

AraHarMeme: A Multimodal Benchmark for Harmful Meme Detection in Arabic

Anonymous ACL submission

Abstract

In this work, we introduce AraHarMeme, a benchmark dataset for harmful content detection in Arabic memes. The dataset consists of 5,313 Arabic social media images collected from diverse platforms and manually annotated into three classes: harmful memes, non-harmful memes, and non-meme content. To the best of our knowledge, AraHarMeme is the first Arabic multimodal benchmark that explicitly addresses harmful meme detection. We conduct experiments using text-only, image-only, supervised multimodal models, as well as large vision-language models evaluated under zero-shot settings. The results show that Arabic-focused text models consistently outperform multilingual baselines, while image-only models are insufficient for reliable harmful content detection. Supervised multimodal models yield improvements when strong Arabic text encoders are combined with appropriate visual backbones, whereas zero-shot large vision-language models remain considerably less effective than supervised approaches on this dataset. The dataset and evaluation protocol will be publicly released to support future research in Arabic multimodal safety and content moderation.

Content Warning. This paper and the introduced dataset contain examples of harmful and offensive content (e.g., hate speech, harassment, or abusive imagery) included solely for the purpose of scientific research. Reader discretion is advised.

1 Introduction

Harmful content on social media remains a persistent challenge for automated moderation systems, particularly as abusive or hateful messages are increasingly expressed in indirect and context-dependent ways (Kiela et al., 2020). While much prior work has focused on textual abuse and hate speech (Talat and Hovy, 2016), a growing portion

of harmful online content is conveyed through images combined with short textual overlays.

Internet memes are a common example of such content. Memes often rely on the interaction between image and text, as well as shared cultural knowledge, to convey meaning. As a result, harmful intent may not be apparent from either modality alone. This observation has motivated multimodal benchmarks for English, most notably the Hateful Memes dataset (Kiela et al., 2020), which demonstrated that unimodal approaches are insufficient for reliable detection. More recent work has further examined multimodal and large vision-language models for meme-based safety tasks, revealing mixed performance despite strong general multimodal capabilities (Alayrac et al., 2022; Liu et al., 2023).

In contrast, harmful meme detection in Arabic remains largely unexplored. Arabic social media content is characterized by extensive dialectal variation, frequent code-switching, and non-standard orthography, all of which complicate automatic analysis (Antoun et al., 2020; Abdul-Mageed et al., 2021). Existing Arabic datasets for abusive or hateful content focus almost exclusively on text-only settings and do not account for the multimodal structure of memes. Although recent Arabic multimodal resources have begun to emerge (Alwajih et al., 2025), none explicitly target harmful meme detection as a content moderation task.

To address this gap, we introduce **AraHarMeme**, a multimodal benchmark for harmful meme detection in Arabic. The dataset consists of 5,313 images collected from social media and annotated into three categories: harmful memes, non-harmful memes, and non-meme content. The inclusion of a non-meme class reflects realistic moderation scenarios in which visually similar but semantically unrelated images are common.

We accompany AraHarMeme with a standardized evaluation protocol and baseline experi-

084 ments spanning unimodal, multimodal, and large
085 language-vision models. Together, these compo-
086 nents provide a foundation for studying harmful
087 meme detection in Arabic and for analyzing the
088 limitations of current modeling approaches in a
089 dialect-rich, multimodal setting. The main contri-
090 butions are summarized as follows:

- 091 • We introduce **AraHarMeme**, a multimodal
092 Arabic benchmark explicitly designed for
093 harmful meme detection.
- 094 • We provide manually curated annotations
095 covering harmful, non-harmful, and non-
096 meme content.
- 097 • We establish strong baseline results, high-
098 lighting key challenges in Arabic harmful
099 meme detection under multimodal settings.
- 100 • We assess the quality and stability of
101 the benchmark through cross-dataset evalu-
102 ation on an existing Arabic meme dataset
103 (AraMeme) (Alam et al., 2024), analyzing
104 the consistency of model performance across
105 related tasks.
- 106 • We will publicly release the dataset and eval-
107 uation protocol to support future research on
108 Arabic multimodal content moderation.

109 2 Related Work

110 Early research on harmful content detection fo-
111 cused primarily on text-based settings, addressing
112 hate speech, abusive language, and offensive con-
113 tent on social media platforms (Talat and Hovy,
114 2016; Davidson et al., 2017). While effective for
115 explicit abuse, text-only approaches struggle when
116 harmful intent is implicit or context-dependent.

117 Memes present a particularly challenging case.
118 As multimodal artifacts combining images with
119 short, informal text, memes often convey mean-
120 ing through interactions between modalities rather
121 than through either modality alone. This has mo-
122 tivated growing interest in multimodal approaches
123 to harmful content detection.

124 2.1 Multimodal Harmful Meme Detection

125 The Hateful Memes benchmark demonstrated
126 that unimodal models are insufficient for reliable
127 meme-based harm detection, establishing multi-
128 modal reasoning as a core requirement for this task
129 (Kiela et al., 2020). Subsequent work explored

130 multimodal fusion architectures that combine vi-
131 sual encoders with pretrained language models,
132 generally outperforming text-only and image-only
133 baselines (Goel and Poswal, 2024; Arya et al.,
134 2024).

135 More recent studies have examined additional
136 aspects of harmful meme detection, including ex-
137 plainability and robustness. Lin et al. (Lin et al.,
138 2024) propose explainable multimodal models
139 that expose conflicting visual and textual evidence,
140 while evaluations of large vision–language models
141 such as Flamingo and LLaVA show mixed results
142 for harmful meme classification despite strong per-
143 formance on general multimodal tasks (Alayrac
144 et al., 2022; Liu et al., 2023). Recent analy-
145 ses further highlight persistent safety limitations
146 of large multimodal models in real-world meme
147 settings (Lee et al., 2025). To address robust-
148 ness concerns, retrieval-augmented and sentiment-
149 aware approaches have been proposed, reporting
150 improvements on hateful or harmful meme bench-
151 marks (Mei et al., 2025; Duan et al., 2025). Most
152 of this work, however, focuses on English data.

153 2.2 Arabic and Low-Resource Multimodal 154 Benchmarks

155 Arabic NLP research has largely concentrated on
156 text-only tasks such as sentiment analysis and hate
157 speech detection, reflecting challenges related to
158 dialectal variation and non-standard orthography
159 (Pramanick et al., 2021). Multimodal resources
160 for Arabic have only recently begun to emerge.
161 AraMeme introduces an Arabic meme dataset fo-
162 cused on propagandistic content but is not de-
163 signed for harmfulness detection as a moderation
164 task (Alam et al., 2024).

165 More recent efforts expand Arabic multimodal
166 coverage beyond memes, including culturally
167 grounded instruction datasets and vision–language
168 benchmarks that highlight performance gaps
169 across dialects (Alwajih et al., 2025). Community-
170 driven shared tasks have also released multimodal
171 Arabic datasets for classification and propaganda-
172 related challenges (Hasanain et al., 2024). De-
173 spite this progress, there is still no benchmark ex-
174 plicitly designed for harmful meme detection in
175 Arabic that supports systematic evaluation across
176 text-only, image-only, supervised multimodal, and
177 zero-shot vision–language models. This gap moti-
178 vates the introduction of AraHarMeme.



Figure 1: Example of an Arabic harmful meme requiring multimodal reasoning. The harmful intent arises from the interaction between the overlaid Arabic text and the image, while neither modality alone is sufficient to determine harmfulness.

3 AraHarMeme Benchmark

AraHarMeme is designed to reflect key linguistic, cultural, and orthographic characteristics of Arabic social-media discourse while capturing the multimodal nature of harmful content in memes. Figure 1 illustrates a representative instance in which harmful intent arises from the interaction between overlaid Arabic text and the image, while neither modality alone is sufficient to determine harmfulness. In this example, the text, glossed as "*Blacks are quiet*", can be interpreted as a descriptive statement in isolation, and the image alone depicts a person holding a machete-like object without explicit harmful framing. When interpreted jointly, however, the combination gives rise to a harmful, derogatory implication.

Each data instance therefore consists of a meme image, OCR-extracted text, and a manually assigned harm label, supporting both binary and multi-class evaluation settings. Figure 2 summarizes the AraHarMeme dataset construction and annotation pipeline.

In addition to multimodal complexity, the dataset exhibits substantial linguistic diversity: as shown in Figure 3 (Appendix C), AraHarMeme covers a wide range of Arabic dialects, including Egyptian, Gulf, Levantine, and Maghrebi Arabic, as well as Arabizi and code-switched Arabic-English content. Orthographic variation is also pervasive in the OCR-extracted text (Figure 4, Appendix C), with frequent non-standard spelling, character elongation, emojis, and visually stylized fonts characteristic of Arabic memes.

3.1 Dataset Collection

The AraHarMeme dataset was collected from multiple publicly accessible online platforms, including Twitter/X, Facebook public pages, Telegram public channels and Reddit communities (e.g., r/ArabMemes). These platforms were selected due to their high prevalence of informal, dialectal, and visually grounded content.

Data collection relied on keyword-based crawling combined with image scraping. Seed keywords were manually curated to cover a broad range of harmful and non-harmful meme content commonly observed in Arabic social media. These included (i) harm-related terms (e.g., insults, harassment, and derogatory expressions), (ii) dialectal and colloquial variants of such terms across major Arabic dialects, and (iii) meme-related cues frequently used in Arabic memes. Representative examples include insult terms such as "كلب" (dog), "حمار" (donkey), and dialectal variants (e.g., "كلبك", "حماررر", "اسكت" (shut up) and "عار" (shame), as well as colloquial spellings and elongations. To reduce bias toward Modern Standard Arabic, keywords were expanded to include Egyptian, Gulf, Levantine, and Maghrebi forms, Arabizi spellings (e.g., kalb (dog) , 7mar (donkey)), and common code-switched Arabic-English expressions. The complete list of keywords is provided in Appendix A.

3.2 Dataset Filtering

Following collection, a structural filtering stage was applied to remove low-quality or unsuitable content while preserving multimodal diversity. Near-duplicate images were identified and removed using perceptual hashing, retaining a single representative instance per duplicate cluster. Images were discarded if they were severely corrupted, extremely low-resolution, or visually irrelevant to meme discourse. Importantly, no filtering decisions at this stage were based on harmfulness, which was determined exclusively during manual annotation.

3.3 OCR Text Extraction

Text embedded in images was extracted using *Tesseract OCR* (v5.3)¹ with Arabic language models. OCR was applied uniformly to all images regardless of font style, text density, or visual lay-

¹<https://github.com/tesseract-ocr/tesseract>

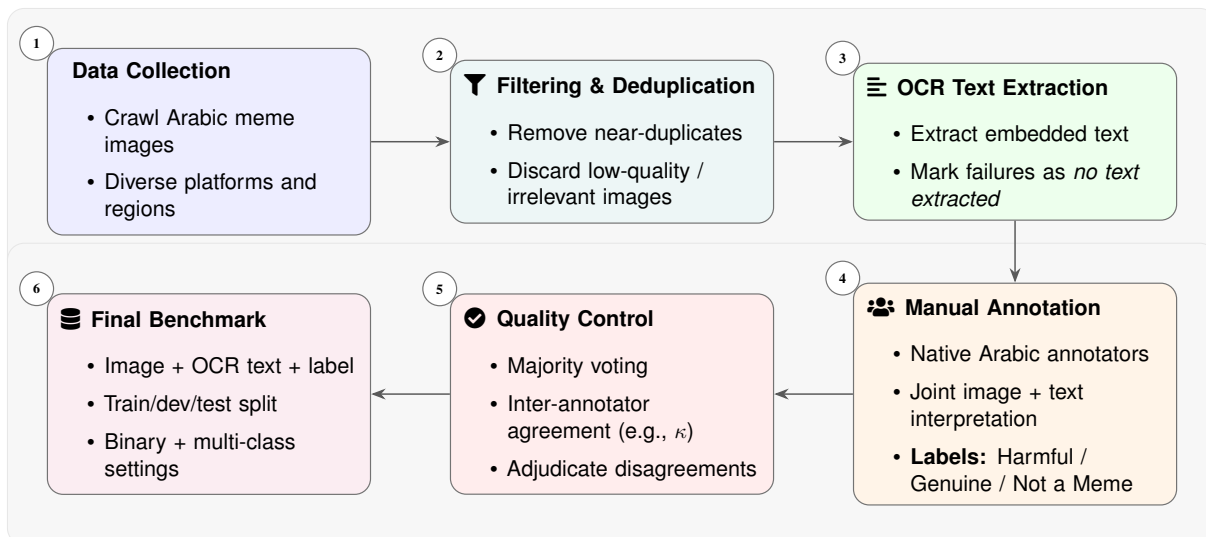


Figure 2: Overview of the AraHarMeme dataset construction and annotation pipeline.

out. We do not perform any automatic OCR post-processing or correction at this stage. In particular, we do not normalize spelling, remove noise, or manually edit OCR outputs during extraction. Instead, any necessary corrections are performed only during the manual annotation phase to ensure that the final text faithfully reflects what is visibly written on the meme.

An OCR result is considered *unsuccessful* if it:

1. contained only isolated characters or fragments that did not form interpretable words,
2. consisted of corrupted or nonsensical strings caused by stylized fonts, heavy noise, or low resolution, or
3. clearly did not correspond to any visible textual content in the image.

Instances meeting these criteria were labeled as "*no text extracted*". This designation was assigned only after manual verification during annotation.

3.4 Manual Annotation

Instances were annotated manually by a team of three native Arabic speakers with academic training in Arabic linguistics and demonstrated expertise in regional Arabic dialects. Prior to annotation, annotators were provided with detailed written guidelines describing label definitions, examples of implicit harm, sarcasm, and culturally grounded references, as well as instructions for jointly interpreting visual and textual content.

Annotators were explicitly instructed to consider both modalities together when assigning labels. All harmfulness judgments were made exclusively at this stage and did not influence earlier filtering or OCR-related decisions. OCR-extracted text was verified and corrected when necessary during annotation, particularly in cases involving dialectal spelling, stylized fonts, or OCR noise. Corrections were strictly limited to faithfully reflecting the visible text in the image and did not involve reinterpretation, paraphrasing, or semantic modification.

The annotation schema consists of three classes:

- **Harmful:** Memes that convey offensive, abusive, harassing, hateful, or otherwise harmful content. Harmfulness may be expressed explicitly (e.g., direct insults or abusive language) or implicitly through image–text interaction, sarcasm, or contextual references.
- **Genuine:** Memes that are non-harmful and contain neutral, benign, or light humorous content without offensive or derogatory intent.
- **Not a Meme:** Images that lack meme-specific characteristics, including standalone photographs, decorative images, screenshots of unrelated content, or images without communicative or semantic intent.

3.5 Quality Control and Agreement

The dataset instances independently annotated by all three annotators. Final labels were determined

using majority voting. In cases where all three annotators disagreed, the instance was reviewed jointly and resolved through discussion in accordance with the annotation guidelines.

Inter-annotator agreement was measured using pairwise Cohens κ , yielding agreement scores of 0.72, 0.64, and 0.85 across annotator pairs, with an average $\kappa = 0.74$, indicating substantial agreement. This high level of agreement reflects consistent interpretation despite the presence of dialectal variation, implicit harm, and multimodal ambiguity.

3.6 Dataset statistics

The final AraHarMeme benchmark consists of **5,323 images**, split into 4,189 training instances and 1,134 test instances. The dataset includes 1,992 harmful memes, 2,039 genuine (non-harmful) memes, and 1,282 non-meme images, as summarized in Table 2.

Appendix C provides additional dataset statistics. Figure 3 shows the distribution of Arabic dialects, with most instances containing dialectal Arabic rather than Modern Standard Arabic. Figure 4 presents the character-level length distribution of OCR-extracted text, which varies substantially across instances due to informal writing styles and OCR-related noise, and Figure 5 reports the class distribution across the three labels.

3.7 Comparison with Prior Arabic Meme Datasets

Although multimodal content such as memes has been widely studied in high-resource languages like English, work on Arabic meme datasets is still emerging. One of the first resources to explicitly target Arabic meme understanding is ArMeme (Alam et al., 2024), a multimodal corpus of approximately 6,000 memes annotated for propagandistic content. ArMeme serves as an initial benchmark for multimodal research in Arabic and has been utilized in multimodal classification and shared tasks on propaganda detection (Haouhat et al., 2024; Shah et al., 2024). An extension of ArMeme, known as MemeXplain, augments the base dataset with bilingual explanations for the propaganda labels, enabling explainable multimodal classification research (Kmainasi et al., 2025).

Despite their value, these resources have limitations relative to AraHarMeme. ArMeme and MemeXplain are primarily geared toward detect-

ing propagandistic content in memes, whereas AraHarMeme adopts a broader conception of harmful content, encompassing general offensive, abusive, and harmful themes beyond propaganda.

Other Arabic multimodal corpora, such as general visual-text pairs for sentiment analysis or cultural representation, exist but do not specifically address memes or harmful multimodal content relevant to content moderation, making AraHarMeme unique in its targeted task focus. Table 1 compares AraHarMeme with existing Arabic multimodal meme datasets.

4 Experiments

4.1 Experimental Setup

All experiments are conducted on the AraHarMeme benchmark described in Section 3. Each instance consists of an image, OCR-extracted Arabic text (when available), and a manually assigned label from three classes: Harmful, Genuine, and Not a Meme. A fixed train-test split is used throughout all experiments. Test instances are strictly excluded from all training and fine-tuning procedures. The same test split is used for supervised learning and zero-shot evaluation to ensure consistent and comparable results across experimental settings. All supervised experiments are implemented using the HuggingFace Transformers library² and executed on a single GPU.

We report Accuracy, Precision, Recall, and F1-score for all supervised models. For zero-shot large language models, we additionally report macro-averaged Precision, Recall, and F1-score to better capture per-class behavior under imbalanced predictions.

4.2 Text-only Models

Text-only experiments are implemented using pretrained transformer models for sequence classification. We evaluate both **Arabic-specific** models, including ARBERT (Abdul-Mageed et al., 2021), MARBERT and MARBERTv2 (Abdul-Mageed et al., 2021), CamelBERT (Inoue et al., 2021), and Asafaya (Safaya et al., 2020), as well as **multilingual baselines** such as XLM-R (Conneau et al., 2020), mBERT (Devlin et al., 2019), and DistilMBERT (Sanh et al., 2019). All models are initialized from their publicly released pretrained checkpoints.

²<https://huggingface.co>

Table 1: Comparison of AraHarMeme with existing Arabic multimodal meme datasets (provided for reference).

Dataset	Size	Language	Coverage	Modalities	Annotation Types	Primary Task(s)
ArMeme (Alam et al., 2024)	~6K	Arabic (MSA + dialects)	Arabic (MSA + dialects)	Image+Text	Propaganda labels	Propaganda detection
MemeXplain (Kmainasi et al., 2025)	~6K	Arabic (MSA + dialects)	Arabic (MSA + dialects)	Image+Text	Propaganda labels + explanations	Explainable propaganda classification
AraHarMeme (ours)	~5K	Dialectal Arabizi, switching	Arabic, code-switching	Image+Text	Harmful / Non-harmful / Not a Meme	Broad harmful meme detection

Statistic	Value
Total memes	5,313
Training set	4,189
Test set	1,134
Genuine (non-harmful)	2,039
Harmful	1,992
Not a meme	1,282

Table 2: Statistics of the AraHarMeme dataset.

Fine-tuning is performed using the HuggingFace Trainer API (Wolf et al., 2020). Models are trained for four epochs with a batch size of 16, using the AdamW optimizer (Loshchilov and Hutter, 2017) with a learning rate of 2×10^{-5} . A linear learning rate scheduler with warm-up over the first 10% of training steps is applied. We use standard cross-entropy loss for all classification experiments. No additional task-specific pretraining or data augmentation is applied. Evaluation is performed on the held-out test set using standard classification metrics.

4.3 Image-only Models

Image-only experiments are conducted using standard convolutional and transformer-based vision architectures, including ResNet-34 and ResNet-50 (He et al., 2016), ViT-Base (Dosovitskiy, 2020), DeiT-Base (Touvron et al., 2021), and ConvNeXt-Tiny (Liu et al., 2022). All models are initialized with ImageNet-pretrained weights.

Image models are trained for four epochs using the same optimizer, learning rate, and batch size as the text-based models to maintain a controlled comparison. Standard image preprocessing and resizing are applied. No textual or multimodal information is used in this setting.

4.4 Supervised Multimodal Models

Supervised multimodal experiments combine a pretrained text encoder with a visual backbone through a late fusion strategy. Textual and visual

representations are extracted independently and concatenated before the final classification layer.

All multimodal models are trained using the same hyperparameters as unimodal models to ensure a controlled comparison: four epochs, batch size 16, the AdamW optimizer (Loshchilov and Hutter, 2017), and a learning rate of 2×10^{-5} . No modality-specific loss weighting, curriculum learning, or auxiliary supervision is applied. This design isolates the contribution of multimodal fusion without introducing additional training complexity.

4.5 Zero-shot Evaluation of Large Language-Vision Models

We evaluate a set of **Large Language-Vision Models (LLVMs)** in a zero-shot setting. All evaluated models are natively multimodal, i.e., capable of processing both visual and textual inputs. To better understand their behavior under different input conditions, we evaluate each model using three configurations: *text-only*, *image-only*, and *text-image*.

The evaluated models include **Gemma-3 (4B, 12B)**(Team et al., 2025), **LLaVA-7B³**, **LLaVA-Phi3-3.8B⁴**, **Moondream-1.8B⁵**, **Qwen2.5-VL-3B⁶**, and **Qwen3-VL-30B⁷**. These models are selected to cover a range of architectures and parameter scales representative of recent large language-vision systems.

All evaluations are conducted in a zero-shot manner. No task-specific fine-tuning, parameter updates, or in-context examples are used. Instead, models are prompted using a fixed instruction prompt that defines the task, the label space, and explicit decision rules. The prompt instructs

³<https://llava-vl.github.io/>

⁴<https://ollama.com/library/llava-phi3:3.8b>

⁵<https://moondream.ai>

⁶<https://github.com/QwenLM/Qwen2.5-VL>

⁷<https://github.com/QwenLM/Qwen3-VL>

the model to classify each input into one of three categories: *Harmful Arabic Meme*, *Not Harmful Arabic Meme*, or *Not a Meme / Not Arabic*. Models are explicitly instructed not to guess missing text and not to rely on external knowledge beyond the visible input. The full prompt used for all zero-shot experiments is provided in Appendix B.

In the *text-only* configuration, models receive only the OCR-extracted text when available. In the *image-only* configuration, only the image is provided. In the *text-image* configuration, both modalities are provided jointly. Model outputs are mapped deterministically to the target label set based on the predicted class number. Inference is performed using deterministic decoding to ensure reproducibility across runs.

5 Results and Discussion

This section analyzes the performance of supervised text-only, image-only, and multimodal models, as well as zero-shot large language-vision models, on the AraHarMeme benchmark. Results for supervised models are reported in Table 3, while zero-shot results are summarized in Table 4.

Text-only models. Among text-only approaches, Arabic-specific transformer models consistently outperform multilingual baselines. As shown in Table 3, ARBERT achieves the highest performance (Acc = 0.820, F1 = 0.819), followed closely by MARBERT (Acc = 0.818, F1 = 0.817). Other Arabic-focused models, including Asafaya and CamelBERT, show comparable performance, whereas multilingual models such as XLM-R, mBERT, and DistilMBERT perform noticeably worse. This gap indicates that Arabic-specific pretraining remains essential for modeling dialectal variation, informal usage, and non-standard orthography characteristic of Arabic social-media content.

Image-only models. Image-only models perform substantially worse than text-based classifiers. Although modern vision architectures such as ConvNeXt-Tiny and DeiT-Base achieve moderate performance (F1 around 0.72--0.73), they consistently lag behind text-only models. Both convolutional and transformer-based vision models exhibit limited recall, suggesting that visual information alone is often insufficient for identifying harmful intent, particularly in cases where harm is implicit or linguistically mediated.

Supervised multimodal models. Supervised multimodal models show more varied behavior. When strong Arabic text encoders are combined with appropriate visual backbones, multimodal fusion yields substantial gains over unimodal baselines. In particular, the ARBERT + ConvNeXt-Tiny model achieves the best overall performance (Acc = 0.951, F1 = 0.944), outperforming both text-only and image-only models. These improvements are reflected in balanced precision and recall, indicating that visual features can complement textual representations when linguistic modeling is sufficiently robust. In contrast, multimodal models built on MARBERT or multilingual encoders yield smaller or inconsistent gains, suggesting that multimodal fusion does not compensate for weaker linguistic representations.

Zero-shot large language-vision models. Table 4 reports zero-shot performance for large language-vision models evaluated under text-only, image-only, and multimodal input configurations. Across all models, zero-shot performance remains substantially lower than that of supervised text-based and multimodal approaches. Providing both image and text generally improves performance relative to unimodal inputs; however, the gains are limited. Qwen3-VL-30B achieves the strongest zero-shot results (Acc = 0.615, F1 = 0.579 in the multimodal setting), but still falls well short of fine-tuned Arabic-specific text models and supervised multimodal architectures. These results indicate that current large language-vision models are not competitive for Arabic harmful meme detection without task-specific adaptation.

To further assess the robustness of these observations, we evaluate the same subset of large language-vision models on the ArMeme dataset. The results, presented in Table 5 (Appendix D), show qualitatively similar behavior: the relative ranking of models is largely preserved, multimodal inputs generally help, and overall zero-shot performance remains limited. This cross-dataset consistency suggests that the performance patterns observed on AraHarMeme are not dataset-specific, but instead reflect broader limitations of current multimodal LLMs for Arabic meme understanding.

6 Conclusion

This paper introduced **AraHarMeme**, a multimodal benchmark for harmful meme detection in

Table 3: Performance comparison of text-only, image-only, and supervised multimodal models on the AraHarMeme benchmark. Macro-averaged metrics are reported for all models. Best results within each block are highlighted in bold.

Model	Acc	P	R	F1
Text-only Models				
ARBERT	0.820	0.835	0.820	0.819
MARBERT	0.818	0.835	0.818	0.817
Asafaya	0.808	0.822	0.808	0.807
CamelBERT	0.808	0.826	0.808	0.805
MARBERTv2	0.792	0.805	0.792	0.789
XLM-R	0.776	0.789	0.776	0.772
mBERT	0.742	0.753	0.742	0.736
DistilMBERT	0.748	0.758	0.748	0.744
Image-only Models				
ConvNeXt-Tiny	0.729	0.726	0.729	0.727
DeiT-Base	0.722	0.723	0.722	0.721
ResNet50	0.703	0.709	0.703	0.705
ViT-Base	0.687	0.693	0.687	0.689
ResNet34	0.686	0.679	0.686	0.676
Supervised Multimodal Models				
ARBERT + ConvNeXt-T	0.951	0.947	0.941	0.944
ARBERT + DeiT-B	0.895	0.887	0.879	0.883
MARBERT + DeiT-B	0.705	0.708	0.681	0.690
MARBERT + ConvNeXt-T	0.700	0.672	0.665	0.666
mBERT + DeiT-B	0.709	0.744	0.649	0.658
DistilMBERT + ConvNeXt-T	0.660	0.660	0.602	0.608
XLM-R + ConvNeXt-T	0.648	0.590	0.595	0.592
MARBERT + ViT-B	0.656	0.613	0.616	0.613
ARBERT + ResNet50	0.615	0.604	0.589	0.593

581 Arabic, designed to reflect the linguistic diver-
582 sity, cultural context, and implicit multimodal na-
583 ture of Arabic social-media content. Extensive
584 experiments show that fine-tuned Arabic-specific
585 text models substantially outperform multilingual
586 baselines, that image-only models are insufficient
587 for reliable harmfulness detection, and that super-
588 vised multimodal models are effective only when
589 paired with strong Arabic-aware linguistic repre-
590 sentations. In contrast, zero-shot large language-
591 vision models remain markedly less competitive,
592 even when provided with both image and text.

593 Future work will explore larger-scale Arabic
594 multimodal pretraining, more fine-grained harm
595 taxonomies, and improved cross-modal fusion
596 strategies that better capture implicit and cultur-
597 ally grounded harmful content. We release Ara-
598 HarMeme and the accompanying evaluation pro-
599 tocol to support further research on Arabic multi-
600 modal safety and content moderation.

Table 4: Zero-shot performance of large language and multimodal LLMs on AraHarMeme. Macro-averaged metrics are reported. Best results within each block are in bold.

Model	Input	Acc	F1	P	R
Gemma-3 12B					
Text		0.484	0.353	0.310	0.411
Image		0.491	0.342	0.495	0.409
Multi		0.506	0.424	0.573	0.448
Gemma-3 4B					
Text		0.365	0.362	0.376	0.403
Image		0.456	0.317	0.538	0.386
Multi		0.462	0.435	0.454	0.438
LLaVA-7B					
Text		0.394	0.329	0.381	0.344
Image		0.261	0.242	0.320	0.329
Multi		0.387	0.314	0.346	0.340
LLaVA-Phi3-3.8B					
Text		0.354	0.269	0.330	0.330
Image		0.387	0.248	0.219	0.342
Multi		0.346	0.270	0.394	0.345
Moondream-1.8B					
Text		0.404	0.194	0.468	0.335
Image		0.307	0.247	0.210	0.328
Multi		0.328	0.233	0.189	0.309
Qwen2.5-VL-3B					
Text		0.317	0.260	0.209	0.378
Image		0.305	0.288	0.472	0.400
Multi		0.378	0.284	0.238	0.353
Qwen3-VL-30B					
Text		0.595	0.548	0.589	0.547
Image		0.553	0.502	0.615	0.507
Multi		0.615	0.579	0.633	0.569

Ethical Considerations

601 AraHarMeme consists exclusively of publicly
602 available meme images and does not contain any
603 user-identifying information. No private accounts,
604 restricted groups, or personal metadata were ac-
605 cessed or collected, and the dataset therefore poses
606 minimal privacy risk.
607

608 As with any dataset involving harmful or offen-
609 sive content, the annotations reflect human judg-
610 ment and may encode cultural or societal biases.
611 We provide clear annotation guidelines and en-
612 courage users of the dataset to remain aware of
613 these limitations when developing or deploying
614 models. Models trained on AraHarMeme may sup-
615 port research on Arabic content moderation and
616 contribute to tools for analysts, journalists, and
617 platform moderators. However, we caution against
618 using such models as fully automated decision-

619 makers without human oversight.

620 Limitations

621 AraHarMeme was collected from publicly acces-
622 sible social media platforms and reflects the vari-
623 ability of real-world Arabic meme content. As a
624 result, the benchmark is sensitive to OCR quality.
625 Although OCR outputs were manually reviewed
626 during annotation, stylized fonts, low-resolution
627 images, emojis, and complex layouts may intro-
628 duce noise in fully automated settings, potentially
629 affecting model performance.

630 In addition, harmfulness annotation is inher-
631 ently subjective. While annotations were per-
632 formed by native Arabic speakers with familiar-
633 ity across dialects, interpretations of offensive or
634 harmful content can vary across regions and com-
635 munities, which may limit generalizability.

636 Use of AI Assistance

637 In the preparation of this work, we used AI-
638 assisted tools to support code development (e.g.,
639 GitHub Copilot) and for limited language-related
640 assistance, such as spelling checks and minor
641 stylistic refinements (e.g., Grammarly and Chat-
642 GPT). All scientific content, analysis, and conclu-
643 sions were produced and verified by the authors.

644 References

645 Muhammad Abdul-Mageed, AbdelRahim Elmadany,
646 and El Moatez Billah Nagoudi. 2021. [ARBERT &
647 MARBERT: Deep bidirectional transformers for ara-
648 bic](#). In *Proceedings of the 59th Annual Meeting of
649 the Association for Computational Linguistics and
650 the 11th International Joint Conference on Natural
651 Language Processing (Volume 1: Long Papers)*, On-
652 line. Association for Computational Linguistics.

653 Firoj Alam, Abul Hasnat, Fatema Ahmad, Md Arid
654 Hasan, and Maram Hasanain. 2024. [Armeme:
655 Propagandistic content in arabic memes](#). In *Pro-
656 ceedings of the 2024 Conference on Empirical
657 Methods in Natural Language Processing*, pages
658 21071--21090.

659 Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc,
660 Antoine Miech, Iain Barr, Yana Hasson, Karel
661 Lenc, Arthur Mensch, Katherine Millican, Malcolm
662 Reynolds, and 1 others. 2022. [Flamingo: a vi-
663 sual language model for few-shot learning](#). *Ad-
664 vances in neural information processing systems*,
665 35:23716--23736.

666 Fakhraddin Alwajih, Samar M. Magdy, Abdel-
667 lah El Mekki, Omer Nacar, Youssef Nafea,
668 Safaa Taher Abdelfadil, Abdulfattah Mohammed

669 Yahya, Hamzah Luqman, Nada Almarwani, Samah
670 Aloufi, Baraah Qawasmeh, Houdaifa Atou, Serry
671 Sibae, Hamzah A. Alsayadi, Walid Al-Dhabyani,
672 Maged S. Al-shaibani, Aya El aatar, Nour Qan-
673 dos, Rahaf Alhamouri, and 18 others. 2025. [Pearl:
674 A multimodal culturally-aware Arabic instruction
675 dataset](#). In *Findings of the Association for
676 Computational Linguistics: EMNLP 2025*, pages
677 23048--23079, Suzhou, China. Association for Com-
678 putational Linguistics.

679 Wissam Antoun, Fady Baly, and Hazem Hajj. 2020.
680 [AraBERT: Transformer-based model for Arabic lan-
681 guage understanding](#). In *Proceedings of the 4th
682 Workshop on Open-Source Arabic Corpora and Pro-
683 cessing Tools, with a Shared Task on Offensive Lan-
684 guage Detection*, pages 9--15, Marseille, France. Eu-
685 ropean Language Resource Association.

686 Greeshma Arya, Mohammad Kamrul Hasan, Ashish
687 Bagwari, Nurhizam Safie, Shayla Islam, Fatima
688 Rayan Awad Ahmed, Aaishani De, Muhammad At-
689 tique Khan, and Taher M Ghazal. 2024. [Multi-
690 modal hate speech detection in memes using con-
691 trastive language-image pre-training](#). *IEEe Access*,
692 12:22359--22375.

693 Alexis Conneau, Kartikay Khandelwal, Naman Goyal,
694 Vishrav Chaudhary, Guillaume Wenzek, Francisco
695 Guzmán, Edouard Grave, Myle Ott, Luke Zettle-
696 moyer, and Veselin Stoyanov. 2020. [Unsupervised
697 cross-lingual representation learning at scale](#). In
698 *Proceedings of the 58th Annual Meeting of the
699 Association for Computational Linguistics*, pages
700 8440--8451, Online. Association for Computational
701 Linguistics.

702 Thomas Davidson, Dana Warmusley, Michael Macy,
703 and Ingmar Weber. 2017. [Automated hate speech
704 detection and the problem of offensive language](#). In
705 *Proceedings of the 11th International AAAI Confer-
706 ence on Web and Social Media*.

707 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and
708 Kristina Toutanova. 2019. [BERT: Pre-training of
709 deep bidirectional transformers for language under-
710 standing](#). In *Proceedings of the 2019 Conference of
711 the North American Chapter of the Association for
712 Computational Linguistics: Human Language Tech-
713 nologies, Volume 1 (Long and Short Papers)*, pages
714 4171--4186, Minneapolis, Minnesota. Association
715 for Computational Linguistics.

716 Alexey Dosovitskiy. 2020. [An image is worth 16x16
717 words: Transformers for image recognition at scale](#).
718 *arXiv preprint arXiv:2010.11929*.

719 Yuxiao Duan, Xiang Zhao, and Hao Guo. 2025.
720 [Sentiment-aware cross-modal semantic interaction
721 model for harmful meme detection](#). *Decision Sup-
722 port Systems*, page 114509.

723 Aman Goel and Abhishek Poswal. 2024. [Mmhs: Multi-
724 modal model for hate speech intensity prediction](#). In
725 *International Conference on Speech and Computer*,
726 pages 95--108. Springer.

727	Abdelhamid Haouhat, Hadda Cherroun, Slimane Bel-	<i>vances in neural information processing systems</i> ,	784
728	laouar, and Attia Nehar. 2024. Modos at araieval	36:34892--34916.	785
729	shared task: Multimodal propagandistic memes		
730	classification using weighted sam, clip and arabi-	Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Fe-	786
731	angpt . In <i>Proceedings of the Second Arabic Natural</i>	ichtenhofer, Trevor Darrell, and Saining Xie. 2022.	787
732	<i>Language Processing Conference</i> , pages 483--488,	A convnet for the 2020s. In <i>Proceedings of the</i>	788
733	Bangkok, Thailand. Association for Computational	<i>IEEE/CVF conference on computer vision and pat-</i>	789
734	Linguistics.	<i>tern recognition</i> , pages 11976--11986.	790
735	Maram Hasanain, Md. Arid Hasan, Fatema Ahmed,	Ilya Loshchilov and Frank Hutter. 2017. Decou-	791
736	Reem Suwaileh, Md. Rafiul Biswas, Wajdi Zag-	pled weight decay regularization. <i>arXiv preprint</i>	792
737	ghouani, and Firoj Alam. 2024. Araieval shared	<i>arXiv:1711.05101</i> .	793
738	task: Propagandistic techniques detection in uni-		
739	modal and multimodal arabic content . In <i>Pro-</i>	Jingbiao Mei, Jinghong Chen, Guangyu Yang, Weizhe	794
740	<i>ceedings of the Second Arabic Natural Language</i>	Lin, and Bill Byrne. 2025. Robust adaptation of	795
741	<i>Processing Conference (ArabicNLP 2024)</i> , pages	large multimodal models for retrieval augmented	796
742	456--466. Association for Computational Linguis-	hateful meme detection . In <i>Proceedings of the</i>	797
743	tics.	<i>2025 Conference on Empirical Methods in Natural</i>	798
744	Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian	<i>Language Processing</i> , pages 23806--23828, Suzhou,	799
745	Sun. 2016. Deep residual learning for image recog-	China. Association for Computational Linguistics.	800
746	nition. In <i>Proceedings of the IEEE conference</i>		
747	<i>on computer vision and pattern recognition</i> , pages	Shraman Pramanick, Shivam Sharma, Dimitar Dim-	801
748	770--778.	itrov, Md. Shad Akhtar, Preslav Nakov, and Tanmoy	802
749	Go Inoue, Bashar Alhafni, Nurpeis Baimukan, Houda	Chakraborty. 2021. MOMENTA: A multimodal	803
750	Bouamor, and Nizar Habash. 2021. The interplay	framework for detecting harmful memes and their	804
751	of variant, size, and task type in Arabic pre-trained	targets . In <i>Findings of the Association for Computa-</i>	805
752	language models . In <i>Proceedings of the Sixth Ara-</i>	<i>tional Linguistics: EMNLP 2021</i> , pages 4439--4455,	806
753	<i>bic Natural Language Processing Workshop</i> , pages	Punta Cana, Dominican Republic. Association for	807
754	92--104, Kyiv, Ukraine (Virtual). Association for	Computational Linguistics.	808
755	Computational Linguistics.		
756	Douwe Kiela, Hamed Firooz, Aravind Mohan, Vedanuj	Ali Safaya, Moutasem Abdullatif, and Deniz Yuret.	809
757	Goswami, Amanpreet Singh, Pratik Ringshia, and	2020. KUISAIL at SemEval-2020 task 12: BERT-	810
758	Davide Testuggine. 2020. The hateful memes chal-	CNN for offensive speech identification in social me-	811
759	lenge: Detecting hate speech in multimodal memes.	dia . In <i>Proceedings of the Fourteenth Workshop on</i>	812
760	<i>Advances in neural information processing systems</i> ,	<i>Semantic Evaluation</i> , pages 2054--2059, Barcelona	813
761	33:2611--2624.	(online). International Committee for Computational	814
762	Mohamed Bayan Kmainasi, Abul Hasnat, Md Arid	Linguistics.	815
763	Hasan, Ali Ezzat Shahroor, and Firoj Alam. 2025.	Victor Sanh, Lysandre Debut, Julien Chaumond, and	816
764	MemeIntel: Explainable detection of propagandistic	Thomas Wolf. 2019. Distilbert, a distilled version	817
765	and hateful memes . In <i>Proceedings of the 2025 Con-</i>	of bert: smaller, faster, cheaper and lighter . <i>arXiv</i>	818
766	<i>ference on Empirical Methods in Natural Language</i>	<i>preprint arXiv:1910.01108</i> .	819
767	<i>Processing</i> , pages 30263--30279, Suzhou, China.		
768	Association for Computational Linguistics.	Uzair Shah, Md. Rafiul Biswas, Marco Agus, Mowafa	820
769	DongGeon Lee, Joonwon Jang, Jihae Jeong, and	Househ, and Wajdi Zaghouani. 2024. Mememind at	821
770	Hwanjo Yu. 2025. Are visionlanguage models safe	araieval shared task: Generative augmentation and	822
771	in the wild? a meme-based benchmark study. In	feature fusion for multimodal propaganda detection	823
772	<i>Proceedings of the 2025 Conference on Empirical</i>	in arabic memes . In <i>Proceedings of the Second Ara-</i>	824
773	<i>Methods in Natural Language Processing (EMNLP)</i> ,	<i>bic Natural Language Processing Conference</i> , pages	825
774	pages 30533--30576. Association for Computational	467--472, Bangkok, Thailand. Association for Com-	826
775	Linguistics.	putational Linguistics.	827
776	Hongzhan Lin, Ziyang Luo, Wei Gao, Jing Ma,	Zeerak Talat and Dirk Hovy. 2016. Hateful symbols or	828
777	Bo Wang, and Ruichao Yang. 2024. Towards ex-	hateful people? predictive features for hate speech	829
778	plainable harmful meme detection through multi-	detection on twitter . In <i>Proceedings of the NAACL</i>	830
779	modal debate between large language models. In	<i>student research workshop</i> , pages 88--93.	831
780	<i>Proceedings of the ACM Web Conference 2024</i> ,		
781	pages 2359--2370.	Gemma Team, Aishwarya Kamath, Johan Ferret,	832
782	Haotian Liu, Chunyuan Li, Qingyang Wu, and	Shreya Pathak, Nino Vieillard, Ramona Mer-	833
783	Yong Jae Lee. 2023. Visual instruction tuning. <i>Ad-</i>	hej, Sarah Perrin, Tatiana Matejovicova, Alexan-	834
		dre Ramé, Morgane Rivière, and 1 others. 2025.	835
		Gemma 3 technical report. <i>arXiv preprint</i>	836
		<i>arXiv:2503.19786</i> .	837
		Hugo Touvron, Matthieu Cord, Matthijs Douze, Fran-	838
		cisco Massa, Alexandre Sablayrolles, and Hervé	839

Jégou. 2021. Training data-efficient image transformers & distillation through attention. In *International conference on machine learning*, pages 10347--10357. PMLR.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, and 1 others. 2020. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 conference on empirical methods in natural language processing: system demonstrations*, pages 38--45.

A Keyword-Based Crawling

Seed keywords were manually curated and grouped into categories to guide keyword-based crawling of Arabic memes. The list below provides representative examples used during data collection.

- **Harm-related insults (MSA and colloquial):**

كلب، حمار، وسخ، قدر، غبي، حقير، نجس، متخلف،
حتالة، زبالة، ساقط

- **Threatening or silencing expressions:**

اسكت، انخرس، اطبق فك، هدد، اقتلك، اضربك،
دمرك، سحقك، نهايتك، اتبه لنفسك

- **Derogatory and humiliating phrases:**

عار، عيب، فضيحة، مهزلة، مسخرة، ضحكة، عيب عليك،
وصمة، نخزي

- **Sarcastic and provocative cues:**

برافو عليك، شاطر، عبقرى زمانك، يا سلام، ذكي جداً،
مثال يحتذى به

- **Dialectal variants and regional slang:**

- **Egyptian:**

حيوان، أهبل، عبيط، معفن، ولا مؤاخدة

- **Levantine:**

حمارك، غيبان، وسخين، بلا نخ

- **Gulf:**

تافه، فاشل، هطف، ما تسوى

- **Maghrebi:**

حماررر، زبلة، بهلول، مسخ

- **Orthographic variations and elongation:**

حماااار، كليب، وسخنخ، غيببي، قذذذر

- **Arabizi spellings (Latin-script Arabic):**

kalb, 7mar, weskh, ghabi, wes5, 3ayb, 5ara,
enta 7mar

- **Code-switched Arabic-English expressions:**

stupid عرب dirty ناس crazy ناس idiot عربي
ناس stupid كلام trash

- **Meme-related discourse markers:**

لما، لما تكون، لما تشوف، هذا لما، توقع، تخيل، القصة
(POV) باختصار

The keyword list was iteratively expanded to include dialectal spellings, colloquial usage, and non-standard orthography commonly observed in Arabic social-media memes. These keywords were used solely to guide crawling and candidate collection; final inclusion and labeling decisions were performed through manual annotation.

B Zero-shot Classification Prompt

This appendix presents the exact instruction prompt used for zero-shot evaluation in this work. The prompt was originally designed for multimodal input, where models receive both an image and OCR-extracted text. The same prompt was applied verbatim in text-only and image-only settings by providing only the available modality and omitting the missing one. No task-specific fine-tuning, in-context examples, or external knowledge were used.

Classification Instructions

You will receive an image and text. Classify the input into one of the following categories.

1 -> **Harmful Arabic Meme:** The image is a meme. The visible text is written exclusively in Arabic or Arabic dialects (e.g., Moroccan, Egyptian, Algerian, Tunisian). The content contains insults, profanity, hate, discrimination, mockery, or degrading humor.

2 -> **Not Harmful Arabic Meme:** The image is a meme. The visible text is written exclusively in Arabic or Arabic dialects. The content is neutral, harmless, or contains light humor.

3 -> **Not a Meme / Not Arabic:** The image is not a meme, contains no readable text, or the visible text is predominantly non-Arabic and unrelated to Arabic meme discourse.

Rules: If any non-Arabic text is present, output 3. Do not guess missing text. Do not use external knowledge.

Extracted text: {caption}

Output format: <class_number>
<explanation>

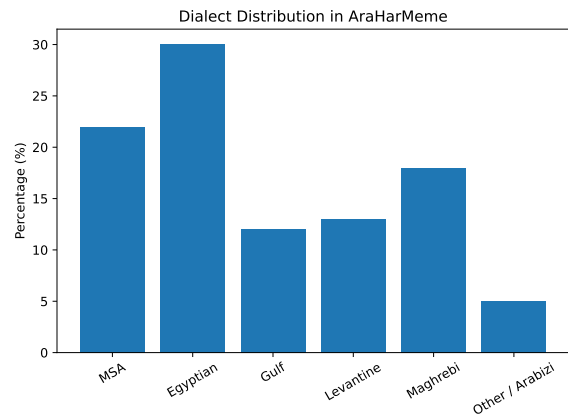


Figure 3: Dialect distribution in the AraHarMeme benchmark.

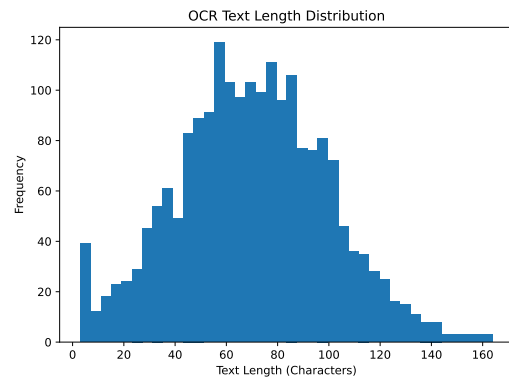


Figure 4: Character-level text length distribution of OCR-extracted meme text.

C Dataset Statistics and Distributions

This appendix provides additional statistics and distributional visualizations for the AraHarMeme benchmark, complementing the dataset description in Section 3.

D Additional Experiments on ArMeme

We further evaluate the curated benchmark by testing a subset of large language models on the existing ArMeme dataset. Table 5 reports the zero-shot performance of this subset of LLMs.

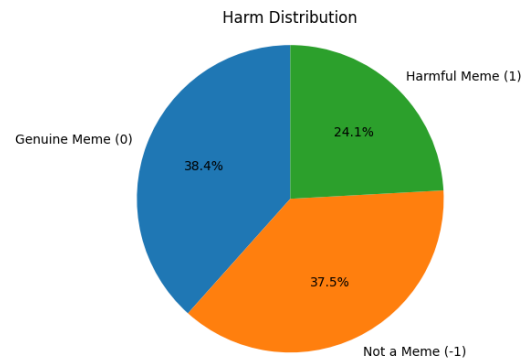


Figure 5: Label distribution across the three annotation categories.

Table 5: Zero-shot performance of selected LLMs on ArMeme. Macro-averaged metrics are reported. Best results within each model block are highlighted in bold.

Model	Input	Acc	F1	P	R
Gemma-3 12B					
	Text	0.123	0.136	0.489	0.271
	Image	0.237	0.138	0.179	0.216
	Multi	0.278	0.229	0.432	0.343
Gemma-3 4B					
	Text	0.258	0.115	0.231	0.250
	Image	0.655	0.200	0.193	0.250
	Multi	0.499	0.238	0.236	0.258
LLaVA 7B					
	Text	0.397	0.206	0.225	0.247
	Image	0.655	0.198	0.164	0.250
	Multi	0.580	0.239	0.228	0.254
LLaVA-Phi3 3.8B					
	Text	0.438	0.222	0.228	0.254
	Image	0.635	0.228	0.244	0.257
	Multi	0.470	0.232	0.232	0.260
Moondream 1.8B					
	Text	0.658	0.198	0.164	0.250
	Image	0.621	0.222	0.219	0.249
	Multi	0.489	0.226	0.220	0.240
Qwen2.5-VL 3B					
	Text	0.303	0.156	0.230	0.244
	Image	0.604	0.240	0.233	0.257
	Multi	0.354	0.186	0.223	0.250
Qwen3-VL 30B					
	Text	0.623	0.264	0.252	0.278
	Image	0.627	0.243	0.242	0.261
	Multi	0.646	0.268	0.267	0.281