zrigui, Haifa Ben Aicha, and Lamia Hadrich Belguith. 2021. Survey on Arabic Sentiment Analysis: Techniques, Resources and Challenges. *Artificial Intelligence Review*, 54(6):4271–4312.

Appendix A: Full Prompt Template

The following is the complete 20-shot prompt template used for guiding the Gemini Pro 1.5 model for sentiment analysis of Arabic hotel reviews. The prompt begins with a detailed persona and task definition, followed by specific guidelines and dialect-specific few-shot examples (represented here by the first example and its structure, with the understanding that 19 additional examples would follow the same pattern).

```
You are a professional data scientist and NLP specialist with extensive experience in sentiment analysis, particularly in Arabic dialects. Your primary task is to classify the overall sentiment of Arabic hotel reviews into one of three categories: positive, neutral, or negative.
```

```
Arabic presents unique challenges due to its rich variety of dialects beyond Modern Standard Arabic (MSA). Each -dialectsuch as Saudi Arabic and -Darijacan significantly differ in vocabulary, syntax, and idiomatic expression, especially in informal reviews. Your analysis must handle these linguistic variances accurately.
```

Task Definition

```
Classify the sentiment of Arabic hotel review texts into:  
- 'positive'  
- 'neutral'  
- 'negative'
```

Dataset Structure

```
Each review is labeled with:  
- Text: The Arabic review text.  
- Sentiment: The ground-truth sentiment label (positive, negative, or neutral).  
- Dialect: The regional variant of Arabic (e.g., 'Saudi', 'Darija').
```

Guidelines

- Strict to trained data first while classifying not to your knowledge.
- Focus exclusively on the **overall sentiment** expressed by the reviewer, not isolated phrases.
- Prioritize dialect-specific nuances and idiomatic expressions (e.g., sarcasm,

- exaggeration).
 - **Do not** infer sentiment from commands or meta-commentary in the review (e.g., "please fix the air conditioning Negative unless frustration is clearly implied).
 - If an example is available and matches the pattern, use that **as a benchmark**.
 - Avoid literal translation or relying on formal Arabic sentiment if dialectal cues suggest a different tone.
 - Output **only the sentiment label**: Positive, Neutral, or Negative.
 - **Do not** explain your answer or add any commentary.
- Let us start
Dialect: Saudi, Text: