# CROSS-LINGUAL MULTIMODAL RETRIEVAL-AUGMENTED GENERATION FOR OPEN QUESTION ANSWERING IN TAMIL AND YORUBA

**Anonymous authors** 

Paper under double-blind review

#### **ABSTRACT**

As large language models (LLMs) with retrieval-augmented generation (RAG) gain traction in multimodal knowledge-base question answering (KBQA), concerns about their transfer to low-resource languages (LRLs) remain unaddressed. We introduce **LR-MMQA**<sup>1</sup>, a benchmark assessing multimodal cross-lingual retrieval and reasoning under the challenges of LRLs. Using a state-of-the-art LLM, we translated the hardest questions from WebQA and MultimodalQA, creating a dataset that stresses cross-evidence aggregation and multi-hop inference. We also introduce XM-RAG, a cross-lingual multimodal RAG pipeline optimized for LRLs, which achieves 38.1 answer accuracy overall, over 6.3 points higher than the next best baseline. Our findings expose significant biases and discrepancies in existing systems, with LR-MMQA highlighting specific failure points. Notably, XM-RAG's performance on LR-MMQA is far below top models on English datasets (WebQA: 64.4, MultimodalQA: 73.48 answer accuracy), demonstrating that current methods still fail at complex, real-world tasks in LRLs. By releasing LR-MMQA and XM-RAG, we provide a resource to evaluate and address these gaps and guide progress toward equitable multimodal KBQA.

#### 1 Introduction

In recent years, Large Language Models (LLMs) have made significant strides in Knowledge Base Question Answering (KBQA) through Retrieval Augmented Generation (RAG) Lewis et al. (2020); Xu et al. (2024); Luo et al. (2024), a paradigm that increasingly leverages multimodal retrieval from vast corpora to demonstrate improved accuracy over text-only methods Suri et al. (2025); Chen et al. (2022); Ling et al. (2025); Yan & Xie (2024). Despite these achievements, retrieval models still struggle to answer knowledge-based questions accurately, fluently, and completely in low-resource languages (LRL) due to limited training data and a lack of high-quality retrieval content in these languages Qi et al. (2025); Rogoz & Lupaşcu (2025). Various methods, such as translate-then-retrieve, have been developed to address this problem Asai et al. (2021) and have been further enhanced by using a multilingual encoder to embed the query in a multilingual semantic space to be used for retrieval from a high-resource language (HRL) corpora Asai et al. (2022). This shifts reliance from the inadequate knowledge present in LRLs to the more comprehensive knowledge in HRLs. This method has been recently augmented via the addition of an image encoder and multimodal retrieval framework, expanding the scope of questions that can be answered correctly Li & Ke (2025).

However, this solution for multimodal KBQA in LRLs is an extremely basic RAG pipeline that underperforms compared to the state-of-the-art seen in high-resource language (HRL) systems Mei et al. (2025). State-of-the-art unimodal frameworks for LRLs exist, but they lack the crucial multimodal processing needed to accurately reflect human communication and information understanding Baltruaitis et al. (2019). While HRLs have advanced significantly in multimodal KBQA, LRL progress is years behind Rogoz & Lupaşcu (2025). Furthermore, evaluation of multimodal retrieval for open KBQA in LRLs is impossible, as there are no datasets with LRL questions that require multimodal understanding and retrieval for their answers Rogoz & Lupaşcu (2025). This lack of a

<sup>&</sup>lt;sup>1</sup>You can find the dataset here: https://huggingface.co/datasets/anonymous132145/LR-MMQA

benchmark makes it impossible to identify and address the shortcomings in current models, thereby perpetuating the performance gap between languages and preventing significant progress.

To enable comprehensive evaluation and critical advances in this area, we introduce LR-MMQA, the first multimodal, cross-lingual open KBQA benchmark for LRLs, featuring 718 questions in Yoruba and Tamil, with ground-truth documents in english. This dataset is designed to require multimodal query understanding as well as multimodal and cross-lingual retrieval for complete answers, with all translations validated by native speakers to ensure accuracy.

We also propose XM-RAG, a novel multimodal RAG baseline. XM-RAG is designed to enable accurate and grounded KBQA for LRLs by directly encoding LRL queries and employing a cross-lingual, multimodal retrieval mechanism from a high-resource knowledge base. The retrieved evidence is reranked using a state-of-the-art learned reranker and then summarized and fused via a refinement and fusion layer. The fused multimodal evidence is then used to generate high-quality answers in the user's original language.

Overall, XM-RAG significantly outperforms existing baselines on LR-MMQA in both accuracy and F1 without fine-tuning. By introducing this framework and benchmark, we aim to enable complete, accurate, and accessible open KBQA across languages in a lightweight and modular fashion. **Our main contributions are:** 

- LR-MMQA, the first LRL KBQA benchmark requiring multimodal query understanding
  and multilingual multimodal retrieval from Tamil and Yoruba. LR-MMQA enables a finer
  analysis of RAG models in low-resource settings, revealing significant weaknesses in crosslingual retrieval, multi-hop reasoning, and answer synthesis, and guiding progress toward
  equitable QA.
- XM-RAG, a multimodal RAG baseline designed for accurate, fluent, and grounded Knowledge Base Question Answering in low-resource languages, leveraging cross-lingual multimodal retrieval from high-resource multimodal knowledge bases.
- We show that XM-RAG significantly improves performance in terms of both accuracy and retrieval for KBQA in both LRLs.

#### 2 Related works

Multimodal Retrieval-Augmented Generation for Knowledge Base Question Answering Despite the many advances of LLMs with multimodal RAG, they still struggle to use external knowledge and unseen data, both of which are necessary for KBQA Zhang et al. (2024). RAG addresses this issue by retrieving external evidence from a corpus of knowledge, in turn increasing accuracy and grounding of model responses to knowledge-base questions Lewis et al. (2020). These results have seen further improvements due to multimodal retrieval. Methods such as Multimodal Multihop, a methodology used to gather data from multiple sources to formulate an answer, show evidence of these promising results when built upon baseline models Yarabelly (2025). Multimodal retrieval allows for the vast multimodal evidence to be leveraged to answer questions that cannot be fully answered with only text. MuRAG Chen et al. (2022) does this by treating images as visual tokens. RA-BLIP Ding et al. (2024) projects retrieved text and images into a shared space before fusion. No matter how these methods treat images, they only retrieve content from the language of the query, meaning the quality of the answers is dependent on the quantity and quality of retrievable evidence present, both of which are lacking in low-resource languages.

Cross-Lingual Retrieval Work has shown that retrieval models struggle in LRLs due to lack of high-quality retrieval content Qi et al. (2025); Rogoz & Lupaşcu (2025). A commonly explored solution to this problem has been a pipeline in which the original LRL query is used to retrieve content from a high-resource knowledge base (KB). That content is used to generate an answer in the original LRL. Quite a lot of work has been done exploring this solution. The most common approach to this is the translate-then-retrieve pipeline, as seen in XOR-RETRIEVE Asai et al. (2021). This approach translates the query to a high resource language, uses this embedding to retrieve text evidence, and feeds it into a multilingual pre-trained model to generate an answer in the LRL. A substantial improvement on this approach can be seen in CORA, where a multilingual embedding of the query is used, reducing errors that arise with machine translation Asai et al. (2022). This direction, leveraging multilingual models like BERT for direct cross-lingual information retrieval, has

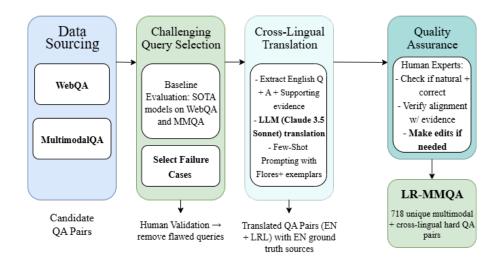


Figure 1: Flowchart of LR-MMQA creation from WebQA and MultimodalQA using Claude 3.5 Sonnet.

been explored in various works, demonstrating effectiveness in matching queries across languages Jiang et al. (2020). The embedded query is then fed into a pre-trained retrieval algorithm and the retrieved evidence is fed into a multilingual auto-regressive generation model to produce an answer. However, all these works only take in, reason on, and retrieve text, significantly limiting the type of questions they can answer accurately.

Multimodal Reasoning and Retrieval Knowledge Base Question Answering Benchmarks Many KBQA benchmarks contain or are entirely composed of questions that require retrieval and understanding across data modalities Chang et al. (2022); Talmor et al. (2021); Marino et al. (2019). However, these datasets solely contain questions and answers in high-resource languages (HRLs). On the other hand, there are KBQA benchmarks that contain questions and answers for LRLs, but only provide the dataset in a single modality, normally text Sawczyn et al. (2024); Rohera et al. (2024); Longpre et al. (2020). Some papers provide publicly accessible benchmark VQA datasets for diverse LRLs, but it should be noted that VQA (answer from the picture) is a field that has been much further explored than the questions in LR-MMQA (answers require multi-hop retrieval). While evaluation of VQA in low-resource languages has been studied Nguyen et al. (2023); Kim et al. (2024), evaluation of multi-hop multimodal reasoning in these languages remains largely unexplored, which is the focus of LR-MMQA.

#### 3 DATASET

#### 3.1 Dataset Overview

**LR-MMQA** is a benchmark designed to rigorously evaluate RAG systems on multimodal and cross-lingual understanding and retrieval for low-resource open KBQA. Curated from two high-resource datasets, it comprises 718 unique multimodal questions specifically selected for their difficulty, representing cases where state-of-the-art models currently fail. Questions were first translated from Standard American English (SAE) into Tamil and Yoruba using an LLM, then post-edited by native-speaker volunteers to ensure fluency and correctness. See Figure 1 for a visual overview of LR-MMQA creation. The questions require multilingual reasoning and retrieval, as the ground-truth documents are not in the query's language. This design simulates a real-world QA environment where comprehensive sources are often unavailable in low-resource languages.

#### 3.2 Data Collection and Preparation

**Data Sourcing** Our dataset is derived from the WebQA dataset and a subset of the MultimodalQA dataset, both of which are standardized collections of open-domain knowledge-seeking queries.

Specifically, these datasets contain queries that require models to retrieve and reason over images, text, or both. Using WebQA and MultimodalQA ensured our analysis was based on authentic, knowledge-seeking questions, enabling question-answering that feels natural and relevant, even when extended to low-resource settings.

Challenging Query Selection To achieve rigorous evaluation, we pre-filtered by running a few existing baselines on WebQA and MultimodalQA. We ran SKURG Yang et al. (2023) and RAMQA Bai et al. (2025) on both datasets, selecting these two publicly available multimodal retrieval frameworks as they represent the highest performing models on the respective datasets. We run these models to establish a higher performance ceiling and identify questions that remain challenging even for state-of-the-art systems, thereby defining a more robust set of "hard" examples. We then selected questions based on their "failure" status between the two models, where a failure represents an average WebQA QA score (composite metric of question accuracy and BARTScore) or MultimodalQA accuracy score of 0. Humans then validated each of these questions to ensure that they contained no errors that made them unanswerable, such as the answer no longer relating to the question. All flawless questions were selected for translation. This comprehensive inclusion directly captures every instance where state-of-the-art systems demonstrably fail, aligning precisely with the benchmark's objective. Table 1 shows a specific breakdown of the selected queries.

Origin	Text	Image	Image + Text
WebQA	106	490	0
MultimodalQA	41	67	14
Total	147	557	14

Table 1: Origin and required retrieval modalities for the selected question-answer pairs. (**Total unique questions: 718**)

#### 3.3 Cross-Lingual Translation and Annotation

**Language Selection** The question-answer pairs in LR-MMQA are in the two low-resource languages (LRLs) of Tamil and Yoruba. Other than being the languages that the team members speak fluently, Tamil and Yoruba are highly semantically diverse, enabling a broad and representative evaluation of cross-lingual capabilities and understanding.

**Translation Protocol** For each selected sample in the English benchmarks, we extract the English Question and English gold answer, along with any supporting materials (i.e., if the dataset includes questions with images). To translate QA pairs from SAE to each of the two LRLs, we employed a few-shot prompting strategy Brown et al. (2020) informed by examples from FLORES+, a highquality dataset containing parallel examples of human translated sentences across languages, including Tamil and Yoruba. Prior work has shown that exemplar-based prompting improves multilingual translation quality in LLMs Lin et al. (2022). We used three exemplar translations from FLORES+ per language. Utilizing Claude 3.5 Sonnet, the LLM was prompted to rewrite the QA pairs into Tamil and Yoruba, informed by the examples. This approach ensures that translations maintain linguistic authenticity and respect the semantic nuances of the target language. Detailed examples of these prompts can be found in Appendix A. Claude 3.5 Sonnet is used because previous work has shown that it has remarkable resource efficiency and outperforms state of the art neural machine translation (NMT) on translation tasks Enis & Hopkins (2024). Furthermore, past work has also shown that LLM translation is superior to NMT in terms of how closely it resembles human translation Sizov et al. (2024), an important quality as questions and answers must appear genuine to truly gauge a model's ability to answer questions in the LRL. Samples of translated queries can be found in the Appendix B.

Structured Entries After translation, the data points for the benchmark were prepared. Every data point in LR-MMQA can be expressed as the tuple  $(Q_{EN},Q_{LRL},E_{MM},A_{EN},A_{LRL},S_{GD})$ , where  $Q_{LRL}$  is the low-resource language question with its parallel English translation  $Q_{EN},E_{MM}$  represents the supporting evidence (images),  $A_{LRL}$  is the low-resource language gold answer with its parallel English translation  $A_{EN}$ , and  $S_{GD}$  are the ground truth documents needed to accurately answer the questions. There is one data point for each of the two languages, meaning LR-MMQA is comprised of 1436 of these data points. A sample data point can be found in Appendix C.

### 3.4 QUALITY ASSURANCE AND VALIDATION

To ensure no errors persist in the dataset, human bilingual experts verified that the LRL QA pairs accurately relate to the supporting evidence, as translation can subtly alter word meanings, potentially rendering a question unanswerable by the provided context or misaligning it with the gold answer. If they noticed any errors in translation, they made necessary corrections.

Additionally, we conducted systematic translation quality evaluation on a representative sample (F). Two native speakers independently assessed 150 translations per language using 10-point Likert scales for adequacy and fluency, achieving substantial inter-annotator agreement (= 0.76 overall) and high quality scores (8.2 adequacy, 8.0 fluency), confirming the effectiveness of our LLM-assisted translation.

### 4 XM-RAG

#### 4.1 OVERVIEW

XM-RAG is a cross-lingual, cross-modal RAG pipeline that performs knowledge-based question answer (KBQA) tasks queried in LRLs by retrieving text and image evidence from high-resource corpora and generating answers directly in the users language. The system follows a modular design to ensure scalability for low resource settings via the use of using off-the-shelf encoders, and utilizing the sufficient ability of multilingual multimodal retrieval to improve KBQA accuracy in unseen LRLs.

#### 4.2 INPUT PROCESSING

Given an LRL question q (optionally with image input), we first run lightweight language identification via FastText to attach a language tag (e.g., <am>, <yo>, <ta>, <zh>, <es>, <en>). While not necessary for evaluation on LR-MMQA, this step promotes the use of XM-RAG in other languages or outside the current testing environment. The question text is encoded without translation using M-CLIP to obtain a unit-normalized query vector  $\mathbf{q}_{\text{text}}$  in a shared embedding space. For input images, we compute visual embedding  $\mathbf{q}_{\text{img}}$  using the same backbone:

$$\mathbf{q}_{\text{text}} = \text{M-CLIP}_{\text{text}}(q), \quad \mathbf{q}_{\text{img}} = \text{M-CLIP}_{\text{vision}}(I_{\text{in}})$$

#### 4.3 Cross-modal Retrieval

We maintain separate FAISS indices for HRL texts (benchmark texts, Wikipedia passages, web snippets) and images. Given  $\mathbf{q}_{\text{text}}$  and  $\mathbf{q}_{\text{img}}$ , we perform k-NN search to obtain:

- Top-K text candidates:  $\{(d_i, s_i^{\text{text}})\}$
- Top-K image candidates:  $\{(v_i, s_i^{img})\}$

using IVF/Flat indexing for scalable retrieval. For evaluation on LRL-MM-QA, HRL ground-truth documents are derived from MultimodalQA and WebQA sources.

#### 4.4 Cross-modal Reranking

We compute cross-modal similarity scores aggregating textual and visual evidence:

$$S(d) = \alpha \cdot s^{\text{text}}(d) + \beta \cdot \max_{v \in \mathcal{N}(d)} s^{\text{img}}(v) + \gamma \cdot \phi(d)$$
 (1)

where  $\mathcal{N}(d)$  are images co-occurring with passage d, and  $\phi(\cdot)$  is a lightweight heuristic (e.g., language match to target tag or answer-type cues). Parameters  $\alpha, \beta, \gamma$  are fixed in our baseline; an optional MLP reranker can be plugged in.

#### 4.5 Multimodal Fusion and Evidence Compression

We perform late fusion by selecting top-n passages and top-m images, compressing them into answer-oriented context:

• BLIP-2 generates bilingual (LRL-tagged) two-sentence visual summaries

272 273 • HRL passages are truncated via rare-caption filtering and deduplication

274

This yields compact, salient context while maintaining lightweight operation.

275 276

#### 4.6 MULTILINGUAL ANSWER GENERATION

277 278 The fused context and original LRL question are fed into a multilingual seq2seq generator (mT5) with appended language tag:

279 280

$$Answer = mT5 \left( \underbrace{[q_{LRL}]}_{tagged input} \oplus evidence_{fused} \right)$$

281 282

The language tag controls output language and prevents spurious HRL translation. The generator attends to both textual evidence and BLIP-2 visual summaries.

283 284 285

# 4.7 Training and Inference Configuration

286 287

XM-RAG operates zero-shot with off-the-shelf components:

288

· All embeddings unit-normalized

289

• FAISS IVF with  $n_{\text{probe}}$  tuned per corpus size

291 292

Strict token caps for latency/memory control

293

· No task-specific fine-tuning

# RESULTS AND ANALYSIS

295 296

297 298

#### 5.1 METRICS

299 300 To assess model performance on LR-MMQA, we used a combination of quantitative and qualitative metrics:

301 302

• Retrieval Metrics: Precision, recall, and F1 scores were used to evaluate the quality of retrieved documents

303 304 305

· Accuracy: Token overlap scores were used to measure the accuracy and coherence of generated outputs.

306 307 308

309

310

LR-MMQA utilizes a Full-Wiki evaluation, meaning the model must retrieve the correct answer from the entire Wikipedia corpus. This differs from datasets like WebQA and MultimodalQA that provide a small set of candidate documents, thus eliminating the need for a separate retrieval step. In the LR-MMQA setting, there are no distractors to contend with, as the primary task is to find and extract the correct information from a full knowledge base in a different language than the query, simulating a real-world scenario.

311 312 313

#### 5.2 Baselines

314 315

316

Given the absence of a publicly available true multimodal, multilingual baseline model that directly addresses all aspects of our research, we selected a set of representative baselines. These models were chosen to provide a comprehensive comparison against our proposed approach by evaluating its performance across different modalities and languages. This multifaceted evaluation serves to expose the specific gaps in current models that our approach aims to address.

Text-Only Cross-Lingual RAG Baseline This baseline represents the current state-of-the-art in Text-Only RAG, but also serves to highlight its core limitation: the inability to process non-textual information. Although it demonstrates advanced cross-lingual generation by retrieving from a highresource language and answering in a target language, its text-only nature makes it fundamentally inadequate for the real-life questions posed in the LR-MMQA benchmark, which requires a model

321 322 323

		Text				Imag	e		I	mage+	Text		All	
Model	QA-Acc	Prec	Rec	F1	QA-Acc	Prec	Rec	F1	QA-Acc	Prec	Rec	F1	QA-Acc	F1
Text-Only RAG	20.8	17.5	16.8	17.1	0.4	0.5	0.4	0.4	0.2	0.8	0.6	0.7	8.3	6.2
Monolingual RAG	22.1	18.6	17.6	18.1	12.3	10.8	11.2	10.9	12.6	10.1	9.7	9.9	18.2	15.0
GPT-4o	18.4	N/A	N/A	N/A	16.4	N/A	N/A	N/A	8.1	N/A	N/A	N/A	16.6	N/A
RAGVL + MT	30.1	34.1	28.7	31.2	32.1	39.9	33.2	36.2	36.2	42.5	34.5	38.1	31.8	35.2
XM-RAG	36.7	37.1	85.6	51.8	38.3	44.5	57.1	50.0	42.9	44.2	67.9	53.6	38.1	50.5

Table 2: Performance of models on LR-MMQA by query modality.

to reason across different modalities. Its inclusion demonstrates that even sophisticated Text-Only models fail when faced with the real-world complexity of multimodal tasks.

Multimodal Monolingual RAG Baseline This baseline tests multilingual limitations of standard multimodal RAG pipelines trained in high-resource languages. It handles multiple modalities but is monolingual. Queries are machine-translated to SAE, and answers translated back to the target language. Comparing against this baseline highlights that existing multimodal systems cannot transfer knowledge across languages effectively and that machine translation is a poor substitute. This demonstrates that current models are inadequate for a task that requires multimodal and multilingual capabilities, which is essential for success when evidence in the query language is limited.

Multimodal RAG Baseline RagVL Chen et al. (2024) RagVL is the SOTA on WebQA and MultimodalQA (Full Wiki), combining vision-language modeling with retrieval via knowledge-enhanced reranking and noise-injected training. We adapt it with the same MT process as the multimodal monolingual baseline, showing that even the strongest multimodal RAG systems fail to generalize across languages.

**SOTA Commercial MLLM - GPT-40** We include GPT-40 as a strong, non-retrieval-augmented baseline representing peak multimodal understanding. Despite its capabilities, it underperforms in a real-world, retrieval-intensive, multilingual setting, illustrating that even advanced commercial models cannot yet handle the challenges of LR-MMQA or difficult knowledge-seeking queries in low-resource languages.

#### 5.3 Main Results

Retrieval Quality Analysis We evaluate retrieval baselines on our benchmark, which consists of WebQA and MMQA questions where state-of-the-art systems fail, so overall performance is predictably low. For context, full-wiki SOTA F1 on the original WebQA and MMQA is 77.64 and 98.9Chen et al. (2024), respectively; despite XM-RAGs strong gains below, there remains a large absolute gap. As shown in Table 2, on text questions XM-RAG recall is extremely high (85.6) but paired with the lowest precision (37.1), producing the second lowest F1 across modalities (51.8). This reflects a consistent pattern where XM-RAG recall exceeds precision, stemming from the dense retriever surfacing many cross-lingual candidates that the reranker cannot fully filter. Even so, XM-RAGs and RagVL's higher precision compared with text-only (F1 17.1) and monolingual baselines (18.1) shows the benefit of both baselines reranker and XM-RAG's multilingual retriever, which better prune off-topic candidates and promote more effective search across languages. An example of XM-RAGs successful multi-hop retrieval, compared to baseline failures, appears in Appendix E.

The monolingual baseline outperforms the text-only baseline primarily on image and image+text questions, which make up a large portion of the dataset, while on text-only questions their metrics are close. On both image and image+text queries, XM-RAG substantially outperforms all baselines across metrics (e.g. F1: 50.0 and 53.6), yielding an overall F1 of 50.5, which is 15.3 points higher than the next best baseline (35.2). Although LR-MMQA is built from failure cases, the gap to original WebQA/MMQA highlights persistent limits in multilingual encoding. Addressing the precisionrecall imbalance in retrieval, alongside improving cross-lingual representations, will be essential for future systems to close this disparity.

Answer Quality Answer accuracy shows a clear gap between the evaluated systems. GPT-40 outperforms the text-only baseline by 8.5 points, likely because some factual knowledge appears in its training data. The text-only baseline performs near zero on image and multimodal questions (0.4 and 0.2), producing unsupported answers that lower overall accuracy. The monolingual base-

line achieves 9.5 points higher overall accuracy than the text-only baseline by incorporating both text and image evidence. XM-RAG reaches 38.1 accuracy, 6.3 points higher than RagVL and 20.3 points higher than any other baseline. These gains are largely due to its multilingual query encoding and cross-modal fusion, which together retrieve and compress more relevant evidence than the baselines aided by MT.

Although retrieval is often successful, accuracy remains low because even XM-RAG struggles on questions requiring reasoning over multiple highly specific pieces of information (28.7 % of all XM-RAG errors; see Appendix G). Text queries often fail when answers depend on comparing fine-grained details such as dates across documents(45.2%), while image queries often fail when comparing attributes like object colors across multiple images (41.3%). These challenges are amplified in cross-lingual retrieval from low-resource languages, where even small translation errors can prevent precise ground-truth items from being surfaced (31.2% of all failures). This problem is exacerbated by Yoruba questions, where tonal information can be lost in text-only embeddings. For comparison, the full-wiki SOTA accuracy on WebQA and MMQA is 64.40 and 73.48 Chen et al. (2024), leaving large gaps of 26.3 and 35.4 points. Thus, while XM-RAG sets a new state of the art for LR-MMQA, current systems remain limited in retrieving and reasoning over evidence. Future progress depends on better integration of retrieved context and addressing the additional difficulties posed by low-resource languages and cross-lingual retrieval. Failure examples can be found in Appendix H, I, J

Qualitative Analysis on Generated Responses Our qualitative analysis of GPT-40 reveals a consistent geographical bias: queries in Yoruba often yield Nigeria-centered responses, while Tamil queries default to Indian contexts. This behavior highlights the geographical bias inherent in training data, a critical shortcoming in LLMs. Since the majority of training data for LRLs is inherently concentrated in a specific region, the model forms a strong statistical association between the language and its dominant culture.

Such bias is not unique to commercial LLMs, as a monolingual RAG pipeline would face the exact same issue, encountering limited data that only contains information about the country where the language of the query is spoken. On the contrary, The multilingual text-only baseline and XM-RAG did not display the language-culture bias on the same set of questions. Their ability to retrieve multilingually from a wider HRL corpus that transcends single-country limitations allowed them to provide answers that were not confined to a single country or cultural frame, leading to more accurate responses. Reference Appendix D for an example of a question and generated responses falling under this description.

# 5.4 ABLATION STUDY

Model	QA-Acc	F1
XM-RAG (full)	38.1	50.5
w/o Cross-Encoder Reranker	37.2	49.4

Table 3: Ablation study comparing XM-RAG with and without the cross-encoder reranker. The best results for each metric are highlighted in bold.

As shown in Table 3, XM-RAG without the cross-encoder reranker achieves an accuracy of 37.2 and a retrieval F1 of only 40.7, a drop of 2.2% compared to the F1 of the full XM-RAG pipeline. As seen in Figure 2, this F1 performance drop occurs consistently in all modalities, highlighting the critical role of the learned reranker in bridging retrieval and reasoning.

The drop in F1 occurs because the reranker contributes to more precise answer grounding: without it, the system tends to select passages or segments that are semantically related but not directly relevant to the query. This results in noisier context, weaker alignment between evidence and the question, and ultimately a degradation in precision, which can be seen in all modalities in Figure 2. Since F1 directly combines both recall and precision, even moderate precision losses cause a decline in overall F1.

Accuracy also declines under this ablation because incorrect or noisy contexts lead the generator to produce answers that are either partially correct or completely off-target. The reranker ensures that

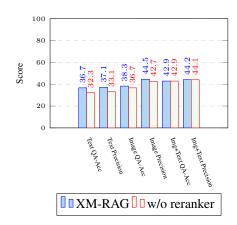


Figure 2: Ablation on reranker: accuracy and precision scores for different modalities.

retrieved evidence is both semantically rich and directly relevant to the query. Without this step, the model's reasoning chain is built on lower-quality foundations, making accurate prediction much less likely.

#### 6 CONCLUSION

In this paper, we introduce the first multimodal KBQA dataset for LRLs, LR-MMQA, as well as XM-RAG, a state-of-the-art baseline model for multimodal KBQA in LRLs. LR-MMQA is a benchmark for Tamil and Yoruba that utilizes questions and images from WebQA and MultimodalQA, two open-domain question-answer-answer datasets. We then translated select queries from both datasets into Tamil and Yoruba, containing 718 unique question-answer pairs for each language. We evaluate our baseline model XM-RAG and compare it with existing open-domain benchmark models. Through XM-RAG's unique combination of features, we achieve SOTA metrics across all modalities compared to other baseline models, which we can attribute to XM-RAG's ability to handle both multilingual and multimodal data.

#### 7 LIMITATIONS

Due to the absence of a multilingual multimodal RAG model for Tamil and Yoruba for KBQA, there may be limited comparison of XM-RAG to other models. It should be noted that LR-MMQA is a relatively small dataset in comparison to WebQA or MultimodalQA, with a particular emphasis on image questions. In the future, more questions and answers should be created with ground-truth text documents to combat this issue and allow further evaluation. Furthermore, LR-MMQA can also be improved with the inclusion of QA pairs in tonal or highly agglutinative languages to create a more inclusive benchmark and better assess a model's performance across diverse languages. Finally, HRL knowledge bases may not fully reflect low-resource scenarios.

#### REPRODUCIBILITY STATEMENT

The code used in this paper can be found here. The steps to reproduce the results are:

- 1. Clone the repository.
- 2. Install dependencies using pip install -r requirements.txt
- Download the LR-MMQA benchmark from https://huggingface.co/datasets/ anonymous132145/LR-MMQA.
- 4. Download the supporting context images from here.
- 5. Follow all instructions in README.md.

After running the code as outlined in the repository, you should be able to reproduce the evaluation metrics reported in Table 2.

# LLM STATEMENT

LLMs were used in this paper to aid and polish writing and experimental code.

### REFERENCES

- Akari Asai, Jungo Kasai, Jonathan Clark, Kenton Lee, Eunsol Choi, and Hannaneh Hajishirzi. XOR QA: Cross-lingual open-retrieval question answering. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL HLT 2021)*, pp. 547–564. Association for Computational Linguistics, 2021. URL https://aclanthology.org/2021.naacl-main.46/.
- Akari Asai, Xinyan Yu, Jungo Kasai, and Hannaneh Hajishirzi. One question answering model for many languages with cross-lingual dense passage retrieval. In *International Conference on Learning Representations (ICLR)*, 2022.
- Yang Bai, Christan Earl Grant, and Daisy Zhe Wang. Ramqa: A unified framework for retrieval-augmented multi-modal question answering, 2025.
- Tadas Baltruaitis, Chaitanya Ahuja, and Louis-Philippe Morency. Multimodal machine learning: A survey and taxonomy. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 41(2): 423–443, 2019. doi: 10.1109/TPAMI.2018.2798607.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners, 2020. URL https://arxiv.org/abs/2005.14165.
- Yingshan Chang, Mridu Narang, Hisami Suzuki, Guihong Cao, Jianfeng Gao, and Yonatan Bisk. WebQA: Multihop and multimodal QA. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 16495–16504. IEEE, 2022.
- Wenhu Chen, Hexiang Hu, Xi Chen, Pat Verga, and William W Cohen. MuRAG: Multimodal retrieval-augmented generator for open question answering over images and text. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 5558–5570. Association for Computational Linguistics, 2022.
- Zhanpeng Chen, Chengjin Xu, Yiyan Qi, and Jian Guo. Mllm is a strong reranker: Advancing multimodal retrieval-augmented generation via knowledge-enhanced reranking and noise-injected training. *arXiv preprint arXiv:2407.21439*, 2024.
- Muhe Ding, Yang Ma, Pengda Qin, Jianlong Wu, Yuhong Li, and Liqiang Nie. RA-BLIP: Multi-modal adaptive retrieval-augmented bootstrapping language-image pre-training. *arXiv* preprint *arXiv*:2410.14154, 2024.
- Maxim Enis and Mark Hopkins. From Ilm to nmt: Advancing low-resource machine translation with claude, 2024. URL https://arxiv.org/abs/2404.13813.
- Zhuolin Jiang, Amro El-Jaroudi, William Hartmann, Damianos Karakos, and Lingjun Zhao. Crosslingual information retrieval with bert. In *Proceedings of the workshop on Cross-Language Search and Summarization of Text and Speech (CLSSTS2020)*, pp. 26–31, 2020.
- Minjun Kim, Seungwoo Song, Youhan Lee, Haneol Jang, and Kyungtae Lim. Bok-vqa: Bilingual outside knowledge-based visual question answering via graph representation pretraining. In *AAAI Conference on Artificial Intelligence*, 2024. URL https://arxiv.org/pdf/2401.06443.

- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Kuttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe van der Wal. Retrieval-augmented generation for knowledge-intensive nlp tasks. *arXiv preprint arXiv:2005.11401*, 2020.
  - Zichao Li and Zong Ke. Cross-modal augmentation for low-resource language understanding and generation. In *Proceedings of the 1st Workshop on Multimodal Augmented Generation via Multimodal Retrieval*, pp. 90–99, 2025.
  - Xi Victoria Lin, Todor Mihaylov, Mikel Artetxe, Tianlu Wang, Shuohui Chen, Daniel Simig, Myle Ott, Naman Goyal, Shruti Bhosale, Jingfei Du, Ramakanth Pasunuru, Sam Shleifer, Punit Singh Koura, Vishrav Chaudhary, Brian O'Horo, Jeff Wang, Luke Zettlemoyer, Zornitsa Kozareva, Mona Diab, Veselin Stoyanov, and Xian Li. Few-shot learning with multilingual generative language models. *arXiv preprint arXiv:2112.10668*, 2022.
  - Zihan Ling, Zhiyao Guo, Yixuan Huang, Yi An, Shuai Xiao, Jinsong Lan, Xiaoyong Zhu, and Bo Zheng. MMKB-RAG: A multi-modal knowledge-based retrieval-augmented generation framework. *arXiv preprint arXiv:2504.06734*, 2025.
  - S. Longpre, Yi Lu, and Joachim Daiber. Mkqa: A linguistically diverse benchmark for multilingual open domain question answering. *Transactions of the Association for Computational Linguistics*, 9:1389–1406, 2020. URL https://aclanthology.org/2021.tacl-1.82.pdf.
  - Haoran Luo, Haihong E, Zichen Tang, Shiyao Peng, Yikai Guo, Wentai Zhang, Chenghao Ma, Guanting Dong, Meina Song, Wei Lin, Yifan Zhu, and Luu Anh Tuan. ChatKBQA: A generate-then-retrieve framework for knowledge base question answering with fine-tuned large language models. In *Findings of the Association for Computational Linguistics: ACL 2024*, pp. 2039–2056. Association for Computational Linguistics, 2024.
  - Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. Ok-vqa: A visual question answering benchmark requiring external knowledge. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 11209–11218, 2019. URL http://openaccess.thecvf.com/content\_CVPR\_2019/html/Marino\_OK-VQA\_A\_visual\_Question\_Answering\_Benchmark\_Requiring\_External\_Knowledge\_CVPR\_2019\_paper. html.
  - Lang Mei, Siyu Mo, Zhihan Yang, and Chong Chen. A survey of multimodal retrieval-augmented generation, 2025. URL https://arxiv.org/abs/2504.08748.
  - Nghia Hieu Nguyen, Duong T.D. Vo, Kiet Van Nguyen, and Ngan Luu-Thuy Nguyen. Openvivqa: Task, dataset, and multimodal fusion models for visual question answering in vietnamese. *Inf. Fusion*, 100:101868, 2023. URL https://arxiv.org/pdf/2305.04183.
  - Jirui Qi, Raquel Fernandez, and Arianna Bisazza. On the consistency of multilingual context utilization in retrieval-augmented generation. *arXiv preprint arXiv:2504.00597*, 2025.
  - Ana-Cristina Rogoz and Marian Lupaşcu. Large multimodal models for low-resource languages: A survey. *arXiv preprint arXiv:2502.05568*, 2025.
  - Pritika Rohera, Chaitrali Ginimav, Akanksha Salunke, Gayatri Sawant, and Raviraj Joshi. L3cube-indicquest: A benchmark question answering dataset for evaluating knowledge of llms in indic context. *ArXiv*, abs/2409.08706, 2024. URL https://aclanthology.org/2024.paclic-1.93.pdf.
  - Albert Sawczyn, Katsiaryna Viarenich, Konrad Wojtasik, Aleksandra Domogaa, Marcin Oleksy, Maciej Piasecki, and Tomasz Kajdanowicz. Developing pugg for polish: A modern approach to kbqa, mrc, and ir dataset construction. *ArXiv*, abs/2408.02337, 2024. URL https://arxiv.org/pdf/2408.02337.
  - Fedor Sizov, Cristina España-Bonet, Josef Van Genabith, Roy Xie, and Koel Dutta Chowdhury. Analysing translation artifacts: A comparative study of LLMs, NMTs, and human translations. In Barry Haddow, Tom Kocmi, Philipp Koehn, and Christof Monz (eds.), *Proceedings of the Ninth Conference on Machine Translation*, pp. 1183–1199, Miami, Florida, USA, November

- 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.wmt-1.116. URL https://aclanthology.org/2024.wmt-1.116/.
- Manan Suri, Puneet Mathur, Franck Dernoncourt, Kanika Gowswami, Ryan A Rossi, and Dinesh Manocha. VisDoM: Multi-document QA with visually rich elements using multimodal retrieval-augmented generation. In *Proceedings of the 2025 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL HLT)*, pp. 6088–6109. Association for Computational Linguistics, 2025.
- Alon Talmor, Ori Yoran, Amnon Catav, Dan Lahav, Yizhong Wang, Akari Asai, Gabriel Ilharco, Hannaneh Hajishirzi, and Jonathan Berant. MultiModalQA: Complex question answering over text, tables and images. In *International Conference on Learning Representations (ICLR)*, 2021.
- Zhentao Xu, Mark Jerome Cruz, Matthew Guevara, Tie Wang, Manasi Deshpande, Xiaofeng Wang, and Zheng Li. Retrieval-augmented generation with knowledge graphs for customer service question answering. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '24, pp. 29052909, New York, NY, USA, 2024. Association for Computing Machinery. ISBN 9798400704314. doi: 10.1145/3626772.3661370. URL https://doi.org/10.1145/3626772.3661370.
- Yibin Yan and Weidi Xie. EchoSight: Advancing visual-language models with wiki knowledge. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pp. 1538–1551. Association for Computational Linguistics, 2024.
- Qian Yang, Qian Chen, Wen Wang, Baotian Hu, and Min Zhang. Enhancing multi-modal multi-hop question answering via structured knowledge and unified retrieval-generation. In *Proceedings of the 31st ACM International Conference on Multimedia*, MM 23, pp. 52235234. ACM, October 2023. doi: 10.1145/3581783.3611964. URL http://dx.doi.org/10.1145/3581783.3611964.
- Mittal Yarabelly. Multimodal multihop source retrieval for web question answering. *arXiv preprint arXiv:2501.04173*, 2025.
- Tao Zhang, Ziqi Zhang, Zongyang Ma, Yuxin Chen, Zhongang Qi, Chunfeng Yuan, Bing Li, Junfu Pu, Yuxuan Zhao, Zehua Xie, Jin Ma, Ying Shan, and Weiming Hu. mr<sup>2</sup>ag: Multimodal retrieval-reflection-augmented generation for knowledge-based vqa. *arXiv preprint arXiv:2411.15041*, 2024. URL https://arxiv.org/abs/2411.15041.

# A TRANSLATION PROMPT

648

649 650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669 670

671

672

673

674

675

676 677

678

679

680 681

#### Yoruba Few-Shot Examples **Tamil Few-Shot Examples** Example 1: Example 1: English: The police were called to the scene. English: After World War I, a new political landscape emerged in the Middle East. Tamil: kaavalthuraiyukku anda idaththirkku vara azhaippu vidukkappattadhu. Yoruba: Lyìn Ogun Àgbáyé Kìíní, ìèlú tuntun kan farahàn ní Àárín Gbùngbùn Il-Ayé. Example 2: English: The new species of butterfly was discovered Example 2: English: The film "Moonlight" won the Academy in the rainforest. Tamil: oru puthiya pattampuuchi inam mazhai kaatil Award for Best Picture. kandupidikkappattadhu. Yoruba: Fîimù "Moonlight" gba Àmì-y Akádmì fún Fîimù Tó Dára Jù L. Example 3: English: The chef prepared a delicious meal using fresh, Example 3: local ingredients. English: A well-known saying is "the early bird catches Tamil: samayalkaarar puthiya, ulloor porutkalai payanthe worm." paduthi oru suvaiyana unavai thayariththaar. Yoruba: Àà àti ìe tí a m jù ni pé "y àár ní mú kòkòrò."

#### **Translation Prompt:**

You are an expert linguist and translator. Your task is to translate a Question-Answer pair from English to the target language. You must maintain the integrity of the question and ensure the translated answer remains a correct, verbatim excerpt from the translated context.

#### **Instructions:**

- 1. Translate the question to the target language.
- 2. Translate the answer to the target language.
- 3. The translated answer must be a direct, literal substring of the translated text (not paraphrased).
- 4. Maintain the original format and structure.
- 5. Ensure all questions and answers are posed as a native speaker would ask and answer.

#### Your Task:

Source Question: {source\_question} Source Answer: {source\_answer}

Table 4: Few-shot exemplars and translation prompt used for creating LR-MMQA.

# **B** TRANSLATION EXAMPLES

English Question	English Answer	Yoruba Question / Answer	Tamil Question / Answer
Who sings the most songs in the world?	Asha Bhosle	Tani o krin plp jùl ní gbogbo àgbáyé? / Asha Bhosle	ulagil adhika paadalgalai paadiyavar yaar? / Aasha Bhosle
How many colors are in the Point Skyhawks logo?	4	Àwn àw mélòó ló wà nínú àmì Point Sky- hawks? / 4	paayind skaihaaks logo- il etthanai niRangal ul- lana? / 4
Danish Viking, who ruled over parts of Friesland between 841 and 873, was the uncle of a Viking leader who raided the British Isles, West Francia, Frisia, and Lotharingia in the 860s and 870s?	Roricus, Rorichus	Viking Denmark, tí ó j ba lórí apá kan ti Fries- land láàrin dún 841 àti 873 j àbúrò bàbá tàbí ìyá fún olórí Viking kan tí ó klu Erékùù Brítánì, Ìw-òòrùn Francia, Frisia, àti Lotharingia ní grùn- ún dún ksàn-án àti grùn- ún dún kwàá? / Roricus, Rorichus	841 muthal 873 varai freeslandin pagudigalai aatchi seidha danish viking, 860kal matrum 870kalil british theevugal, merku Francia, frisia matrum lotharingia-kolaiyaditha oru viking thalaivarin maamaa yaar? / Rorikus, Rorichus

Table 5: Examples of English questions and answers with Yoruba and Tamil translations.

# C SAMPLE PERTURBED DATA POINT

Field	Content
$Q_{EN}$ (English Question)	If a partial seizure spreads to the cortex, it can result in what type of tonic-clonic seizure?
$Q_{LRL}$ (Tamil Question)	paguthi valippu moolaiyin puranukku paravinaal, adhu endha vagaiyana tonic-clonic valippaga maaralaam?
$A_{EN}$ (English Answer)	Grand mal
$A_{LRL}$ (Tamil Answer)	grand maal
$E_{MM}$ (Supporting Context)	N/A
$S_{GD}$ (Ground Truth Documents)	Seizure types Wikipedia (title only) Generalized tonicclonic seizure Wikipedia (title only) Note: Full URLs and snippets omitted for space.

Table 6: Example data point from LR-MMQA showing English and Tamil QA pairs, supporting context (if applicable), and gold source titles.

# D SAMPLE BIASED GPT OUTPUT

Field	Content
$Q_{EN}$ (English Question)	Who sings the most songs in the world?
$Q_{LRL}$ (Yoruba Question)	Ta ni o krin awn orin jù l ni ayé?
A <sub>GD</sub> (Gold Answer)	Asha Bhosle
A <sub>GPT</sub> ( <b>GPT Answer</b> )	Fela Kuti

Table 7: Example of a biased GPT output where the model incorrectly localized the answer to a Nigerian artist, despite the gold answer being *Asha Bhosle*.

# E SAMPLE XM-RAG OUTPUT WITH RETRIEVED SOURCES

Field	Content
$Q_{EN}$ (English Question)	Which color is found on both the Estonia and Poland Pavilion at Expo 2010?
$Q_{LRL}$ (LRL Question)	Àw wo ni a rí lórí Ilé Ìgbìm Estonia àti Poland ní Ìfihàn Àgbáyé 2010?
$A_{EN}$ (English Answer)	Brown is found on both the Estonia and Poland Pavilion at Expo 2010.
$A_{LRL}$ (LRL Answer)	Àw búráùnì ni a rí lórí Ilé Ìgbìm Estonia àti Poland ní Ìfihàn Àgbáyé 2010.
$A_{XM}$ (XM-RAG Answer)	Àw búráùnì náà ló wà lórí Ilé Ìgbìm Estonia àti Poland ní Ìfihàn Àgbáyé 2010.
$S_{GD}$ (Gold Sources)	Estonia Pavilion at Expo 2010 in Shanghai urges action to save the cities.
	Polish Pavilion at Shanghai World Expo 2010.
$S_{XM}$ (Retrieved Sources)	Estonia Pavilion at Expo 2010 in Shanghai urges action to save the cities.
	• Polish Pavilion at Shanghai World Expo 2010.
	Polish Pavilion / WWA Architects.
	Estonian pavilion for Shanghai EXPO 2010 - Identity.
	• Expo 2010 pavilions.

Table 8: Example XM-RAG successful multi-hop reasoning through answering and retrieval.

# F TRANSLATION QUALITY EVALUATION

Language	Sample Size	Adequacy (1-10)	<b>Fluency (1-10)</b>	Inter-Annotator Agreement
Tamil	150	8.3	8.1	$\kappa = 0.78$
Yoruba	150	8.1	7.9	$\kappa = 0.74$
Overall	300	8.2	8.0	$\kappa = 0.76$

Table 9: Translation quality evaluation results for LR-MMQA dataset. Two native speakers independently assessed translations using 10-point Likert scales for adequacy (semantic correctness) and fluency (naturalness). Inter-annotator agreement measured using Cohen's kappa ( $\kappa$ ).

**Evaluation Protocol:** Two native speakers independently rated each translation on 10-point Likert scales. Disagreements resolved through discussion with a third annotator.

**Rating Scale:** Adequacy: 1 = completely incorrect meaning, 10 = perfect semantic preservation. Fluency: 1 = completely unnatural, 10 = native-like naturalness.

# G FAILURE MODE ANALYSIS

Failure Type	Text-Only	Image-Only	Image+Text	Overall
Cross-lingual Retrieval	32.1%	28.4%	35.7%	31.2%
Visual Understanding	0.0%	41.3%	29.8%	25.6%
Multi-hop Reasoning	45.2%	18.9%	21.4%	28.7%
Answer Generation	22.7%	11.4%	13.1%	14.5%

Table 10: Distribution of failure modes across question types through systematic human categorization of error cases from XM-RAG outputs on LR-MMQA. Authors performed this task.

# H CROSS-LINGUAL RETRIEVAL FAILURE EXAMPLE

Field	Content	
$Q_{EN}$ (English Question)	How many people are in the painting of Sappho and Phaon by Jacques-Louis David?	
$Q_{LRL}$ (LRL Question)	Jaak-luuyi Devid varainda Saappo matrum Paayon oviathil eththanai per irukkiraargal?	
$A_{EN}$ (English Answer)	3 people are in the painting of Sappho and Phaon by Jacques-Louis David.	
$A_{LRL}$ (LRL Answer)	Jaak-luuyi Devid varainda Saappo matrum Paayon oviathil 3 per irukkiraargal.	
A <sub>XM</sub> (XM-RAG Answer)	Jaak-luuyi Devittin oviangal patri thagaval kidaikkavillai.	
$S_{GD}$ (Gold Sources)	• Jacques-Louis David - Sappho and Phaon - WGA6092	
$S_{XM}$ (Retrieved Sources)	<ul> <li>Jacques-Louis David paintings overview</li> <li>French neoclassical art collection</li> <li>David historical paintings</li> <li>18th century French artists</li> <li>*Cross-lingual encoding failed to match "Saappo" with "Sappho"*</li> </ul>	

Table 11: Example XM-RAG cross-lingual retrieval failure due to semantic drift in proper name encoding.

# I VISUAL UNDERSTANDING FAILURE EXAMPLE

Content	
Looking at Zocalo from Torre Latino Americana how many yellow buildings can be seen?	
Ni wiwo Zocalo lati Torre Latino Americana, ile pupa melo ni a le ri?	
Looking at Zocalo from Torre Latino Americana, one yellow building is visible.	
Ni wiwo Zocalo lati Torre Latino Americana, ile pupa kan ni a le ri.	
Ni wíwo Zocalo lati Torre Latino Americana, awon ile pupo ni a le ri.	
Zocalo and surroundings as seen from Torre Lati- noamericana, Mexico City	
<ul> <li>Zocalo and surroundings as seen from Torre Latinoamericana, Mexico City</li> <li>Mexico City aerial views</li> <li>Torre Latinoamericana observation deck</li> <li>Historic center Mexico City</li> <li>Zocalo plaza architecture</li> </ul>	

Table 12: Example XM-RAG visual understanding failure in fine-grained object counting and color identification.

# J MULTI-HOP REASONING FAILURE EXAMPLE

Field	Content		
$Q_{EN}$ (English Question)	Which candidate announced their run for president later; Gary Johnson or Michelle Bachmann?		
$Q_{LRL}$ (LRL Question)	Athipar pathavikkaana thangal pottiyai yaar pinthi arivit- thaar; Kaeri Jaansanaa allathu Mishel Baakmanaa?		
$A_{EN}$ (English Answer)	Michelle Baakmanaa		
$A_{LRL}$ (LRL Answer)	Mishel Baakman		
A <sub>XM</sub> (XM-RAG Answer)	Kaeri Jaansan		
$S_{GD}$ (Gold Sources)	<ul> <li>Gary Johnson 2012 presidential campaign - Wikipedia</li> <li>Michele Bachmann 2012 presidential campaign - Wikipedia</li> </ul>		
$S_{XM}$ (Retrieved Sources)	<ul> <li>Gary Johnson 2012 presidential campaign - Wikipedia</li> <li>Michele Bachmann 2012 presidential campaign - Wikipedia</li> <li>2012 Republican primary candidates</li> <li>Presidential campaign announcements 2011</li> <li>Gary Johnson political career</li> </ul>		

Table 13: Example XM-RAG multi-hop reasoning failure in temporal comparison synthesis across retrieved documents.

# K EXPERIMENTAL SETUP AND COMPUTING RESOURCES

The experiments were conducted using a dedicated GPU cluster for training and inference on large models. Below are the specifications and details:

**GPU Resources:** The main experiments were performed on a GPU cluster equipped with 2x NVIDIA A100 SXM GPUs with 251 GB memory each. These GPUs provided high-throughput tensor core acceleration suitable for challenging multimodal and multilingual KBQA tasks.

**CPU Resources:** The cluster included 16 vCPUs, which were used for data preprocessing, baseline evaluations, and lightweight model inference tasks alongside GPU computations.

**Memory:** The GPU cluster had sufficient system RAM to manage large datasets and multimodal inputs efficiently. The 251 GB GPU memory per card allowed for batch processing and minimized data offloading during model execution.

**Storage:** Experiments utilized high-speed SSD storage on the cluster to handle the 1,436 KBQA examples from the benchmark, including multimodal inputs such as images and structured knowledge representations.

# **Experiment Details:**

• RAMQA and SKURG: Running RAMQA and SKURG on the benchmark took approximately 5 hours for WebQA and 3 hours for MultimodalQA. These times reflect the complexity of reasoning across multiple hops and modalities.

• Baseline Models: Running all other baseline models on the same benchmark is estimated to take an additional 2-4 hours each, considering the relatively small dataset size (1,436 examples) but challenging multimodal multi-hop questions. This estimate accounts for persample inference times, preprocessing overhead, and model loading times on the cluster.

**Total Computing Time:** In total, including running RAMQA, SKURG, and all baseline models, the experiments required roughly 10-15 GPU hours and approximately 20-25 CPU hours on the cluster for preprocessing and supporting tasks. This configuration ensured that all models could be executed efficiently while handling the high memory and computational demands of multimodal KBQA reasoning tasks.