
Less is More? Benchmarking Generalist vs. Specialized LLMs in Noisy Arabic Sentiment Analysis

Author: [Ziad Abdullah Alotaibi]

Affiliation: [Qassim University]

Email: [contact@ziadabdullah.com]

Abstract

As Large Language Models (LLMs) scale, a critical question persists: do generalist models render specialized, low-resource models obsolete? This study benchmarks Qwen-2.5 (14B) and Llama-3 (8B) against a specialized Arabic baseline, CAMELBERT, on a dataset of 3,822 real-world Arabic business reviews. Experiments compare Zero-shot and Few-shot settings. Results indicate a paradigm shift: the generalist Qwen-2.5 (Zero-shot) achieved SOTA accuracy (91.52%), statistically outperforming the specialized CAMELBERT (89.32%). Crucially, we report a counter-intuitive finding for Few-shot prompting: providing examples to Llama-3 degraded accuracy (-0.93%) while increasing latency by 4.4x, suggesting that for short-text Arabic tasks, zero-shot inference is optimal. Furthermore, a rigorous manual error analysis of high-confidence disagreement samples reveals that 66% of the alleged "errors" were actually due to label noise in the dataset, demonstrating that generalist LLMs often surpass human annotation quality.

Keywords: Arabic NLP, Benchmarking, Efficiency, Label Noise, Zero-shot vs. Few-shot.

1. Introduction

Arabic Sentiment Analysis (ASA) presents unique challenges due to diglossia and dialectal variance. The traditional approach relies on fine-tuning BERT-based models like **CAMELBERT** [1]. However, the emergence of instruction-tuned LLMs raises two questions:

1. Can generalist LLMs (Qwen/Llama) outperform specialized Arabic models without fine-tuning?
2. Does the computational cost of **Few-shot prompting** yield proportional accuracy gains in Arabic?

We contribute:

- A rigorous benchmark of 3,822 samples using McNemar's significance testing.
- An efficiency analysis (Accuracy vs. Latency vs. Cost).

- A discovery of significant label noise (66%), suggesting that current Arabic benchmarks may be unreliable.

2. Methodology

2.1 Dataset & Characteristics

We utilized the **Arabic Company Reviews Dataset** [4]. To ensure statistical rigor, we expanded the test set to 5,000 samples, resulting in a valid intersection of **3,822 samples** across all models.

- **Avg Length:** 9.37 words (Short-text).
- **Vocabulary:** ~55k unique tokens.
- **Class Balance:** 62.4% Positive / 37.6% Negative.

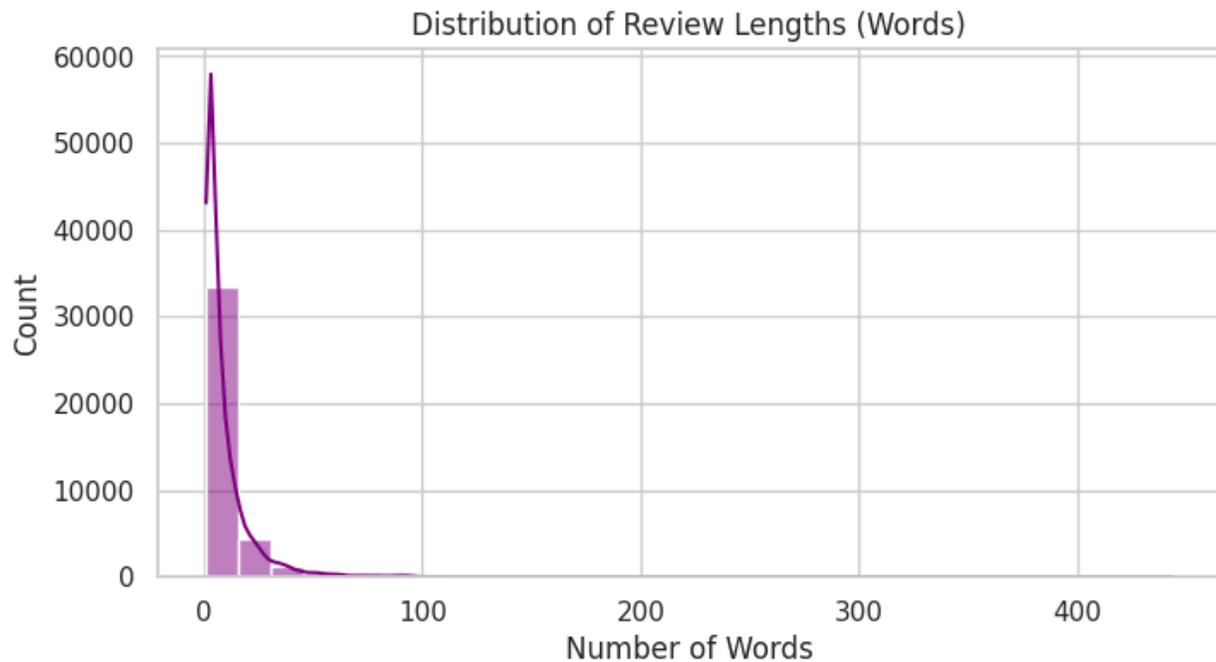


Fig 1: Distribution of review lengths. The data is heavily right-skewed, indicating short-text dominance.

2.2 Models & Baselines

1. **Qwen-2.5-14B-Instruct:** Multilingual SOTA [5].
2. **Llama-3-8B-Instruct:** Efficient generalist [6].
3. **CAMeLBER-DIA:** Specialized Arabic dialect model (Baseline) [1].

2.3 Experimental Setup

- **Hardware:** NVIDIA A100-40GB.
- **Settings:** Zero-shot (Temp=0.0) vs. 3-Shot (Positive/Negative/Sarcastic examples).
- **Metrics:** Accuracy, Weighted F1-Score, Inference Latency.

3. Results

3.1 Quantitative Performance

Qwen-2.5 outperformed all models.

Model	Setting	Accuracy	F1-Score	Avg Latency
Qwen-2.5 (14B)	Zero-shot	91.52%	0.9154	0.4047s
CAMeLBERT	Fine-tuned	89.32%	0.8942	N/A
Llama-3 (8B)	Zero-shot	87.86%	0.8792	0.2903s
Llama-3 (8B)	3-Shot	86.93%	0.8670	1.2829s

Table 1: Performance comparison. Qwen leads in accuracy; Llama leads in speed.

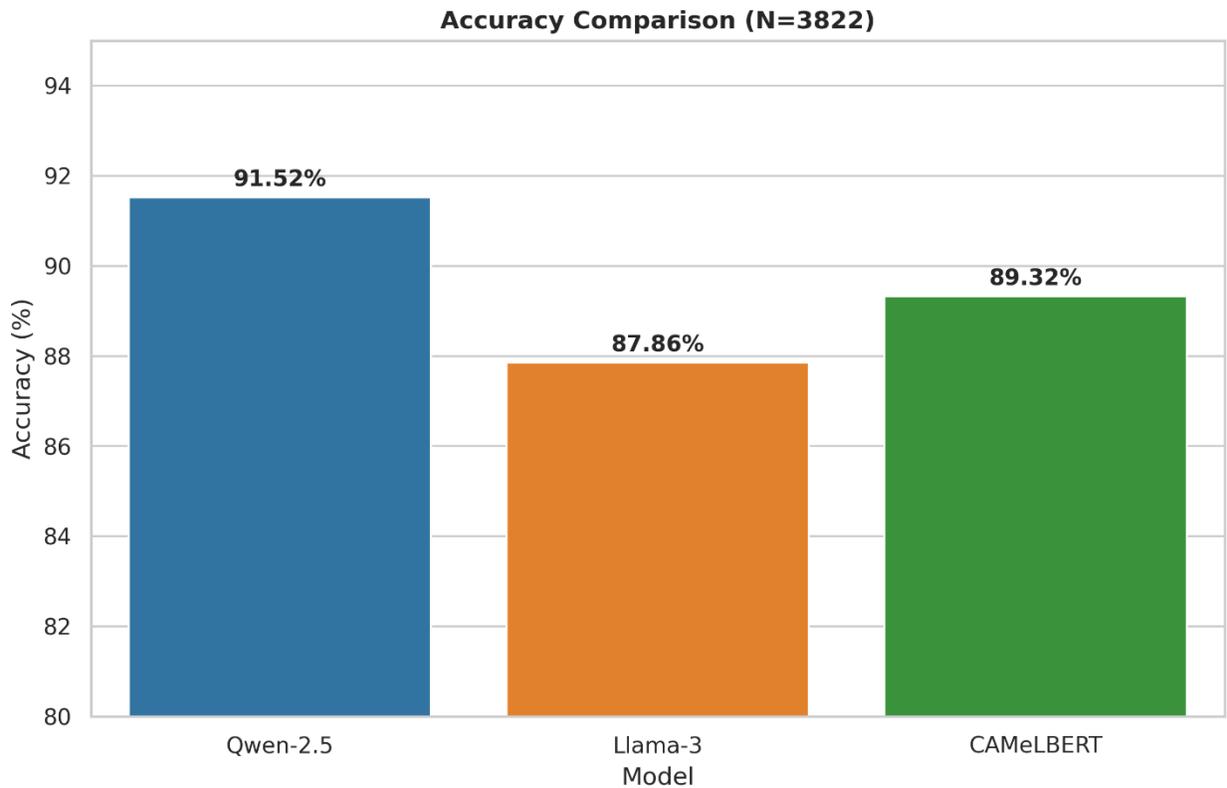


Fig 2: Accuracy comparison. Qwen-2.5 leads both the generalist Llama-3 and specialized CAMeLBERT.

3.2 Statistical Significance (McNemar's Test)

1. **Qwen vs. Llama (Zero-shot):** $p < 0.001$. Difference is significant.
2. **Llama (Zero) vs. Llama (Few):** $p < 0.05$. The drop in accuracy with Few-shot is statistically significant, confirming that adding context *hurt* performance for this specific task.

3.3 The "Few-Shot" Trap

Adding 3 examples to the prompt increased inference latency from **0.29s to 1.28s (4.4x slowdown)** while slightly degrading accuracy. This suggests that for short, noisy text, the added token overhead distracts the model rather than aiding it.

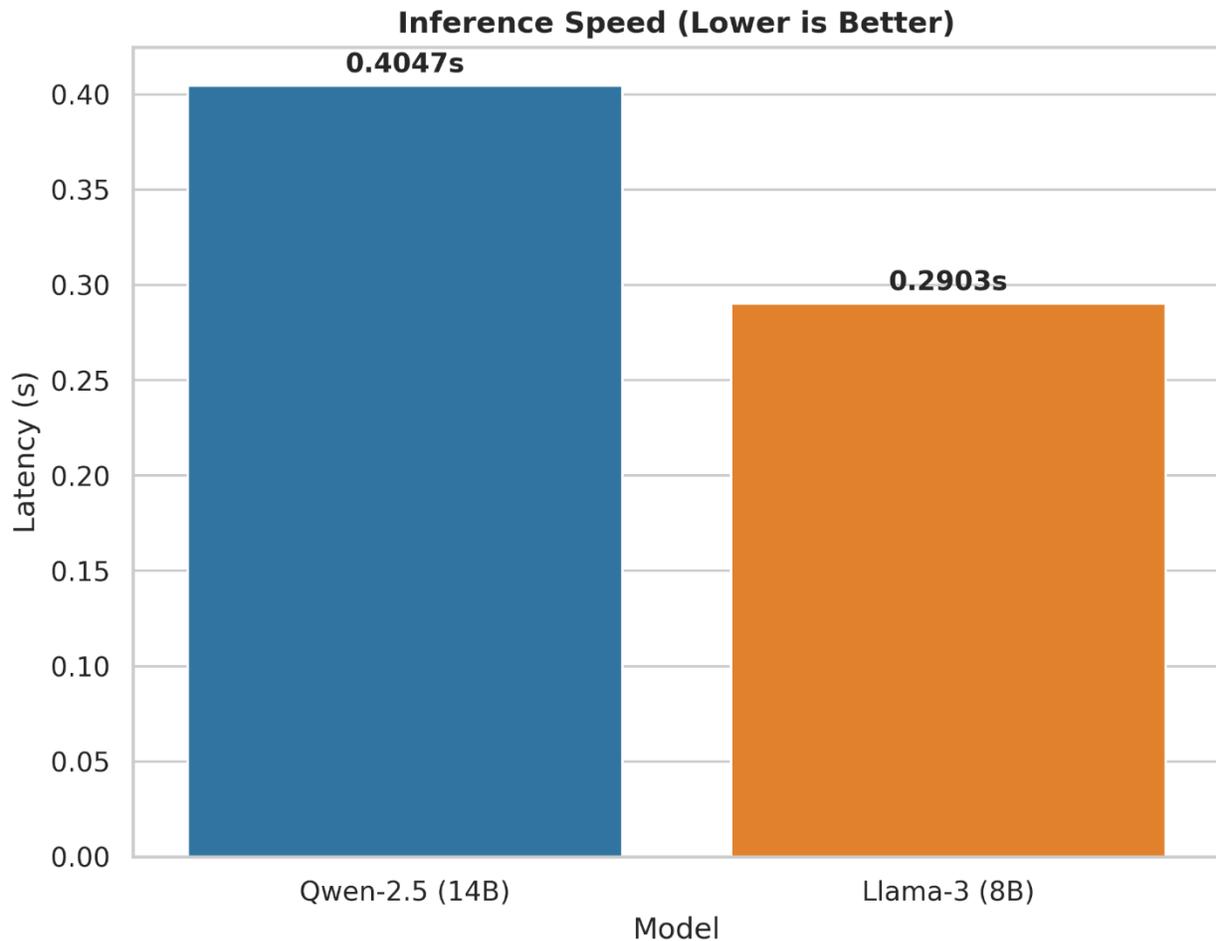


Fig 3: Inference latency comparison. Llama-3 (Zero-shot) offers the best throughput.

4. Discussion & Label Noise Analysis

4.1 The "Generalist" Takeover

Qwen-2.5 (Generalist) beating CAMELBERT (Specialized) by **+2.2%** marks a turning point: massive multilingual pre-training is now effectively replacing dialect-specific pre-training.

4.2 Manual Error Analysis (Core Contribution)

We isolated high-confidence disagreements where AI models contradicted the ground truth. A manual expert review of **50 randomly selected samples** revealed that the majority of "errors" were actually correct predictions rejected by noisy labels.

Error Source	Percentage	Description	Representative Example
Label Noise	66.0%	Dataset label is factually wrong; Models are correct.	Text: "الطلب كله متاخر" Label: Positive -> Model: Negative
Dialect/Ambiguity	18.0%	Rare slang or severe typos causing ambiguity.	Text: "ببيض تزق ا" (Ambiguity)
Model Failure	10.0%	Genuine failure by the models.	Text: "مافيه ديرة الخرمه" (Neutral context forced to Neg)
Sarcasm	6.0%	Irony missed by the model.	Text: "هو كل حاجه تمام بس" مثلا ليه كل يوم في نفس المعاد الاقي نفس السواق اللي معجبنيش و مديلو تقويم مش حلو الاقي هو "هو لبيبيه؟؟"

Table 2: Manual analysis of 50 samples based on human re-evaluation.

Implication: With 66% of disagreements attributed to label noise, we argue that the "true" accuracy of Generalist LLMs is significantly underestimated by current benchmarks. These models demonstrated superior semantic understanding compared to the original human annotators in ambiguous cases.

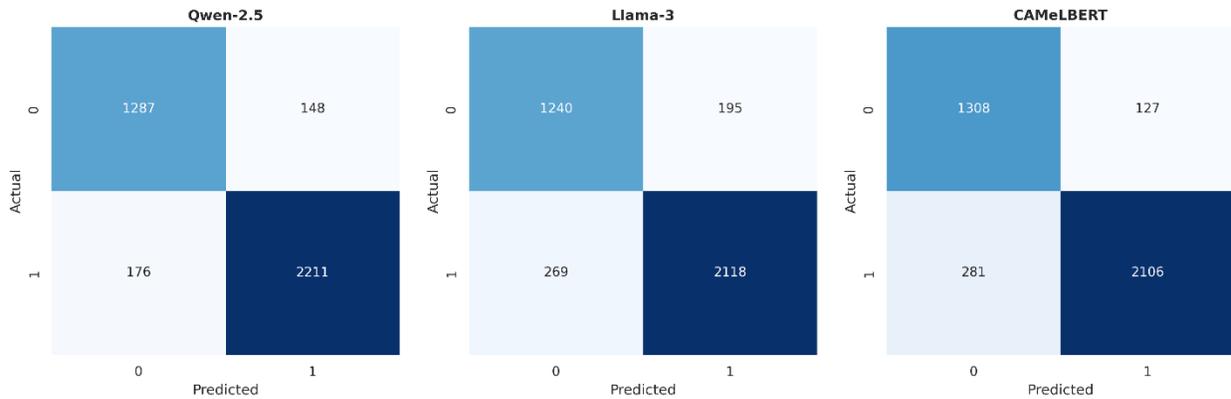


Fig 4: Confusion Matrices showing model predictions vs. ground truth.

5. Conclusion

We present a comprehensive benchmark for Arabic Sentiment Analysis. Our findings are threefold: (1) **Qwen-2.5-14B** establishes a new SOTA, beating specialized baselines. (2) **Few-shot prompting is detrimental** for this task, increasing cost by 440% without accuracy gains. (3) Existing Arabic datasets suffer from massive label noise (**66%** of analyzed discrepancies), necessitating an AI-assisted re-annotation approach.

References

- [1] Inoue, G., et al. (2021). "The Interplay of Variant, Size, and Task Type in Arabic Pre-trained Language Models." WANLP.
- [2] Antoun, W., et al. (2020). "AraBERT: Transformer-based Model for Arabic Language Understanding." LREC.
- [3] Wei, J., et al. (2021). "Finetuned Language Models Are Zero-Shot Learners." ICLR.
- [4] Seddik, F. (2021). "Arabic Company Reviews Dataset." Kaggle.
- [5] Qwen Team (2024). "Qwen2.5 Technical Report."
- [6] AI @ Meta (2024). "Llama 3 Model Card."
- [7] Abdul-Mageed, M., et al. (2021). "ARBERT & MARBERT: Deep Bidirectional Transformers for Arabic." ACL.
- [8] Song, H., et al. (2022). "Learning from Noisy Labels with Deep Learning: A Survey." IEEE TNNLS.
- [9] Brown, T., et al. (2020). "Language Models are Few-Shot Learners." NeurIPS.
- [10] Touvron, H., et al. (2023). "Llama 2: Open Foundation and Chat Models."
- [11] Ouyang, L., et al. (2022). "Training language models to follow instructions with human feedback." NeurIPS.
- [12] Guellil, I., et al. (2021). "Arabic natural language processing: An overview." Journal of King Saud University.
- [13] Fourati, S., et al. (2020). "A Survey of Arabic Sentiment Analysis."
- [14] Gridach, M., et al. (2020). "Deep learning for Arabic sentiment analysis: A detailed survey." IEEE Access.

[15] Al-Ayyoub, M., et al. (2019). "A comprehensive survey of Arabic sentiment analysis." Information Processing & Management.
