

MVL-SIB: A Massively Multilingual Vision-Language Benchmark for Cross-Modal Topical Matching

Fabian David Schmidt^{1*}, Florian Schneider^{2*}, Chris Biemann², Goran Glavaš¹

¹Center for Artificial Intelligence and Data Science, University of Würzburg, Germany

²Language Technology Group, University of Hamburg, Germany

fabian.schmidt@uni-wuerzburg.de, florian.schneider-1@uni-hamburg.de

Dataset: MVL-SIB

Abstract

Existing multilingual vision-language (VL) benchmarks often only cover a handful of languages. Consequently, evaluations of large vision-language models (LVLMs) predominantly target high-resource languages, underscoring the need for evaluation data for low-resource languages. To address this limitation, we introduce MVL-SIB, a massively multilingual vision-language benchmark that evaluates both cross-modal and text-only topical matching across 205 languages—over 100 more than the most multilingual existing VL benchmarks encompass. We then benchmark a range of open-weight LVLMs together with GPT-4o(-mini) on MVL-SIB. Our results reveal that LVLMs struggle in cross-modal topic matching in lower-resource languages, performing no better than chance on languages like N’Koo. Our analysis further reveals that VL support in LVLMs declines disproportionately relative to textual support for lower-resource languages, as evidenced by comparison of cross-modal and text-only topical matching performance. We further observe that open-weight LVLMs do not benefit from representing a topic with more than one image, suggesting that these models are not yet fully effective at handling multi-image tasks. By correlating performance on MVL-SIB with other multilingual VL benchmarks, we highlight that MVL-SIB serves as a comprehensive probe of multilingual VL understanding in LVLMs.¹

1 Introduction

Large Vision-Language Models (LVLMs) extend Large Language Models (LLMs) to take images as inputs, leveraging their advanced language capabilities for vision-language (VL) tasks like image captioning and visual question answering (VQA). However, LVLMs are typically trained mainly on

*Equal contribution.

¹Code: <https://github.com/floschne/mvl-sib>

Images-To-Sentences (I2S)

Which sentence best matches the topic of the images? The images and the sentences each belong to one of the following topics: “entertainment”, “geography”, “health”, “politics”, “science and technology”, “sports”, or “travel”. Choose one sentence from A, B, C, or D. Output only a single letter!

Images



Sentences

- A. „Maroochydore führte am Ende die Rangfolge an, mit sechs Punkten Vorsprung vor Noosa als Zweiter.“
- B. „Es wurden keine schweren Verletzungen gemeldet, jedoch mussten mindestens fünf der zur Zeit der Explosion Anwesenden aufgrund von Schocksymptomen behandelt werden.“
- C. „Finnland ist ein großartiges Reiseziel für Bootstouren. Das „Land der tausend Seen“ hat auch Tausende von Inseln – in den Seen und in den Küstenarchipelen.“
- D. „Es ist auch nicht erforderlich, dass Sie eine lokale Nummer von der Gemeinde erhalten, in der Sie leben. Sie können eine Internetverbindung über Satellit in der Wildnis v on Chicken in Alaska erhalten und eine Nummer auswählen, die vorgibt, dass Sie im sonnigen Arizona sind.“

Your answer letter:

Figure 1: Cross-modal topic matching ‘Images-To-Sentence’ for German with $k=5$ reference images.

English data, leading to significant limitations despite the base LLMs’ multilingual abilities. They may fail to follow instructions or struggle to interpret text within images in Non-English languages ([Schneider and Sitaram, 2024](#); [Tang et al., 2024](#)). Although many multilingual VL benchmarks exist, they typically cover at most 10 languages ([Bugliarello et al., 2022](#); [Liu et al., 2021](#); [Tang et al., 2024, *inter alia*](#)). Only concurrent work has scaled VL evaluation to 100 languages using machine translation (MT) with human post-editing ([Vayani et al., 2024](#)). Nevertheless, benchmarks constructed using semi-manual MT cannot support truly low-resource languages adequately, as current MT models lack the necessary quality for these languages. Moreover, existing benchmarks primarily assess lower-level VL semantics through concrete text-image relationships, such as those found in VQA. This underscores the need for VL benchmarks that cover truly low-resource languages and evaluate more abstract VL interactions.

To address these challenges, we introduce the massively multilingual vision-language SIB (MVL-SIB) dataset, which extends the topic labels of the multi-way parallel sentences from SIB-200 ([Adelani et al., 2024](#)) by associating each topic with hand-selected images. MVL-SIB evaluates

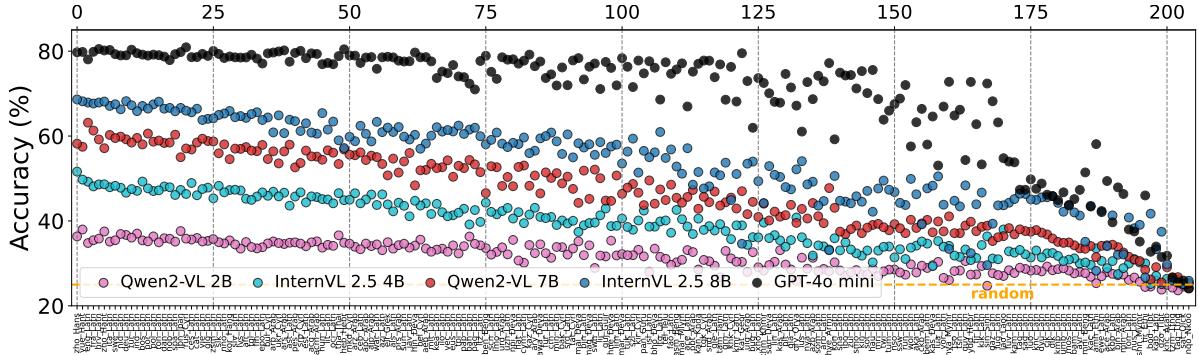


Figure 2: **Images-To-Sentences** @ $k=3$. The English prompt describes the cross-modal topic matching task, lists all topics, and provides both $k=3$ reference images and 4 sentences in the corresponding language $\{\text{eng_Latn}, \dots, \text{nqo_Nkoo}\}$. LVLMs must select the sentence of 4 options that topically fits $k=3$ reference images. The sentences spanning 205 languages and 7 topics are drawn from SIB-200 (Adelani et al., 2024), while images for the topics were hand-selected (cf. Appendix A.1). An example prompt is shown in Appendix A.7.2; further details are in §4. **Plot.** The x-axis orders the languages of the candidate sentences $\{\text{eng_Latn}, \dots, \text{nqo_Nkoo}\}$, respectively, by descending performance (y-axis). The top x-axis indicates the running index of each language L_i ($i \in \{1, \dots, 205\}$).

cross-modal image-text topic matching in 205 languages: LVLMs must select one of 4 candidate sentences that best matches the topic of the reference images (‘images-to-sentence’, cf. Figure 1) or, conversely, choose one of 4 candidate images corresponding to the topic of the reference sentences (‘sentences-to-image’). Figure 2 displays the ‘images-to-sentence’ performance for across all languages in MVL-SIB, sorted in descending order. Notably, GPT-4o-mini performs robustly on the top 125 languages. However, beyond that range its performance declines sharply, falling to chance levels in the lowest-resource languages, such as N’Koo. Bridging these gaps is crucial for developing genuinely inclusive VL technology.

Contributions. 1) MVL-SIB supports parallel VL evaluation in 205 languages on professionally translated texts, a 105 more languages than any other VL benchmark. The tasks, images-to-sentence and sentences-to-images with prefix and postfix images in context, respectively, allows for fine-grained analysis of VL interactions. We also define corresponding text-only tasks by replacing the images with the topic label to compare the VL support (by language) of LVLMs against the text-only support of their underlying LLMs. Both tasks allow to vary the number of included images to analyze the shift from single to multi-image support in LVLMs. 2) We thoroughly evaluate LVLMs on cross-modal image-text topic matching, finding that task performance is closely associated with both model size and the size of pre-training corpora of the respective languages. We further find that only GPT-4o-mini seizes on multiple references in both

cross-modal tasks. Open-weights LVLMs, moreover, favor one of the two tasks, highlighting the asymmetry in their VL support. 3) We analyze the relationship between stand-alone text and vision-language support in LVLMs by also benchmarking LVLMs on text-only topic matching. The performance gap between matching sentences to reference images or the topic tends to be larger the better the LVLM supports the underlying language. Conversely, the spread in performance between picking the fitting image or topic for reference sentences increases the worse the vision-language support of the LVLM for the evaluated language is. 4) We correlate MVL-SIB with established multilingual VL benchmarks on the languages shared between the respective pairs of datasets, showing that the MVL-SIB tasks align well with all VL tasks except OCR. We further show that images-to-sentence and sentence-to-image probe distinct aspects of VL interaction, as certain benchmarks correlate more strongly with one task than with the other. This analysis shows that MVL-SIB is a reliable and comprehensive VL benchmark for the lower-resource languages that are not covered by other datasets.

2 Related Work

Multilingual Vision-Language Models. Researchers have extended more English-centric LVLMs like BLIP-2 or LLaVA by continuing to train on multilingual data. Google’s PaLI models were the earliest closed-weight models trained on multilingual captions and VQA data; their open-weight PaliGemma followed a comparable strategy. Meanwhile, modern LLMs (e.g., Qwen 2.5,

Llama 3, Gemma 2, Aya) have improved in multilingual tasks but frequently fail to respond consistently in non-English languages, particularly in low-resource settings (Schneider and Sitaram, 2024). However, foundational models tend to focus on higher-resource languages and do not fully account for broader linguistic contexts in vision-language tasks. mBLIP is the first open multilingual LVLM trained on image captions and a small set of translated instruct data in 98 languages (Geigle et al., 2024a). Pangea incorporates multicultural dimensions by blending machine-translated data, existing multilingual resources, and synthetic data, on 39 languages (Yue et al., 2025). Most recently, Geigle et al. (2025) studied the composition of training data for multilingual adaptation of LVLMs, observing that only 25-50% of the data needs to be in English. The authors apply their findings when training ‘Centurio’, which achieves state-of-the-art performance on 14 multilingual VL benchmarks.

Multilingual Vision-Language Benchmarks. Existing datasets span VQA, natural language inference (NLI), image captioning, outlier detection, and culturally grounded QA:

VQA. xGQA extends the questions of the GQA dataset into 8 languages (5 scripts), but the answers in English (Pfeiffer et al., 2022). MaXM offers short-form QA fully in 7 languages, pairing culturally aligned images with same-language QA pairs (Changpinyo et al., 2023). MTVQA focuses on text-heavy VQA in 9 languages (Tang et al., 2024).

Culturally-grounded VQA. CVQA collects culturally diverse images and queries in 31 languages with multiple country-specific variants (Romero et al., 2025). The concurrent ALM-Bench covers 100 languages via MT with GPT-4o followed by human post editing with both generic and cultural multiple-choice and ‘true or false’ questions, as well as free-form VQA (Vayani et al., 2024).

Visual Reasoning & NLI. XVNLI evaluates cross-lingual visual NLI on 5 languages (Buggiarello et al., 2022). Binary reasoning tasks include MaRVL (Liu et al., 2021) and M5B-VGR (Schneider and Sitaram, 2024), each using linguistically specific images and textual statements. M5B-VLOD presents an outlier detection challenge, where a statement holds true for all but one image (Schneider and Sitaram, 2024).

Multiple-Choice QA. Babel-ImageNet (Geigle et al., 2024c) translates ImageNet labels into nearly 300 languages for multiple-choice object classifica-

tion. M3Exam and xMMMU also feature multiple-choice VQA in 9 and 7 languages, respectively.

MVL-SIB fills the gaps in existing multilingual VL benchmarks. It provides test data professionally translated to 205 languages, covering over 105 more languages than other benchmarks for which MT cannot synthesize reliable data for. Other VL datasets that eschew MT, such as culturally-grounded VQA benchmarks, typically construct more language-specific non-parallel data, that does not support comparative evaluation across languages. The cross-modal topic matching tasks can be also framed text-only (cf. §3.2), replacing the images that represent topics with the explicit topic labels. Thereby, MVL-SIB enables to ablate the vision-language support from the textual support for a language. Finally, the benchmark allows to vary the number of images provided LVLMs to analyze the support for multi-image reasoning.

3 Dataset and Tasks

3.1 Dataset

For MVL-SIB, we extend the following Flores-based datasets to create a massively multilingual, multi-way parallel VL benchmark for identifying topical associations between images and sentences.

Flores. Flores is a machine translation benchmark containing 3,001 sentences from English Wikipedia paragraphs (Team et al., 2022), professionally translated into over 200 languages.²

SIB-200. Adelani et al. (2024) grouped the coarse topical annotations for sentences in the DEV and DEVTEST subsets of Flores into 7 higher-level topics.³ The resulting SIB-200 dataset is a benchmark for topical classification with 1,004 parallel examples for 205 language variants.

MVL-SIB. For each topic, we first manually collect 10 permissively licensed images that distinctly represent the topic (e.g., sports) with minimal overlap or ambiguity (cf. Appendix A.1). We verify that all LVLMs in our study correctly classify the topics for the images when prompted (cf. Appendix A.7.1). We next create 3 different MVL-SIB instances from each of the 1,004 SIB sentences, totaling 3,012 MVL-SIB instances. For each MVL-SIB instance, we couple the respective SIB sen-

²Flores splits 3,001 sentences into DEV (997), DEVTEST (1,012), and TEST (992) sets. The TEST set was not released.

³The topics are entertainment, geography, health, politics, science/technology, sports, and travel.

tence with (1) a random selection of 5 positive images (same topic) and 4 additional sentences from the same category as the original sentence, as well as (2) 3 negative images and sentences randomly sampled from different topics compared to the starting SIB sentence. The set of sampled sentences and images by instance is maintained across languages.

3.2 Cross-modal & Text-only Topic Matching

We formulate both cross-modal and text-only topic matching tasks based on MVL-SIB. In every task, we present the model with the list of topics that images and sentences may be associated with.³ Otherwise, it would be unclear along which dimension the model should match images and sentences. The portion of the prompt that introduces the task is provided in English, while the sentences to be topically aligned with images are presented in one of the 205 languages included in MVL-SIB. LLMs reliably perform tasks described in English, even when task-related information is conveyed in other languages (Muennighoff et al., 2023; Romanou et al., 2025). This ensures a fair comparison across all 205 languages, where MT would not accurately preserve the meaning of the prompts. We detail prompts for each task in Appendix A.2.

Cross-modal Topic Matching. Using our text-image samples (cf. §3.1), we define two cross-modal topic matching tasks: Images-To-Sentence (I2S) and Sentences-To-Image (S2I). In I2S, the model must select, from 4 candidates, the sentence that matches the topic of k reference images. Conversely, in S2I, the model chooses, from 4 options, the image that shares the topic with k reference sentences. In both tasks, we present the model with $k \in \{1, 3, 5\}$ references, respectively. These tasks evaluate the model’s ability to align high-level visual and textual cues on topics.

Text-only Topic Matching. We construct two tasks by replacing the images in I2S and S2I, that represent the topics, with their corresponding labels (e.g., sports). The resulting unimodal tasks, Topic-To-Sentence (T2S) and Sentences-To-Topic (S2T), mirror the cross-modal tasks, I2S and S2I, respectively. For T2S, we evaluate only $k=1$, since repeating the topic label adds no information. These baseline tasks allow us to delineate between language support and vision-language understanding in LVLMs.

MVL-SIB offers 4 crucial advantages over prior benchmarks. 1) It supports evaluation in 205 languages, covering over 100 more languages than

existing benchmarks for which MT models fail to synthesize reliable evaluation data. 2) MVL-SIB supports ablating language understanding and multimodal reasoning of LVLMs by comparatively evaluating the mirroring text-only and cross-modal topic matching tasks (cf. §3.2). 3) MVL-SIB enables intricate analysis of single- and multi-image VL interactions in LVLMs by allowing topics to be represented by varying numbers of images in cross-modal tasks. 4) MVL-SIB comprises higher-level VL reasoning tasks, pairing varied images and diverse texts to test nuanced VL understanding.

4 Experimental Setup

Models. We test state-of-the-art LVLMs Qwen2-VL (Wang et al., 2024), InternVL 2.5 (Chen et al., 2024), Centurio-Qwen (Geigle et al., 2025), and GPT-4o(-mini) across available sizes.⁴ Smaller LVLMs are evaluated on all languages, while larger ones (26B+) are tested on subsets (cf. §5.3). For cross-modal topic matching, we also evaluate on mSigLIP-base (Zhai et al., 2023). Trained explicitly for semantic similarity on multilingual image-caption pairs, the ViT represents a strong baseline.⁵ Its prediction denotes the choice that has the highest average cosine similarity to the k references.

Image preprocessing. We downsample the images to 640×480 pixels, as the tasks rely on higher-level visual cues for topically associating images and texts rather than finer image details. This can significantly reduce the number of visual tokens input to LVLMs, enabling more efficient inference.

Hyperparameters. We decode text greedily with temperature set to 0.0 to ensure reproducibility.

Metric. We compute the share of prompts for which responses begin with the right letter. If the label is "A", a response such as "A." is also correct.

5 Results and Discussion

We categorize the languages in MVL-SIB based on their ‘resourceness’. To do so, we reorganize the language groups from Joshi et al. (2020) into four tiers. We first rank the tiers w.r.t. Wikipedia size and then merge (i) the two highest-resource tiers and (ii) the two lowest-resource tiers.⁶ This both

⁴We provide details on the LVLMs in Appendix A.3.

⁵The model is available on Huggingface at: [google/siglip-base-patch16-256-multilingual](https://huggingface.co/google/siglip-base-patch16-256-multilingual).

⁶Sorting by Wikipedia size (in number of pages) swaps tiers 1 and 2; we then merge Tier 0 with the new Tier 1, as well as tiers 4 and 5.

Resourceness		High			Mid			Low								
		ENGLISH		TIER 4 (26)			TIER 3 (32)			TIER 2 (96)			TIER 1 (51)			
<i>k</i> References		1	3	5	1	3	5	1	3	5	1	3	5	1	3	5
<i>Images-To-Sentence:</i> Select 1 of 4 sentences topically matching <i>k</i> reference images																
mSigLIP-base		57.7	64.6	66.4	53.3	58.6	59.7	51.4	56.2	57.1	38.9	41.2	41.7	36.1	37.6	38.0
Qwen2-VL 2B		36.3	34.8	34.9	35.5	35.3	34.1	34.5	34.3	33.2	31.0	30.6	30.0	29.5	29.0	28.6
Qwen2-VL 7B		65.8	63.1	58.9	57.7	56.5	51.7	55.4	54.5	49.6	44.3	44.4	40.5	39.6	39.7	36.5
InternVL 2.5 4B		52.5	49.2	48.1	50.3	46.6	47.7	48.6	45.3	46.1	38.7	37.1	37.3	35.4	34.5	34.5
InternVL 2.5 8B		67.7	67.9	68.7	64.6	64.9	65.7	61.2	60.8	61.6	51.0	51.4	51.8	46.1	46.0	46.3
Centurio Qwen		54.8	60.0	62.4	54.2	59.2	60.6	53.4	58.1	58.9	46.6	48.9	49.2	43.0	44.2	44.7
GPT-4o-mini		68.3	78.1	77.4	71.6	79.0	78.1	72.0	78.9	77.7	63.5	68.0	66.4	56.9	60.3	58.7
<i>Sentences-To-Image:</i> Select 1 of 4 images topically matching <i>k</i> reference sentences																
mSigLIP-base		56.3	66.0	69.6	51.8	61.6	64.0	49.1	58.3	60.2	36.0	40.4	41.2	32.9	36.3	36.9
Qwen2-VL 2B		41.9	43.1	43.4	41.6	42.5	42.7	40.8	42.4	42.4	33.7	35.6	35.5	31.0	32.7	32.8
Qwen2-VL 7B		71.7	70.4	68.6	65.5	65.5	63.5	64.4	65.3	64.1	50.3	52.5	43.5	45.9	46.6	
InternVL 2.5 4B		47.7	44.5	43.0	38.0	40.3	40.4	36.7	39.6	40.3	30.7	34.4	35.7	28.8	32.1	33.7
InternVL 2.5 8B		66.2	69.0	68.7	57.5	62.5	61.6	52.9	58.5	58.1	43.4	49.8	49.7	39.7	45.9	46.2
Centurio Qwen		35.3	36.1	35.6	31.1	32.9	33.3	31.0	32.8	33.1	28.7	29.7	29.8	28.1	28.7	28.7
GPT-4o-mini		77.5	86.4	89.1	77.2	86.5	88.6	77.1	86.1	88.4	68.4	79.8	82.7	61.7	74.0	77.2

Table 1: **Cross-modal Topic Matching:** LVLMs must select the candidate sentence (image) from 4 choices that topically align with *k* reference images (sentences). Prompts provided in §A.2. Languages are tiered by Wikipedia sizes (cf. §5). Number of languages in parentheses. **Metric:** share of responses starting with correct option letter. Details in §4. In each column, the best model is emphasized in **bold**, the second-best model is underlined.

Resourceness		High			Mid			Low									
		ENGLISH		TIER 4 (26)			TIER 3 (32)			TIER 2 (96)			TIER 1 (51)				
Task		Topic-To Sentence	Sentences To-Topic														
<i>k</i> References		1	1	3	5	1	1	3	5	1	1	3	5	1	1	3	5
Qwen2-VL 2B		56.7	86.1	96.5	98.2	49.1	78.4	92.0	95.2	45.2	73.7	89.6	93.6	36.4	53.7	69.3	76.4
Qwen2-VL 7B		85.7	89.1	95.3	97.5	81.8	84.1	93.3	96.1	80.5	84.1	93.5	96.6	63.7	67.8	81.3	86.7
InternVL 2.5 4B		81.4	90.6	98.0	99.1	72.7	85.2	95.8	97.9	68.2	83.7	95.5	97.7	50.2	67.8	83.4	87.9
InternVL 2.5 8B		87.0	91.7	98.1	<u>99.0</u>	83.0	86.3	96.3	98.3	79.1	82.4	94.1	96.7	65.4	66.8	81.8	86.9
Centurio Qwen		85.4	89.9	96.7	97.7	83.6	88.0	95.8	97.7	82.6	87.7	96.0	97.9	70.4	73.1	86.8	90.7
GPT-4o-mini		88.5	92.4	98.1	99.1	89.3	91.6	98.4	99.3	89.3	91.6	98.4	99.3	80.9	82.1	93.3	95.7

Table 2: **Text-only Topic Matching:** LVLMs must select the candidate sentence (topic) of 4 choices that aligns topically with the reference topics (*k* reference sentences). See Table 1 for further details.

better reflects current corpus availability for LLM pre-training (Xue et al., 2021; Kudugunta et al., 2023), and aligns with downstream performance (cf. Appendix A.3.1). We isolate English from Tier 4, since it is the pivotal language in NLP. The full per-language results by task and model are provided in Appendix A.7.

5.1 Cross-modal Topic Matching

Images-To-Sentence (I2S). The upper segment of Table 1 displays the results for I2S, in which the LVLMs must pick the candidate sentence that topically matches the *k* reference images.

English. The performance on I2S with English candidate sentences scales well with model size. The small Qwen2-VL 2B performs only slightly better than chance (25% vs. ca. 35%). The comparably sized Qwen2-VL 7B, InternVL 2.5 8B, and Centurio-Qwen 8B peak around 62% to 68% at var-

ious *k*. These models nevertheless non-negligibly trail GPT-4o-mini (78.1%). Among LVLMs, only InternVL 2.5 8B, Centurio-Qwen, and GPT-4o-mini benefit from multiple reference images. When the number of references *k* increases from 3 to 5, GPT-4o-mini declines slightly in performance, while InternVL and Centurio-Qwen continue to improve marginally (ca. +1%) and more notably (ca. +3-6%), respectively. All other LVLMs deteriorate materially with more reference images (ca. -3-4%). mSigLIP indeed is a strong baseline, trailing only GPT-4o-mini and InternVL 2.5 8B at *k*=5. The ViT yields large gains with 4 more images (+9.7%).

Tiers. The performance gap of other languages to English correlates well with their resource levels by language tier. When presented with candidate sentences in non-English high-resource languages (cf. Tier 4), GPT-4o-mini even performs better slightly better. For very low-resource languages in

Tier 1, such as N’Koo or Tamazight, all models fail to perform better than chance (cf. Appendix A.7.2). Among LVLMs, only GPT-4o-mini remains overall robust for topically matching sentences of low-resource languages to images, whereas other models drop severely in performance (ca. 15–20%). While mSigLIP still performs well, it declines more notably than LVLMs on lower-resource languages.

Sentences-To-Image (S2I). The lower part of Table 1 presents the results for S2I. Here, the models select the candidate image among four options that topically fits the k reference sentences.

English. Performance again correlates well with model capacity. However, in S2I, only GPT-4o-mini significantly seizes on additional references (+13%) to excel with 89%, while models like Qwen2-VL 7B and InternVL 2.5 8B exhibit peak performance at $k=1$ and $k=3$, respectively, that taper slightly with more sentence references. Centurio-Qwen performs only slightly better than random (25% vs. ca. 35%). In S2I, mSigLIP is again very strong, second only to GPT-4o-mini at $k=5$. The encoder once more seizes sizable gains from additional references (+13.3%).

Tiers. In non-English evaluations, the overall trend remains similar, though absolute performance is lower. The gap between high- and low-resource language tiers is evident, as all models yield higher scores across Tiers 4 and 3. GPT-4o-mini maintains robust performance even in the most challenging Tier 1 when provided multiple references (ca. 75%).

In sum, both the training protocol and the model size collectively determine whether models favor I2S or S2I. For instance, InternVL 2.5 8B outperforms Qwen2-VL 7B on I2S across the board, while trailing on S2I. Moreover, only GPT-4o-mini consistently seizes on additional references and remains largely robust to the lowest-resource languages on both tasks. This likely stems from insufficient training to enable open-weight LVLMs to perform higher-level VL reasoning with diverse texts and multiple images successfully. For instance, Centurio Qwen was mostly trained on data that prefixes images to text, resulting in low performance when images are postfixed to the context.

5.2 Text-only Topic Matching

Table 2 lists the results for text-only topic matching, in which the images are exchanged with their topic label. These tasks denominate ‘upper-bounds’ for their cross-modal counterparts to enable ablations

of language support and VL support in LVLMs.

Topic-to-Sentence (T2S). In this task, LVLMs choose the sentence that best suits the reference topic. Barring Qwen2-VL 2B, all models perform well on the task. Notably, 5 images should capture the underlying topic well for all models (cf. Appendix A.7.1). Despite that, the gap between text-only and vision-language tasks (cf. Table 1) is sizable across all open-weights models for ‘English’ (ca. 20%+). In English, GPT-4o-mini and InternVL 2.5 8B achieve the highest accuracies, indicating strong topic comprehension. For non-English languages, while the overall scores are reduced, high-resource languages benefit from richer training signals compared to their low-resource counterparts – models like Qwen2-VL 2B and Centurio-Qwen 8B show a more pronounced drop in the latter, underscoring the impact of language resources.

Sentences-to-Topic (S2T). In S2T, where models choose the topic that best aligns with k reference sentences, performance scales both with model size and the number of references. The gains from additional context (3–5 k) are particularly notable for larger LVLMs. In English, GPT-4o-mini improves markedly from 92.4% at $k=1$ to 98.1% at $k=3$. In non-English languages, similar patterns emerge: high-resource languages consistently yield higher accuracies than low-resource ones, with GPT-4o-mini and InternVL 2.5 8B exhibiting the most stable improvements across varying k . This reinforces the role of model capacity and training data diversity in effective cross-lingual topic matching.

Comparing the results of cross-modal and text-only topic matching sheds further light on the VL interactions for lower-resource languages in LVLMs. The performance gap between matching sentences to reference images and matching them to textual topics tends to narrow regardless of the number of references, likely reflecting their limited textual support in LVLMs. In contrast, the discrepancy between selecting an image versus a topic for reference sentences becomes much more pronounced, especially at $k=5$. These findings suggest that VL support degrades more sharply than textual support for lower-resource languages in LVLMs.

5.3 Further Analyses

Task Correlation. To compare MVL-SIB with other tasks, we access the results for Qwen2-VL and InternVL 2.5 models on several established multilingual VL benchmarks from Geigle et al.

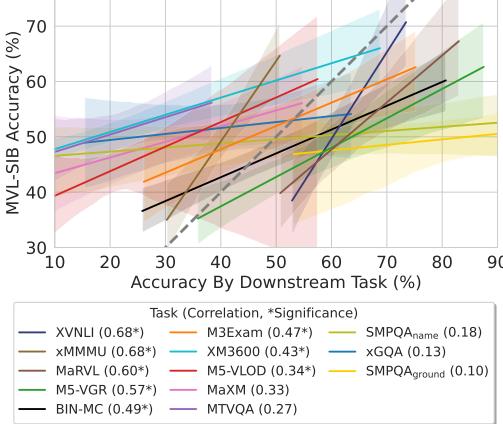


Figure 3: I2S with $k=3$.

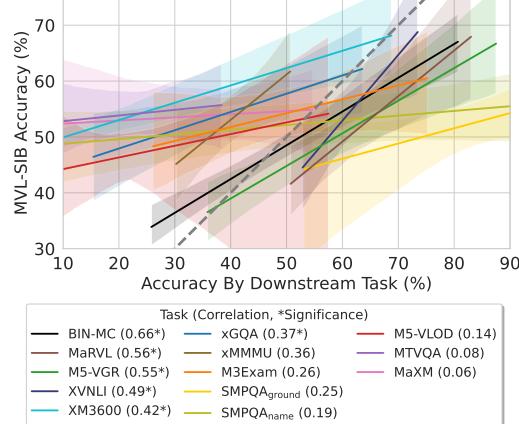


Figure 4: S2I with $k=3$.

Correlations Between MVL-SIB & Multilingual VL Benchmarks. Pearson correlation coefficients obtained by regressing MVL-SIB performance against performance on multilingual VL tasks on languages common to both datasets, respectively. An asterisk (*) indicates whether the coefficient is statistically significant at $p \leq 0.05$.

(2025).⁷ Next, we align the results across languages for the MVL-SIB tasks and the other benchmarks. Finally, we plot the linear regressions of performance on I2S and S2I with $k=3$ against the performance on the VL benchmarks, respectively, pooled over models, in Figures 3 and 4.⁸

MVL-SIB positively correlates with all tasks in both the I2S and S2I evaluations. However, both the magnitude and statistical significance of these linear relationships vary across VL benchmarks. Both I2S and S2I exhibit the strongest connections with XVNLI (visual inference), BIN-MC (multiple-choice image classification), MaRVL, and M5-VGR (both visually-grounded boolean reasoning). Since these tasks restrict valid answers to a small set of fixed options (i.e., choice letters or ‘yes/no/maybe’), LVLMs must engage in higher-level vision-language disambiguation rather than relying solely on lower-level visual cues to solve these tasks. In addition, xMMMU, M3Exam, and M5-VLOD are significantly related only to I2S, whereas xGQA is significantly aligned solely with S2I. The former tasks are structurally similar to I2S, typically presenting one or more images that LVLMs are given as visual context to answer multiple-choice questions. We hypothesize that only S2I is significantly correlated with xGQA, since S2I is more analogous to object-centric benchmarks. In S2I, LVLMs likely leverage targeted semantic cues (e.g., keywords or phrases) from the

⁷We omit Centurio-Qwen, since it degrades on the S2I task. We further provide details on all the multilingual vision-language benchmarks we correlate MVL-SIB with in Appendix A.4

⁸Note that we include a constant in our regression model to bridge task-specific scales of results.

reference sentences to better disambiguate candidate images by topic. This behavior aligns with lower-level VL tasks such as xGQA (cross-lingual short-form QA) or BIN-MC (multiple-choice object classification), where texts and images are more deliberately connected. In contrast, MTVQA, SMPQA-name, and SMPQA-ground show only weak or statistically insignificant correlations with the MVL-SIB tasks. Since these tasks require LVLMs to comprehend text embedded in images, a low-level, fine-grained VL task, they differ substantially from the higher-level VL reasoning evaluated by MVL-SIB.

Overall, the regression analysis indicates that I2S and S2I capture distinct yet complementary aspects of VL understanding, collectively exhibiting a strong relationship with a broad set of VL tasks. This aligns with our main results (cf. Table 1), which show that different LVLMs may favor one MVL-SIB task over the other. This renders MVL-SIB as a suitable benchmark for evaluating *universal* VL understanding (across 205 languages). It also enables to ablate performance across the key axes of analysis, the task formulation (I2S vs. S2I), the language vs. vision-language support for 205 languages, and the number of images in context.⁹

Larger LVLMs. To evaluate larger LVLMs on MVL-SIB, we construct language-tier subsets that reliably estimate performance while mitigating excessive computational overhead (cf. §5). For both I2S and S2I, we identify the three languages in each tier that best replicate the average performance of

⁹In S2I, candidates could comprise more than a single image.

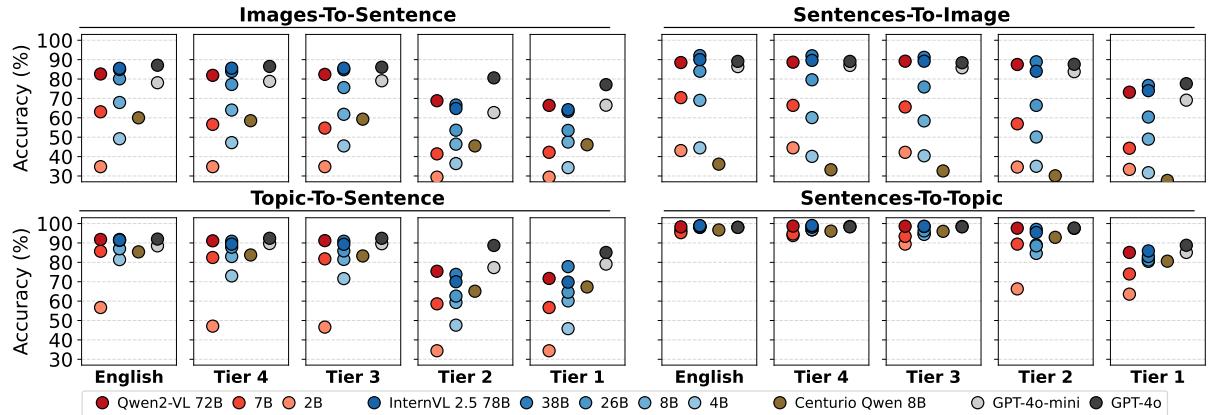


Figure 5: **Larger LVLMs on Subsampled Tiers.** We extract 3 languages per tier that mimic avg. performance full language groups (cf. §5.3) and evaluate LVLMs across all model sizes on {12S,S2I,T2S,S2T} @ $k=3$ (cf. §3.2).

the tier. First, we compute the average performance per language tier for both 12S and S2I with $k=3$, pooling results across models (cf. §5). Then, for each tier, we select the three languages whose performance deviates least from the tier mean. The languages chosen for each tier by task are detailed in Appendix A.6. Finally, we test GPT-4o, InternVL 2.5 {26,32,72}B, and Qwen 2.5 VL {32,70}B on these subsets to assess how well larger models perform across language tiers.

Figure 5 displays the results for both 12S and S2I with $k=3$.¹⁰ We observe that larger models catch up to GPT-4o on 12S and even outperform it on S2I for higher-resource languages. Since open-weight LVLMs have more limited language support for low-resource languages, GPT-4 and GPT-4o-mini outperform all other models on 12S for languages in tier 1 and 2. Increasing model capacity yields the largest gains on S2I, for which the models exceed GPT-4o up to the lowest-resource Tier 1. Moreover, unlike smaller models, all larger LVLMs (+26B) more effectively leverage multiple image references to improve performance (cf. Appendix A.6). They reap the largest benefit when the number of references increases from 1 to 3 (ca. +3%), with further improvements at $k = 5$ (ca. +1.5%). These results suggest that larger LVLMs have fundamentally better VL support irrespective of the evaluated language (cf. §5.2). The results on text-only topic matching further underscore this notion (cf. lower segment of Figure 5). Larger LVLMs near perfectly match both the reference topic to the correct sentence (ca. 90%) and references sentences to the right topic (ca. 98%). Notably, as model size increases, the performance gap between

¹⁰Qwen2-VL 72B frequently responded with the first letter of the correct topic. If that letter uniquely identifies the correct choice, the answer is considered correct.

cross-modal and text-only matching narrows. Collectively, these results further indicate that VL support relative to text-only support improves with increasing model capacity in LVLMs.

6 Conclusion

We present the massively multilingual vision-language benchmark MVL-SIB for cross-modal (and text-only topic matching) in 205 languages that offers key advantages over prior multilingual VL benchmarks. Notably, it covers over 100 additional languages without relying on machine translation. MVL-SIB allows for a clear separation between textual language support and vision-language support in LVLMs by comparing performance on mirrored cross-modal and text-only tasks. Moreover, it allows us to study how LVLMs handle single-image versus multi-image formulations of cross-modal topic matching by varying the number of images provided. In our comparative evaluation of state-of-the-art LVLMs on MVL-SIB, we find that model performance is strongly correlated with both model size and the volume of available pre-training data for each language. However, all LVLMs experience a dramatic performance drop on the lowest-resource languages. Our analysis further reveals that vision-language support deteriorates disproportionately relative to language support, highlighting the need to incorporate low-resource languages into VL training. Moreover, providing multiple images does not benefit open-weight LVLMs in cross-modal topic matching, suggesting that LVLMs are not yet fully effective in multi-image tasks. Lastly, we validate that MVL-SIB correlates well with existing multilingual VL benchmarks, underscoring its reliability as a source of evaluation data for 205 languages.

7 Limitations

Our work faces three primary limitations. First, although a vast number of LVLMs exist, we selected a representative subset based on key criteria. Specifically, the LVLMs in our study (Qwen2-VL and InternVL, with the exception of Centurio) span a range of parameter counts typical of LLMs. Additionally, we include GPT-4o-mini in the full evaluation and GPT-4o on the subsampled language tiers. Evaluating MVL-SIB across all four tasks I2S, S2I, T2S, and S2T (cf. §3.2) at various $k \in \{1, 3, 5\}$ over 205 languages (i.e., evaluations per model and task, or 2050, in sum per model) becomes computationally intractable. This accumulates to $3 \times 205 = 615$ evaluations per model (205 for T2S as only $k=1$ reference topic exists) or $3 \times 205 + 205 = 2050$ evaluations in total. We therefore both provide subsets of the language tiers to evaluate on and demonstrate that evaluation only requires 1K instances to reliably estimate task performance. Second, while we strove to choose a diverse set of images to capture the full semantic range of each topic, further diversification is possible by sourcing additional images. However, due to the limited availability of openly licensed images, some topics (e.g., politics and entertainment) are represented predominantly by images that embody the topic in a more Western-centric cultural context. Hand-selecting images by topic for each language or, more broadly, cultural groups would not scale to 205 languages and would hinder the comparability of results. Our results furthermore confirm that models just as well perform on a broad range of languages spanning diverse cultural backgrounds as on English (cf. Figure 2). At the same time, LVLMs perform best on Western-centric images, mitigating any variation that would originate from using more culture-specific images. Finally, for the topic geography, we manually selected images that are representative within the context of SIB, as the broader definition of geography is too diffuse to capture visually.

Acknowledgements

We used AI assistance (chatGPT o3-mini) to polish the writing and the tables of the manuscript as well as to refine the code for our visualizations.

References

- David Adelani, Hannah Liu, Xiaoyu Shen, Nikita Vassilyev, Jesujoba Alabi, Yanke Mao, Haonan Gao, and En-Shiun Lee. 2024. **SIB-200: A Simple, Inclusive, and Big Evaluation Dataset for Topic Classification in 200+ Languages and Dialects**. In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 226–245, St. Julian’s, Malta.
- Željko Agić and Natalie Schluter. 2018. **Baselines and Test Data for Cross-Lingual Inference**. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, pages 3890–3894, Miyazaki, Japan.
- Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. 2015. **A Large Annotated Corpus for Learning Natural Language Inference**. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing, EMNLP*, pages 632–642, Lisbon, Portugal.
- Emanuele Bugliarello, Fangyu Liu, Jonas Pfeiffer, Siva Reddy, Desmond Elliott, Edoardo Maria Ponti, and Ivan Vulić. 2022. **IGLUE: A Benchmark for Transfer Learning across Modalities, Tasks, and Languages**. In *International Conference on Machine Learning*, pages 2370–2392, Baltimore, MD, USA.
- Soravit Changpinyo, Linting Xue, Michal Yarom, Ashish Thapliyal, Idan Szpektor, Julien Amelot, Xi Chen, and Radu Soricut. 2023. **MaXM: Towards Multilingual Visual Question Answering**. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 2667–2682, Singapore.
- Zhe Chen, Jiannan Wu, Wenhui Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, et al. 2024. **InternVL: Scaling Up Vision Foundation Models and Aligning for Generic Visual-Linguistic Tasks**. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 24185–24198, Seoul, Korea.
- J. Deng, W. Dong, R. Socher, L. Li, Kai Li, and Li Fei-Fei. 2009. **ImageNet: A large-scale hierarchical image database**. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255, Miami, FL, USA.
- William Gaviria Rojas, Sudnya Diamos, Keertan Kini, David Kanter, Vijay Janapa Reddi, and Cody Coleman. 2022. **The Dollar Street Dataset: Images Representing the Geographic and Socioeconomic Diversity of the World**. In *Advances in Neural Information Processing Systems*, volume 35, pages 12979–12990, New Orleans, LA, USA.
- Gregor Geigle, Abhay Jain, Radu Timofte, and Goran Glavaš. 2024a. **mBLIP: Efficient Bootstrapping of Multilingual Vision-LLMs**. In *Proceedings of the 3rd Workshop on Advances in Language and Vision Research (ALVR)*, pages 7–25, Bangkok, Thailand.

- Gregor Geigle, Florian Schneider, Carolin Holtermann, Chris Biemann, Radu Timofte, Anne Lauscher, and Goran Glavaš. 2025. **Centurio: On Drivers of Multilingual Ability of Large Vision-Language Model**. *arXiv preprint arXiv:2501.05122*.
- Gregor Geigle, Radu Timofte, and Goran Glavaš. 2024b. **African or European Swallow? Benchmarking Large Vision-Language Models for Fine-Grained Object Classification**. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 2653–2669, Miami, Florida, USA.
- Gregor Geigle, Radu Timofte, and Goran Glavaš. 2024c. **Babel-ImageNet: Massively Multilingual Evaluation of Vision-and-Language Representations**. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5064–5084, Bangkok, Thailand.
- Drew A. Hudson and Christopher D. Manning. 2019. **GQA: A New Dataset for Real-World Visual Reasoning and Compositional Question Answering**. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*, pages 6700–6709, Long Beach, CA, USA.
- Pratik Joshi, Sebastin Santy, Amar Budhiraja, Kalika Bali, and Monojit Choudhury. 2020. **The State and Fate of Linguistic Diversity and Inclusion in the NLP World**. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6282–6293, Online.
- Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A. Shamma, Michael S. Bernstein, and Li Fei-Fei. 2017. **Visual Genome: Connecting Language and Vision Using Crowdsourced Dense Image Annotations**. *Int. J. Comput. Vision*, 123(1):32–73.
- Sneha Kudugunta, Isaac Caswell, Biao Zhang, Xavier Garcia, Christopher A. Choquette-Choo, Katherine Lee, Derrick Xin, Aditya Kusupati, Romi Stella, Ankur Bapna, and Orhan Firat. 2023. **MADLAD-400: A Multilingual And Document-Level Large Audited Dataset**. In *Advances in Neural Information Processing Systems*, volume 36, pages 67284–67296, New Orleans, LA, USA.
- Fangyu Liu, Emanuele Bugliarello, Edoardo Maria Ponti, Siva Reddy, Nigel Collier, and Desmond Elliott. 2021. **Visually Grounded Reasoning across Languages and Cultures**. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 10467–10485, Punta Cana, Dominican Republic.
- Niklas Muenninghoff, Thomas Wang, Lintang Sutawika, Adam Roberts, Stella Biderman, Teven Le Scao, M Saiful Bari, Sheng Shen, Zheng Xin Yong, Hailey Schoelkopf, Xiangru Tang, Dragomir Radev, Alham Fikri Aji, Khalid Almubarak, Samuel Albanie, Zaid Alyafeai, Albert Webson, Edward Raff, and Colin Raffel. 2023. **Crosslingual Generalization through Multitask Finetuning**. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15991–16111, Toronto, Canada.
- Jonas Pfeiffer, Gregor Geigle, Aishwarya Kamath, Jan-Martin O. Steitz, Stefan Roth, Ivan Vulić, and Iryna Gurevych. 2022. **xGQA: Cross-Lingual Visual Question Answering**. In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 2497–2511, Dublin, Ireland.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021. **Learning Transferable Visual Models From Natural Language Supervision**. In *International Conference on Machine Learning*, pages 8748–8763, Online.
- Angelika Romanou, Negar Foroutan, Anna Sotnikova, Zeming Chen, Sree Harsha Nelaturu, Shivalika Singh, Rishabh Maheshwary, Micol Altomare, Mohamed A. Haggag, Sneha A, Alfonso Amayuelas, Azril Hafizi Amirudin, Viraat Aryabumi, Danylo Boiko, Michael Chang, Jenny Chim, Gal Cohen, Aditya Kumar Dalmia, Abraham Diress, Sharad Duwal, Daniil Dzehnaliou, Daniel Fernando Erazo Florez, Fabian Farestam, Joseph Marvin Imperial, Shayekh Bin Islam, Perttu Isotalo, Maral Jabbarishiviari, Börje F. Karlsson, Eldar Khalilov, Christopher Klamm, Fajri Koto, Dominik Krzeminski, Gabriel Adriano de Melo, Syrielle Montariol, Yiyang Nan, Joel Niklaus, Jekaterina Novikova, Johan Samir Obando Ceron, Debjit Paul, Esther Ploeger, Jebish Purbey, Swati Rajwal, Selvan Sunitha Ravi, Sara Rydell, Roshan Santhosh, Drishti Sharma, Marjana Prifti Skenduli, Arshia Soltani Moakhar, Bardia Soltani Moakhar, Ran Tamir, Ayush Kumar Tarun, Azmine Toushik Wasi, Thenuka Ovin Weerasinghe, Serhan Yilmaz, Mike Zhang, Imanol Schlag, Marzieh Fadaee, Sara Hooker, and Antoine Bosselut. 2025. **INCLUDE: Evaluating Multilingual Language Understanding with Regional Knowledge**. In *The Thirteenth International Conference on Learning Representations*, Singapore.
- David Romero, Chenyang Lyu, Haryo Akbarianto Wi-bowo, Teresa Lynn, Injy Hamed, Aditya Nanda Kishore, Aishik Mandal, Alina Dragonetti, Artem Abzaliev, Atnafu Lambebo Tonja, Bontu Fufa Balcha, Chenxi Whitehouse, Christian Salamea, Dan John Velasco, David Ifeoluwa Adelani, David Le Meur, Emilio Villa-Cueva, Fajri Koto, Fauzan Farooqui, Frederico Belcavello, Ganzorig Batnasan, Gisela Vallejo, Grainne Caulfield, Guido Ivetta, Haiyue Song, Henok Biadgign Ademtew, Hernán Maina, Holy Lovenia, Israel Abebe Azime, Jan Christian Blaise Cruz, Jay Gala, Jiahui Geng, Jesus-German Ortiz-Barajas, Jinheon Baek, Joelyn Dunstan, Laura Alonso Alemany, Kumaranage Ravindu Yasas Nagasinghe, Luciana Benotti, Luis Fernando D’Haro, Marcelo Viridiano, Marcos Esteche-Garitagotia, Maria Camila Buitrago

Cabrera, Mario Rodríguez-Cantelar, Mélanie Jouiteau, Mihail Mihaylov, Mohamed Fazli Mohamed Imam, Muhammad Farid Adilazuarda, Munkhjargal Gochoo, Munkh-Erdene Otgonbold, Naome Etori, Olivier Niyomugisha, Paula Mónica Silva, Pranjal Chitale, Raj Dabre, Rendi Chevi, Ruochen Zhang, Ryandito Diandaru, Samuel Cahyawijaya, Santiago Góngora, Soyeong Jeong, Sukannya Purkayastha, Tatsuki Kuribayashi, Teresa Clifford, Thanmay Jayakumar, Tiago Timponi Torrent, Toqueer Ehsan, Vladimir Araujo, Yova Kementchedjhieva, Zara Burzo, Zheng Wei Lim, Zheng Xin Yong, Oana Ignat, Joan Nwatu, Rada Mihalcea, Thamar Solorio, and Alham Fikri Aji. 2025. **CVQA: Culturally-diverse Multilingual Visual Question Answering Benchmark**. In *Advances in Neural Information Processing Systems*, volume 37, pages 11479–11505.

Florian Schneider and Sunayana Sitaram. 2024. **M5 – A Diverse Benchmark to Assess the Performance of Large Multimodal Models Across Multilingual and Multicultural Vision-Language Tasks**. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 4309–4345, Miami, Florida, USA.

Hai-Long Sun, Da-Wei Zhou, Yang Li, Shiyin Lu, Chao Yi, Qing-Guo Chen, Zhao Xu, Weihua Luo, Kaifu Zhang, De-Chuan Zhan, and Han-Jia Ye. 2024. **Parrot: Multilingual Visual Instruction Tuning**. *arXiv preprint arXiv:2406.02539*.

Jingqun Tang, Qi Liu, Yongjie Ye, Jinghui Lu, Shu Wei, Chunhui Lin, Wanqing Li, Mohamad Fitri Faiz Bin Mahmood, Hao Feng, Zhen Zhao, Yanjie Wang, Yuliang Liu, Hao Liu, Xiang Bai, and Can Huang. 2024. **MTVQA: Benchmarking Multilingual Text-Centric Visual Question Answering**. *arXiv preprint arXiv:2405.11985*.

NLLB Team, Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Hefernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Rogers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. 2022. **No Language Left Behind: Scaling Human-Centered Machine Translation**. *arXiv preprint arXiv:2207.04672*.

Ashish V. Thapliyal, Jordi Pont Tuset, Xi Chen, and Radu Soricut. 2022. **Crossmodal-3600: A Massively Multilingual Multimodal Evaluation Dataset**. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 715–729, Abu Dhabi, United Arab Emirates.

Ashmal Vayani, Dinura Dissanayake, Hasindri Watawana, Noor Ahsan, Nevasini Sasikumar, Omkar

Thawakar, Henok Biadgign Ademtew, Yahya Hmaiti, Amandeep Kumar, Kartik Kuckreja, Mykola Maslych, Wafa Al Ghallabi, Mihail Mihaylov, Chao Qin, Abdelrahman M. Shaker, Mike Zhang, Mahardika Krisna Ihsani, Amiel Esplana, Monil Gokani, Shachar Mirkin, Harsh Singh, Ashay Srivastava, Endre Hamerlik, Fathinah Asma Izzati, Fadillah Adamsyah Maani, Sebastian Cavada, Jenny Chim, Rohit Gupta, Sanjay Manjunath, Kamila Zhumakhanova, Feno Heriniaina Rabevohitra, Azril Amirudin, Muhammad Ridzuan, Daniya Kareem, Keitan More, Kunyang Li, Pramesh Shakya, Muhammad Saad, Amirpouya Ghasemaghaei, Amirkbek Djanibekov, Dilshod Azizov, Branislava Jankovic, Naman Bhatia, Alvaro Cabrera, Johan Obando-Ceron, Olympiah Otieno, Fabian Farestam, Muztoba Rabbani, Sanoojan Baliah, Santosh Sanjeev, Abduragim Shtanchaev, Maheen Fatima, Thao Nguyen, Amrin Kareem, Toluwaní Aremu, Nathan Xavier, Amit Bhatkal, Hawau Toyin, Aman Chadha, Hisham Cholakkal, Rao Muhammad Anwer, Michael Felsberg, Jorma Laaksonen, Thamar Solorio, Monojit Choudhury, Ivan Laptev, Mubarak Shah, Salman Khan, and Fahad Khan. 2024. **All Languages Matter: Evaluating LMMs on Culturally Diverse 100 Languages**. *arXiv preprint arXiv:2411.16508*.

Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhi-hao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. 2024. **Qwen2-VL: Enhancing Vision-Language Model’s Perception of the World at Any Resolution**. *arXiv preprint arXiv:2409.12191*.

Ning Xie, Farley Lai, Derek Doran, and Asim Kadav. 2019. **Visual Entailment: A Novel Task for Fine-grained Image Understanding**. In *32nd Conference on Neural Information Processing Systems, ViGIL Workshop*, Montreal, Canada.

Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2021. **mT5: A Massively Multilingual Pre-trained Text-to-Text Transformer**. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics (NAACL): Human Language Technologies*, pages 483–498, Online.

An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jianxin Yang, Jin Xu, Jingren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Xuejing Liu, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan,

Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, Zhifang Guo, and Zhihao Fan. 2024. [Qwen2 Technical Report](#). *arXiv preprint arXiv:2407.10671*.

Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, Cong Wei, Botao Yu, Ruibin Yuan, Renliang Sun, Ming Yin, Boyuan Zheng, Zhenzhu Yang, Yibo Liu, Wenhao Huang, Huan Sun, Yu Su, and Wenhua Chen. 2024. **MMU: A Massive Multi-discipline Multimodal Understanding and Reasoning Benchmark for Expert AGI**. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9556–9567, Seattle, WA, USA.

Xiang Yue, Yueqi Song, Akari Asai, Seungone Kim, Jean de Dieu Nyandwi, Simran Khanuja, Anjali Kantharuban, Lintang Sutawika, Sathyanarayanan Ramamoorthy, and Graham Neubig. 2025. **Pangea: A Fully Open Multilingual Multimodal LLM for 39 Languages**. In *The Thirteenth International Conference on Learning Representations*, Singapore.

Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. 2023. **Sigmoid Loss for Language Image Pre-Training**. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11975–11986, Paris, France.

Wenxuan Zhang, Mahani Aljunied, Chang Gao, Yew Ken Chia, and Lidong Bing. 2023. **M3Exam: A Multilingual, Multimodal, Multilevel Benchmark for Examining Large Language Models**. In *Advances in Neural Information Processing Systems*, volume 36, pages 5484–5505, New Orleans, LA, USA.

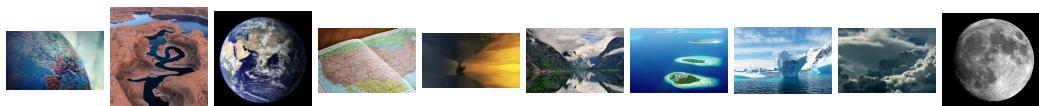
A Appendix

A.1 Images Per Topic

Entertainment



Geography



Health



Politics



Science & Technology



Sports



Travel



A.2 Prompts

Images-To-Sentences (I2S)

Which sentence best matches the topic of the images? The images and the sentences each belong to one of the following topics: "entertainment", "geography", "health", "politics", "science and technology", "sports", or "travel". Choose one sentence from A, B, C, or D. Output only a single letter!

Images



Sentences

- A. ` `` Maroochydore führte am Ende die Rangfolge an, mit sechs Punkten Vorsprung vor Noosa als Zweitem.`` ``
- B. ` `` Es wurden keine schwere Verletzungen gemeldet, jedoch mussten mindestens fünf der zur Zeit der Explosion Anwesenden aufgrund von Schocksymptomen behandelt werden.`` ``
- C. ` `` Finnland ist ein großartiges Reiseziel für Bootstouren. Das „Land der tausend Seen“ hat auch Tausende von Inseln – in den Seen und in den Küstenarchipelen.`` ``
- D. ` `` Es ist auch nicht erforderlich, dass Sie eine lokale Nummer von der Gemeinde erhalten, in der Sie leben. Sie können eine Internetverbindung über Satellit in der Wildnis von Chicken in Alaska erhalten und eine Nummer auswählen, die vorgibt, dass Sie im sonnigen Arizona sind.`` ``

Your answer letter:

Sentences-To-Images (S2I)

Which image best matches the topic of the sentences? The sentences and the images each belong to one of the following topics: "entertainment", "geography", "health", "politics", "science and technology", "sports", or "travel". Choose one image from A, B, C, or D. Output only a single letter!

Sentences

- ``` Maroochydore führte am Ende die Rangfolge an, mit sechs Punkten Vorsprung vor Noosa als Zweitem.```
- ``` Die Schlagmänner der mittleren Reihe, Sachin Tendulkar und Rahul Dravid, zeigten gute Leistungen und erzielten eine Partnerschaft mit 100 Runs.```
- ``` Da pro Tag nur achtzehn Medaillen zur Verfügung stehen, hat es ein Anzahl an Ländern nicht auf das Podium geschafft.```
- ``` Wintersportarten sind in den nördlichen Regionen am beliebtesten und Italiener nehmen an internationalen Wettkämpfen und olympischen Spielen teil.```
- ``` Nach dem Rennen bleibt Keselowski mit 2.250 Punkten Spitzensreiter in der Fahrerwertung.

Images

A.



B.



C.



D.



Your answer letter:

Topic-To-Sentence (T2S)

Which sentence best matches the topic "sports"? The sentences each belong to one of the following topics: "entertainment", "geography", "health", "politics", "science and technology", "sports", or "travel". Choose one sentence from A, B, C, or D. Output only a single letter!

Sentences

- A. ````Maroochydore führte am Ende die Rangfolge an, mit sechs Punkten Vorsprung vor Noosa als Zweitem.````
- B. ````Es wurden keine schweren Verletzungen gemeldet, jedoch mussten mindestens fünf der zur Zeit der Explosion Anwesenden aufgrund von Schocksymptomen behandelt werden.````
- C. ````Finnland ist ein großartiges Reiseziel für Bootstouren. Das „Land der tausend Seen“ hat auch Tausende von Inseln – in den Seen und in den Küstenarchipelen.````
- D. ````Es ist auch nicht erforderlich, dass Sie eine lokale Nummer von der Gemeinde erhalten, in der Sie leben. Sie können eine Internetverbindung über Satellit in der Wildnis von Chicken in Alaska erhalten und eine Nummer auswählen, die vorgibt, dass Sie im sonnigen Arizona sind.````

Your answer letter:

Sentences-To-Topics (S2T)

Which topic best matches the sentences? The sentences belong to one of the following topics: "entertainment", "geography", "health", "politics", "science and technology", "sports", or "travel". Choose one topic from A, B, C, or D. Output only a single letter!

Sentences

- `` `Maroochydore führte am Ende die Rangfolge an, mit sechs Punkten Vorsprung vor Noosa als Zweitem.` ``
- `` `Die Schlagmänner der mittleren Reihe, Sachin Tendulkar und Rahul Dravid, zeigten gute Leistungen und erzielten eine Partnerschaft mit 100 Runs.
- `` `Da pro Tag nur achtzehn Medaillen zur Verfügung stehen, hat es ein Anzahl an Ländern nicht auf das Podium geschafft.` ``
- `` `Wintersportarten sind in den nördlichen Regionen am beliebtesten und Italiener nehmen an internationalen Wettkämpfen und olympischen Spielen teil.` ``
- `` `Nach dem Rennen bleibt Keselowski mit 2.250 Punkten Spitzenreiter in der Fahrerwertung.` ``

Topics

- A. sports
- B. health
- C. travel
- D. science and technology

Your answer letter:

A.3 Further Details

Models. We test state-of-the-art instruction fine-tuned LVLMs across various sizes. Smaller models are evaluated on all languages, while larger LVLMs (26B+) are tested on subsets of the MVL-SIB languages (cf. §5.3).

GPT-4o. We evaluate MVL-SIB on GPT-4o-mini-2024-07-18 and GPT-4o-2024-08-06. We set the image detail in the API to ‘low’, since our tasks require high-level reasoning that does not depend on finer image details.¹¹ Prior works show that GPT-4o is the best-performing multilingual LVLM (Schneider and Sitaram, 2024; Vayani et al., 2024).

Qwen2-VL. Qwen2-VL ties a 675M parameter vision-transformer (ViT) into Qwen2 LLMs (Wang et al., 2024). An MLP compresses adjacent 2×2 visual tokens embedded by the ViT into one token representation, which is then input to the LLM.

InternVL 2.5. Depending on the model size, InternVL uses Qwen2.5 or InternLM as its LLM backbone (Chen et al., 2024). The model embeds images either by a 6B or by a distilled 300M ViT pretrained with CLIP (Radford et al., 2021). The resulting image patch encodings are downsampled by factor 4 and fed through an MLP to the LLM.

Centurio. Centurio is the latest massively multilingual LVLM trained on 100 languages (Geigle et al., 2025), outperforming alternatives like Parrot (Sun et al., 2024) or Pangea (Yue et al., 2025). It employs Qwen2.5 as its LLM (Yang et al., 2024) and SigLIP S0400/384 as its ViT (Zhai et al., 2023). The model mixes resolutions by stacking the encodings of the full image and those of 2×2 tiles along the features. The combined embedding is then projected via an MLP to the LLM’s input space.

Besides architectures, the LVLMs crucially differ in dataset mixtures on which they were trained. Centurio translates image-caption, VQA, OCR, and a few multi-image datasets to 100 languages with NLLB (Team et al., 2022) to mix 50:50 with the original English data. Qwen2-VL and InternVL, however, were trained on much larger, more diverse datasets that comprise sizable multi-image comparison and video understanding datasets. Moreover, assuming that the LLMs of Qwen2-VL, InternVL, and GPT-4o were pretrained on Flores, their performance would be overly optimistic.

¹¹We use GPT-4o-mini because evaluating GPT-4o would be too expensive.

A.3.1 Performance by Model over Languages grouped by Language Tier

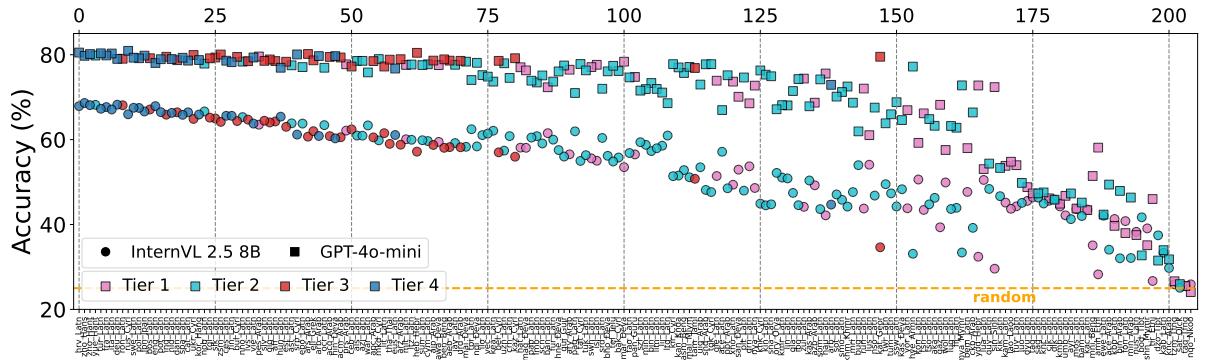


Figure 6: **Images-To-Sentences** @ $k=3$. The English prompt describes the cross-modal topic matching task, lists all topics, and provides both $k=3$ reference images and 4 sentences in the corresponding language $\{\text{eng_Latn}, \dots, \text{nqo_Nkoo}\}$. LVLMs must select the sentence of 4 options that topically fits $k=3$ reference images. The sentences spanning 205 languages and 7 topics are drawn from SIB-200 (Adelani et al., 2024), while images for the topics were hand-selected (cf. Appendix A.1). An example prompt is shown in Appendix A.7.2; further details are in §4.

Plot. The x-axis orders the languages of the candidate sentences $\{\text{eng_Latn}, \dots, \text{nqo_Nkoo}\}$, respectively, by descending performance (y-axis). The top x-axis indicates the running index of each language L_i ($i \in \{1, \dots, 205\}$).

Tiers. The languages are grouped by tiers derived from Joshi et al. (2020) (cf. §5).

A.3.2 Calibration Analysis of Cross-Modal Topic Matching

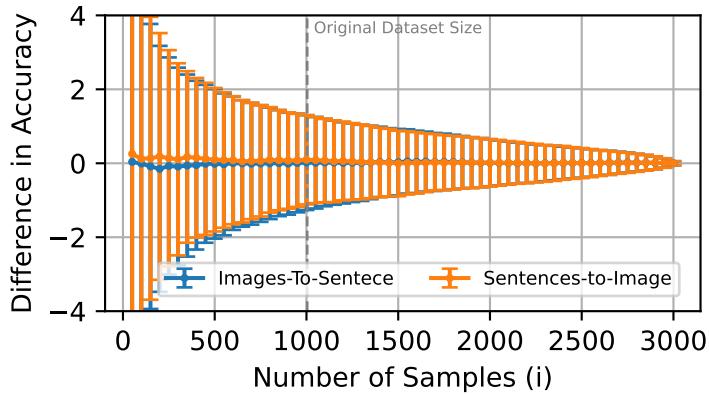


Figure 7: **Calibration Analysis of Cross-Modal Topic Matching for InternVL 2.5 8B.** **Analysis:** To assess the reliability of cross-modal topic matching with fewer samples than our full dataset (1004 samples per language), we randomly select 500 trajectories. We then compute performance metrics on cumulative subsets, incrementing by 50 examples at each step. The difference in performance between the full dataset and each subset is calculated to quantify the deviation at each sample size. **Plot:** The plot displays the average absolute spread in performance (averaged over all languages) along with the standard deviation for InternVL 2.5 8B. We restrict ourselves to a single open-weight LVLM, since the analysis yields identical results across all combinations of LVLMs and language tiers.

Insights: Our analysis shows that performance stabilizes rapidly, with deviations of only about 1% observed at 1,004 instances – same size as the subset from which the dataset was created. This indicates that reliable evaluation of cross-modal topic matching can be achieved with far fewer than 3,012 samples.

A.4 Overview of Multilingual Vision-Language Benchmarks

xGQA. The xGQA dataset (Pfeiffer et al., 2022) is a cross-lingual visual question-answering resource. It extends the well-known English-only GQA dataset (Hudson and Manning, 2019) by providing manual translations of the questions in the balanced *test-dev* set. The dataset contains 9666 questions available in eight languages across five scripts, while the answers remain in English. In addition, it features 300 unique images from Visual Genome (Krishna et al., 2017).

MaXM. MaXM, introduced by Changpinyo et al. (2023), is a VQA dataset covering seven languages written in five scripts. In this dataset, both the questions and their corresponding answers are presented in the same language. The images are drawn from a subset of the XM3600 (Thapliyal et al., 2022) dataset and are selected to correspond to regions where the question-answer pair’s language is spoken, ensuring both linguistic and cultural diversity.

XVNLI. The XVNLI dataset (Bugliarello et al., 2022) introduces the task of Cross-lingual Visual Natural Language Inference, where a model must determine if a textual hypothesis *entails*, *contradicts*, or is *neutral* with respect to a visual premise. This dataset spans five languages across three scripts and includes 357 unique images from Visual Genome. It is built upon a combination of the text-only SNLI (Bowman et al., 2015) dataset and its cross-lingual (Agić and Schluter, 2018) and cross-modal (Xie et al., 2019) counterparts.

MaRVL. The MaRVL dataset (Liu et al., 2021) is designed to benchmark models on Multicultural Reasoning over Vision and Language. Each sample consists of two images, a textual statement, and a binary (true/false) answer grounded in the images. Covering five languages across three scripts, MaRVL includes 4914 culturally diverse images that align with the respective languages. The images in each sample are selected to reflect the culture of the annotator who composed the textual statement in their native language.

XM3600. The XM3600 dataset (Thapliyal et al., 2022) is an extensive multilingual image captioning resource encompassing 36 languages. It contains 261375 captions across 13 scripts, with 100 unique images per language. The images are chosen to reflect the cultural background of the language, ensuring both cultural and linguistic diversity. All

captions were manually produced by professional native speakers rather than being automatically generated. Due to the dataset’s large size, we evaluate XM3600 using a randomly selected subset of 500 images per language.

Babel-ImageNet (multiple-choice) (BIN-MC). Babel-ImageNet (Geigle et al., 2024c) translates ImageNet’s (Deng et al., 2009) labels into nearly 300 languages, allowing us to assess whether models can recognize and correctly link diverse ImageNet objects to their labels in the target language. Given the computational cost, we focus on languages that appear in only one or two other datasets, in addition to English and a select few high-resource languages, and we use 10 images per class instead of 50. We follow Geigle et al. (2024b) and frame the task as a multiple-choice problem by mining hard negative options from the complete label pool. This approach avoids the ambiguity inherent in traditional open-ended VQA formats. Negatives are selected based on the English labels, filtering out candidates not translated by Babel-ImageNet into the target language, and ultimately choosing the three most similar negative labels available.

SMPQA. Geigle et al. (2025) introduce SMPQA (Synthetic Multilingual Plot QA) as a test dataset for evaluating multilingual OCR capabilities for bar plots and pie charts, covering 11 languages and various scripts and resource levels.

M5B-VGR. The M5B-VGR dataset, presented by (Schneider and Sitaram, 2024), is a visually grounded reasoning benchmark akin to MaRVL. Each sample comprises two images, a textual statement, and a binary (true/false) answer based on the images. It spans 12 languages across 7 scripts and features culturally diverse photos from regions where the respective languages are spoken. The images are sampled from the Dollar Street (Gaviria Rojas et al., 2022) dataset, with 120 samples provided per language.

M5B-VLOD. The M5B-VLOD (Visio-Linguistic Outlier Detection) dataset, also introduced by (Schneider and Sitaram, 2024), consists of samples containing five images paired with a textual statement that is true for all but one image. The task is to identify the outlier image that does not match the statement. This dataset covers the same 12 languages as M5B-VGR, with images sampled in a similar manner from the same source, and provides

Task	Dataset	Visual Input	Textual Input	Target Output	Metric	#Lang.
Captioning	XM3600	Single-Image	Prompt (English)	Caption (Target Language)	CIDEr	36
Multiple-Choice Visual Question Answering	BabelImageNet-MC	Single-Image	Question (Target Language)	Letter of the correct Choice	Relaxed Accuracy	20
Text-Heavy Multiple-Choice Visual Question Answering	M3Exam MMMU xMMMU	Single or Multi-Image	Question (Target Language) Question (English) Question (Target Language)	Letter of the correct Choice	Relaxed Accuracy	7 1 7
Text-Heavy Visual Question Answering	MTVQA SMPQA - Name	Single-Image	Question (Target Language)	Word or Phrase (Target Language)	Exact Accuracy	9 11
Text-Heavy Visually Grounded Reasoning	SMPQA - Ground	Single-Image	Question (Target Language)	'yes' / 'no'	Exact Accuracy	11
Visio-Linguistic Outlier Detection	M5B-VLOD	Multi-Image	Hypothesis (Target Language)	Letter of the correct Choice	Relaxed Accuracy	12
Visual Natural Language Inference	XVNLI	Single-Image	Hypothesis (Target Language)	'yes' / 'no' / 'maybe'	Exact Accuracy	5
Visual Question Answering	MaXM xGQA	Single-Image	Question (Target Language)	Word or Phrase (Target Language) Word or Phrase (English)	Exact Accuracy	6 8
Visually Grounded Reasoning	M5B-VGR MaRVL	Multi-Image	Hypothesis (Target Language)	'yes' / 'no'	Exact Accuracy	12 6

Table 3: Summary of multilingual vision-language benchmarks we correlate MVL-SIB against. Relaxed denotes responses that start with the correct option letter (cf. 4).

120 samples per language.

MTVQA. The MTVQA dataset, introduced by (Tang et al., 2024), features text-heavy visual question answering tasks. It includes expert human annotations in 9 diverse languages, comprising 6778 question-answer pairs across 2116 images. The images predominantly contain text in the corresponding language, with questions and answers closely tied to that text. These images are sourced from various publicly available datasets.

CVQA. The CVQA dataset, introduced by (Romero et al., 2025), is a multilingual and culturally nuanced VQA benchmark that includes a broad array of languages, many of which are underrepresented in NLP. It consists of 10000 questions spanning 30 countries and 31 languages, forming 39 distinct country-language pairs (for instance, Spanish appears in 7 different splits corresponding to 7 Spanish-speaking countries). The images were manually collected by human annotators to accurately depict the culture associated with each country-language pair. Each sample includes one image and a question in the respective language. Although the test set is not publicly available, the authors permit up to 5 daily leaderboard submissions for evaluation.

M3Exam. The M3Exam dataset, presented by (Zhang et al., 2023), contains real-world exam questions in 9 languages, available as either text-only or multimodal samples. For our evaluation, we only include samples that require at least one image. Due to the limited number of samples for Swahili and Javanese, we focus on the remaining 7 languages. The selected samples consist of

multiple-choice questions in the target language, accompanied by up to 8 images that may appear in both the question and the answer options, with the number of choices ranging from 4 to 8 per sample. **xMMMU.** xMMMU, introduced by (Yue et al., 2025), comprises college-level multiple-choice VQA samples in seven languages. It was automatically translated using GPT4o from a randomly selected subset of 300 questions from the MMMU (Yue et al., 2024) validation split.

A.5 Prompts

We list the prompts for each dataset in our test suite used for all models in Figure 8.

SMPQA

{QUESTION}\nAnswer the question using a single word or phrase.

CVQA

{QUESTION}\nThere are several options:\nA. {OPTION A}\nB. {OPTION B}\nC. {OPTION C}\nD. {OPTION D}\nAnswer with the option's letter from the given choices directly.

xMMMU

{QUESTION}\nThere are several options:\nA. {OPTION A}\nB. {OPTION B}\nC. {OPTION C}\nD. {OPTION D}\nAnswer with the option's letter from the given choices directly.

MTVQA

{QUESTION}\nAnswer the question using a single word or phrase.\nAnswer in {LANGUAGE}.

M3Exam

{QUESTION}\nOptions:\nA. {OPTION A}\nB. {OPTION B}\nC. {OPTION C}\nD. {OPTION D}\nAnswer with the option's letter from the given choices directly.

BIN-MC

Which of these choices (in English) is shown in the image?\nChoices:\nA. {CHOICE A}\nB. {CHOICE B}\nC. {CHOICE C}\nD. {CHOICE D}\nAnswer with the letter from the given choices directly.

xGQA

{QUESTION}\nAnswer the question using a single word or phrase.\nAnswer in English.

MaXM

{QUESTION}\nAnswer the question using a single word or phrase.\nAnswer in {LANGUAGE}.

MaRVL

Given the two images , is it correct to say “{HYPOTHESIS}”? Answer yes or no.’

XVNLI

Is it guaranteed true that “{HYPOTHESIS}”? Yes, no, or maybe? Answer in English.

M5-VGR

Given the two images , is it correct to say “{HYPOTHESIS}”? Answer yes or no.’

M5-VLOD

Based on the 5 images ordered from top-left to bottom-right, which image does not match the hypothesis “{HYPOTHESIS}”? Choose one from [A, B, C, D, E] and only output a single letter:

XM3600

Briefly describe the image in {LANGUAGE} in one sentence.

Figure 8: Prompts used for the different datasets of our test suite. For M3Exam and xMMMU, the questions contain images at individual positions, and also the options can consist of images. In total, a sample of M3Exam can contain up to 8 images and 8 options, and a sample of xMMMU can contain up to 4 images and 4 options.

A.6 Full Results For Subsets by Task, Model, and Language

Lang	Tier	Qwen2-VL			InternVL 2.5					Centurio Qwen		GPT-4o	
		Qwen2-VL 2B	Qwen2-VL 7B	Qwen2-VL 72B	InternVL 2.5 4B	InternVL 2.5 8B	InternVL 2.5 26B	InternVL 2.5 38B	InternVL 2.5 78B	Centurio Qwen 8B	GPT-4o-mini	GPT-4o	
acm_Arab	3	34.9	55.4	80.8	50.0	61.0	74.1	80.2	82.2	54.4	72.0	N/A	
aka_Latn	2	28.7	38.0	57.2	34.6	46.3	51.3	58.3	57.0	42.4	61.4	N/A	
apc_Arab	3	34.9	54.8	80.0	49.1	61.3	74.2	78.8	82.1	54.4	71.5	N/A	
arb_Arab	4	35.2	57.1	81.4	51.7	62.2	75.8	80.8	82.8	55.2	73.2	N/A	
azj_Latn	2	34.6	53.8	80.5	49.5	58.6	70.9	78.3	80.0	51.8	72.8	N/A	
bul_Cyr1	3	35.0	56.6	79.3	50.3	65.3	77.1	79.8	81.6	53.9	71.5	N/A	
ces_Latn	4	36.2	58.2	79.5	52.0	65.6	78.5	79.5	81.5	53.2	72.7	N/A	
eng_Latn	4	36.3	65.8	79.8	52.5	67.7	N/A	79.7	81.7	54.8	68.3	N/A	
fin_Latn	4	33.0	57.1	79.1	49.8	64.1	76.5	79.2	81.1	54.1	70.2	N/A	
hat_Latn	1	32.1	50.4	77.8	43.0	54.7	66.1	74.1	74.2	54.9	69.3	N/A	
hau_Latn	1	29.1	39.2	56.3	32.4	41.8	45.2	58.0	57.5	43.9	70.2	N/A	
min_Arab	2	27.0	30.4	58.4	27.8	33.6	36.3	53.4	47.8	37.1	44.8	N/A	
umb_Latn	1	29.1	34.8	52.5	31.1	42.5	45.2	50.8	52.0	35.4	47.4	N/A	

Table 4: Subsets of language tiers for I2S @ k=1.

Lang	Tier	Qwen2-VL			InternVL 2.5					Centurio Qwen		GPT-4o	
		Qwen2-VL 2B	Qwen2-VL 7B	Qwen2-VL 72B	InternVL 2.5 4B	InternVL 2.5 8B	InternVL 2.5 26B	InternVL 2.5 38B	InternVL 2.5 78B	Centurio Qwen 8B	GPT-4o-mini	GPT-4o	
acm_Arab	3	35.2	55.1	82.4	45.6	60.6	74.0	85.5	86.7	59.7	79.3	86.0	
aka_Latn	2	27.9	38.3	63.2	34.2	46.3	52.9	61.1	60.6	42.6	63.3	79.3	
apc_Arab	3	34.7	54.3	82.2	45.1	60.5	73.9	84.1	85.7	59.4	78.5	85.8	
arb_Arab	4	35.2	56.4	81.8	46.9	60.9	76.0	85.0	87.1	58.7	79.7	86.7	
azj_Latn	2	33.1	53.4	82.3	46.2	60.9	71.4	82.7	83.5	57.9	78.6	86.5	
bul_Cyr1	3	34.7	54.6	82.7	45.9	64.4	78.8	84.8	84.5	58.8	79.5	86.4	
ces_Latn	4	35.6	56.8	82.2	48.9	65.6	78.4	83.6	85.3	57.3	78.6	86.3	
eng_Latn	4	34.8	63.1	82.6	49.2	67.9	80.1	84.9	85.5	60.0	78.1	87.1	
fin_Latn	4	33.5	56.4	81.9	45.8	65.6	77.2	83.2	84.4	59.7	78.3	86.5	
hat_Latn	1	30.7	52.7	78.9	42.1	58.1	67.9	77.6	79.3	58.0	77.1	85.1	
hau_Latn	1	28.7	39.0	61.6	30.2	42.2	46.8	59.8	59.8	44.5	75.6	85.1	
min_Arab	2	27.4	32.6	61.0	28.9	32.1	36.4	55.9	50.2	35.9	46.3	76.1	
umb_Latn	1	28.9	34.8	58.7	30.5	42.2	45.8	52.6	53.5	35.8	46.8	61.1	

Table 5: Subsets of language tiers for I2S @ k=3.

Lang	Tier	Qwen2-VL			InternVL 2.5					Centurio Qwen		GPT-4o	
		Qwen2-VL 2B	Qwen2-VL 7B	Qwen2-VL 72B	InternVL 2.5 4B	InternVL 2.5 8B	InternVL 2.5 26B	InternVL 2.5 38B	InternVL 2.5 78B	Centurio Qwen 8B	GPT-4o-mini	GPT-4o	
acm_Arab	3	33.7	48.7	83.2	47.6	60.6	73.9	85.8	87.1	61.1	77.3	N/A	
aka_Latn	2	27.0	34.0	64.8	34.5	47.0	52.6	60.6	60.3	43.2	60.7	N/A	
apc_Arab	3	32.7	48.9	83.9	46.2	61.2	73.6	84.7	87.3	60.3	77.4	N/A	
arb_Arab	4	33.3	50.6	83.2	49.3	62.3	75.6	85.9	87.7	60.2	78.2	N/A	
azj_Latn	2	32.2	47.9	82.2	45.5	63.0	72.8	83.6	85.3	59.4	78.1	N/A	
bul_Cyr1	3	33.2	50.1	83.1	48.5	66.0	77.6	85.2	86.5	60.2	78.3	N/A	
ces_Latn	4	34.9	52.2	82.2	48.6	66.0	78.8	84.7	86.4	58.9	79.3	N/A	
eng_Latn	4	34.9	58.9	83.9	48.1	68.7	N/A	84.5	87.0	62.4	77.4	N/A	
fin_Latn	4	32.6	51.1	82.2	46.8	66.6	76.3	83.7	86.3	60.7	77.0	N/A	
hat_Latn	1	31.1	49.5	80.0	41.3	59.4	68.7	77.0	80.4	58.5	75.6	N/A	
hau_Latn	1	28.4	35.9	63.3	31.4	42.8	46.2	60.5	60.4	44.7	73.9	N/A	
min_Arab	2	27.4	30.1	62.0	28.5	32.3	36.3	55.7	49.8	36.4	43.9	N/A	
umb_Latn	1	28.0	31.7	59.1	31.6	42.7	46.6	52.3	53.3	35.7	45.4	N/A	

Table 6: Subsets of language tiers for I2S @ k=5.

Lang	Tier	Qwen2-VL			InternVL 2.5					Centurio Qwen		GPT-4o	
		Qwen2-VL 2B	Qwen2-VL 7B	Qwen2-VL 72B	InternVL 2.5 4B	InternVL 2.5 8B	InternVL 2.5 26B	InternVL 2.5 38B	InternVL 2.5 78B	Centurio Qwen 8B	GPT-4o-mini	GPT-4o	
acm_Arab	3	46.6	81.1	91.1	71.5	80.5	84.9	91.3	89.6	83.2	89.9	92.3	
aka_Latn	2	32.3	53.2	68.1	44.4	59.6	64.4	68.6	66.5	63.5	79.9	87.7	
apc_Arab	3	46.0	80.9	91.3	70.0	80.4	84.6	90.9	89.1	82.7	89.4	92.5	
arb_Arab	4	48.1	82.7	91.4	73.1	81.3	86.5	91.7	90.2	84.9	89.9	92.8	
azj_Latn	2	42.6	80.2	90.2	65.5	76.9	82.0	89.1	87.3	82.0	89.1	91.9	
bul_Cyr1	3	47.1	83.3	91.2	73.4	83.5	88.2	90.6	89.2	83.9	89.3	92.5	
ces_Latn	4	50.8	82.8	90.6	75.3	85.3	89.8	90.8	90.0	84.6	90.1	92.2	
eng_Latn	4	56.7	85.7	91.8	81.4	87.0	91.3	91.8	91.5	85.4	88.5	92.0	
fin_Latn	4	42.3	82.0	91.4	70.5	82.3	87.1	90.2	88.1	81.8	89.2	92.3	
hat_Latn	1	40.2	70.8	87.3	57.9	74.1	79.2	84.9	83.7	82.0	87.3	91.6	
hau_Latn	1	32.6	52.1	66.5	40.0	53.8	57.5	70.7	66.2	66.5	67.3	91.9	
min_Arab	2	28.3	42.3	68.0	32.8	41.5	41.4	63.7	56.4	49.8	62.8	86.4	
umb_Latn	1	30.3	47.1	61.4	39.3	52.4	57.2	N/A	59.6	53.4	62.9	71.7	

Table 7: Subsets of language tiers for T2S @ k=1.

Lang	Tier	Qwen2-VL			InternVL 2.5					Centurio Qwen		GPT-4o	
		Qwen2-VL 2B	Qwen2-VL 7B	Qwen2-VL 72B	InternVL 2.5 4B	InternVL 2.5 8B	InternVL 2.5 26B	InternVL 2.5 38B	InternVL 2.5 78B	Centurio Qwen 8B	GPT-4o-mini	GPT-4o	
aeb_Arab	3	43.1	62.9	80.1	36.8	53.2	68.7	79.3	77.6	29.8	76.3	N/A	
arb_Arab	4	44.3	66.5	80.2	37.7	55.2	71.9	81.4	79.9	30.8	77.7	N/A	
ben_Beng	3	39.8	64.4	81.2	41.2	50.9	69.0	79.1	77.2	33.0	76.4	N/A	
eng_Latn	4	41.9	71.7	N/A	47.7	66.2	N/A	81.7	81.3	35.3	77.5	N/A	
fao_Latn	2	33.6	52.2	75.1	30.9	43.9	56.8	73.6	69.0	29.2	75.8	N/A	
kac_Latn	1	28.8	35.4	51.5	26.4	37.8	44.7	51.6	46.9	26.7	46.6	N/A	
kas_Deva	2	33.3	47.1	72.4	30.2	42.4	53.0	69.0	59.3	27.7	66.5	N/A	
lit_Latn	3	38.1	63.6	78.9</									

Qwen2-VL			InternVL 2.5					Centurio Qwen		GPT-4o		
Lang	Tier	Qwen2-VL 2B	Qwen2-VL 7B	Qwen2-VL 72B	InternVL 2.5 4B	InternVL 2.5 8B	InternVL 2.5 26B	InternVL 2.5 38B	InternVL 2.5 78B	Centurio Qwen 8B	GPT-4o-mini	GPT-4o
aeb_Arab	3	44.1	65.3	89.6	38.5	58.6	76.0	92.0	89.0	31.7	85.4	88.3
arb_Arab	4	44.8	67.6	89.1	39.2	59.8	78.9	92.4	90.1	32.8	86.8	89.1
ben_Beng	3	40.4	67.0	88.9	43.4	58.1	74.5	90.8	89.2	33.9	85.7	87.9
eng_Latn	4	43.1	70.4	88.5	44.5	69.0	84.0	92.1	90.0	36.1	86.4	89.1
fao_Latn	2	35.0	55.2	87.5	34.9	50.6	68.9	89.5	85.9	30.5	86.0	89.1
kac_Latn	1	31.6	41.0	69.7	29.2	46.3	57.5	72.3	68.0	27.4	64.2	70.5
kas_Deva	2	35.1	50.1	87.0	31.9	48.9	60.6	86.5	78.0	28.9	80.4	86.1
lit_Latn	3	42.1	64.7	89.1	39.2	58.5	77.4	91.0	89.4	32.2	86.2	89.0
lua_Latn	1	33.7	46.2	74.6	33.0	50.8	61.5	78.5	77.1	27.8	70.4	79.3
mal_Mlym	2	33.7	65.3	88.0	38.1	50.8	69.7	90.7	88.2	31.0	85.1	87.8
srp_Cyr1	4	44.1	65.1	88.7	39.7	59.4	78.3	91.8	89.8	32.6	87.1	89.5
tur_Latn	4	44.6	66.7	88.3	41.3	61.0	81.7	91.8	88.8	34.3	87.0	88.5
wol_Latn	1	34.8	45.8	75.2	32.9	49.7	62.3	79.2	77.0	28.1	72.7	83.1

Table 9: Subsets of language tiers for S2I @ k=3.

Qwen2-VL			InternVL 2.5					Centurio Qwen		GPT-4o		
Lang	Tier	Qwen2-VL 2B	Qwen2-VL 7B	Qwen2-VL 72B	InternVL 2.5 4B	InternVL 2.5 8B	InternVL 2.5 26B	InternVL 2.5 38B	InternVL 2.5 78B	Centurio Qwen 8B	GPT-4o-mini	GPT-4o
aeb_Arab	3	44.0	65.4	91.6	39.0	57.6	77.0	93.6	91.8	32.8	88.9	N/A
arb_Arab	4	44.3	65.7	90.8	40.0	59.9	79.4	94.0	91.9	33.5	88.7	N/A
ben_Beng	3	41.1	63.2	89.3	45.9	59.4	76.4	93.0	91.5	34.8	88.4	N/A
eng_Latn	4	43.4	68.6	89.6	43.0	68.7	N/A	93.7	91.7	35.6	89.1	N/A
fao_Latn	2	35.3	56.6	90.4	36.0	50.8	71.8	92.8	89.5	30.7	87.8	N/A
kac_Latn	1	31.8	42.3	74.8	31.6	46.2	60.6	79.2	73.9	27.0	68.4	N/A
kas_Deva	2	35.3	48.8	89.0	34.1	50.0	64.7	91.3	82.9	29.4	83.2	N/A
lit_Latn	3	42.1	63.5	91.1	39.5	57.7	78.1	93.8	92.1	32.8	88.2	N/A
lua_Latn	1	34.6	48.4	80.2	34.6	49.7	63.1	85.1	82.8	28.1	75.9	N/A
mal_Mlym	2	31.6	62.7	88.9	39.0	51.8	69.3	93.2	91.7	31.5	87.1	N/A
srp_Cyr1	4	44.5	63.5	90.7	41.2	58.4	79.4	93.6	91.4	33.3	89.1	N/A
tur_Latn	4	45.6	64.8	90.3	40.9	59.9	81.4	93.3	91.3	34.8	88.5	N/A
wol_Latn	1	34.4	46.3	82.1	34.8	50.9	64.3	86.1	82.2	27.9	77.9	N/A

Table 10: Subsets of language tiers for S2I @ k=5.

Qwen2-VL			InternVL 2.5					Centurio Qwen		GPT-4o		
Lang	Tier	Qwen2-VL 2B	Qwen2-VL 7B	Qwen2-VL 72B	InternVL 2.5 4B	InternVL 2.5 8B	InternVL 2.5 26B	InternVL 2.5 38B	InternVL 2.5 78B	Centurio Qwen 8B	GPT-4o-mini	GPT-4o
aeb_Arab	3	77.9	85.2	91.0	87.5	83.5	83.9	91.5	91.9	88.8	91.5	N/A
arb_Arab	4	81.5	86.5	91.7	88.9	86.1	88.0	92.8	93.4	89.8	92.7	N/A
ben_Beng	3	68.8	85.2	91.2	83.7	79.5	81.6	91.2	90.8	86.4	91.1	N/A
eng_Latn	4	86.1	89.1	91.0	90.6	91.7	92.0	92.4	93.5	89.9	92.4	N/A
fao_Latn	2	52.1	72.1	86.3	71.0	72.8	76.9	86.5	83.9	80.3	90.0	N/A
kac_Latn	1	41.3	51.5	58.8	55.2	54.6	55.7	57.4	58.0	54.7	58.2	N/A
kas_Deva	2	50.1	72.4	83.5	71.7	65.1	63.4	80.4	74.2	75.0	82.9	N/A
lit_Latn	3	67.2	81.7	90.1	76.6	77.7	84.6	89.6	89.6	84.2	91.1	N/A
lua_Latn	1	49.4	58.5	65.2	61.3	61.0	59.0	64.8	68.5	63.3	65.8	N/A
mal_Mlym	2	49.0	83.0	89.7	76.9	69.1	74.3	88.7	87.2	80.6	90.5	N/A
srp_Cyr1	4	77.3	86.1	91.0	84.8	85.6	87.0	91.5	92.3	88.2	91.7	N/A
tur_Latn	4	76.4	85.0	91.3	85.1	87.5	89.4	91.7	91.6	88.2	91.5	N/A
wol_Latn	1	50.4	59.1	66.8	63.3	64.0	62.0	66.6	67.6	65.5	69.5	N/A

Table 11: Subsets of language tiers for S2T @ k=1.

Qwen2-VL			InternVL 2.5					Centurio Qwen		GPT-4o		
Lang	Tier	Qwen2-VL 2B	Qwen2-VL 7B	Qwen2-VL 72B	InternVL 2.5 4B	InternVL 2.5 8B	InternVL 2.5 26B	InternVL 2.5 38B	InternVL 2.5 78B	Centurio Qwen 8B	GPT-4o-mini	GPT-4o
aeb_Arab	3	94.1	93.9	98.7	97.2	95.3	96.3	98.8	99.0	96.6	98.4	98.4
arb_Arab	4	94.7	94.4	98.7	97.3	96.7	98.0	99.7	99.1	95.7	98.6	98.4
ben_Beng	3	86.9	93.9	98.7	95.3	93.0	94.7	98.6	98.4	96.5	98.5	98.3
eng_Latn	4	96.5	95.3	98.3	98.0	98.1	98.7	98.7	98.9	96.7	98.1	98.0
fao_Latn	2	69.4	86.6	97.6	88.7	87.1	92.6	97.2	97.2	94.1	98.2	98.3
kac_Latn	1	57.6	68.0	80.1	75.0	77.7	75.5	79.4	79.9	75.3	78.6	79.1
kas_Deva	2	68.6	88.4	96.5	86.9	81.2	82.3	95.4	91.4	90.5	95.9	96.7
lit_Latn	3	87.1	92.4	98.6	92.2	94.9	97.6	98.2	98.3	94.9	98.6	98.2
lua_Latn	1	66.7	76.9	86.2	82.1	83.0	82.1	86.5	88.3	81.4	87.4	92.2
mal_Mlym	2	61.1	93.3	98.7	91.2	85.5	90.7	98.3	97.7	94.0	98.3	98.1
srp_Cyr1	4	93.2	94.4	98.7	96.2	96.5	97.6	98.7	98.9	96.5	98.3	98.4
tur_Latn	4	92.5	94.9	98.8	96.6	97.2	98.6	99.6	99.0	96.0	98.6	98.1
wol_Latn	1	66.5	77.1	88.9	84.5	84.9	92.7	92.9	95.2	91.1	89.3	95.1

Table 12: Subsets of language tiers for S2T @ k=3.

Qwen2-VL			InternVL 2.5					Centurio Qwen		GPT-4o		
Lang	Tier	Qwen2-VL 2B	Qwen2-VL 7B	Qwen2-VL 72B	InternVL 2.5 4B	InternVL 2.5 8B	InternVL 2.5 26B	InternVL 2.5 38B	InternVL 2.5 78B	Centurio Qwen 8B	GPT-4o-mini	GPT-4o
aeb_Arab	3	97.2	97.4	99.7	98.7	97.7	98.8	99.7	99.7	98.3	99.5	N/A
arb_Arab	4	97.6	97.4	99.6	98.6	98.5	99.3	99.7	99.7	97.6	99.3	N/A
ben_Beng	3	92.6	96.9	99.5	97.6	96.8	97.9	99.6	99.5	97.7	99.2	N/A
eng_Latn	4	98.2	97.5	99.2	99.1	99.0	99.5	99.4	99.6	97.7	99.1	N/A
fao_Latn	2	79.9	91.8	99.4	92.7	92.9	96.3	99.2	99.5	96.9	99.0	N/A
kac_Latn	1	68.0	75.0	87.4	81.5	86.0	83.9	87.1	87.3	80.9	86.3	N/A
kas_Deva	2	78.4	92.4	99.0	92.3	87.9	87.5	99.0	96.6	93.8	98.7	N/A
lit_Latn	3	93.8	96.0	99.4	95.9	98.2	99.3	99.4	99.7	97.2	99.4	N/A
lua_Latn	1	75.9	83.3	91.9	88.1	90.1	89.1	92.4	93.5	85.2	93.1	N/A
mal_Mlym	2	67.6	96.7	99.6	95.5	90.2	95.7					

A.7 Full Performance by Task, Model, and Language

A.7.1 Images-To-Topics

Topics

Model	Entertainment	Geography	Health	Politics	Science & Tech.	Sports	Travel
QwenVL-2.5-2B	99.4	83.1	100.0	99.9	100.0	100.0	95.7
QwenVL-2.5-7B	100.0	92.9	100.0	100.0	99.8	100.0	100.0
InternVL-2.5-4B	99.5	84.5	100.0	100.0	99.7	100.0	100.0
InternVL-2.5-8B	100.0	90.0	100.0	100.0	97.5	100.0	100.0
Centurio-Qwen	99.6	90.6	100.0	100.0	99.2	100.0	100.0

Table 14: **Image-To-Topics**. For a reference image, the model must pick the correct topic out of 4 choices.

A.7.2 Images-To-Sentences

Lang.	QwenVL-2.2B					QwenVL-2.7B					InternVL-2.54B					InternVL-2.5.5B					CentriQ-Wen					40-mini				
	1	3	5	1	3	5	1	3	5	1	3	5	1	3	5	1	3	5	1	3	5	1	3	5	1	3	5	1	3	5
avg	31.8	31.4	30.7	46.6	46.4	42.4	40.9	39.0	39.3	53.1	53.3	53.7	47.8	50.5	51.1	64.2	69.2	67.7												
eng_Eng	36.3	34.8	34.5	63.6	58.1	52.5	49.2	48.1	67.7	67.9	67.8	54.8	60.9	62.4	68.3	78.1	77.4													
ace_Arab	27.7	27.9	27.4	33.4	33.3	31.7	29.5	28.9	29.4	34.6	34.1	30.4	38.9	37.5	38.3	48.4	49.4	46.6												
ace_Arab	34.9	35.2	33.7	55.4	55.1	48.7	46.0	45.8	46.0	60.6	60.6	54.4	59.7	61.1	72.0	79.3	77.3													
aca_Arab	35.2	35.0	32.8	56.4	55.2	50.0	50.8	45.8	48.3	61.8	60.8	61.0	54.1	59.1	60.9	72.6	79.5	77.3												
aca_Arab	34.6	33.5	32.2	53.2	52.0	46.7	48.2	44.2	51.5	59.2	58.3	59.3	55.4	60.9	61.3	71.8	78.6	76.3												
aca_Arab	35.2	34.1	33.1	55.9	55.4	49.2	49.8	46.2	48.0	61.6	58.7	60.8	53.9	59.2	60.5	72.6	78.6	76.8												
aca_Latn	28.7	27.0	27.0	38.0	38.3	34.0	34.4	34.2	35.4	46.5	46.3	47.0	42.4	42.6	43.2	61.4	63.0	63.7												
als_Arab	33.4	33.5	32.4	54.9	54.4	51.0	51.0	43.9	42.3	42.9	50.1	61.0	62.8	51.1	56.6	58.3	52.7	78.5	77.4											
ara_Arab	34.9	34.7	32.7	54.8	54.3	48.9	49.1	45.1	46.2	61.3	60.5	61.2	54.4	59.4	60.3	71.6	78.5	77.3												
arb_Arab	35.2	35.2	33.3	57.1	56.4	50.4	51.7	46.9	49.3	62.2	60.9	62.3	55.2	58.7	60.2	73.0	79.7	78.2												
arb_Arab	28.7	28.7	28.7	34.7	34.7	34.0	35.5	34.0	35.5	43.5	44.7	46.3	46.2	46.3	47.0	62.8	72.9	71.6												
arb_Arab	34.9	34.9	34.9	56.4	56.4	50.4	50.4	45.3	45.3	59.3	59.3	59.3	55.4	60.4	60.4	72.8	79.5	77.3												
ary_Arab	33.8	32.4	31.6	52.2	50.7	45.7	45.8	43.1	45.5	56.8	56.5	57.5	54.0	57.4	58.4	70.3	76.5	74.9												
ary_Arab	35.2	33.9	33.3	56.2	55.6	49.5	49.0	45.5	46.9	61.4	58.8	61.6	55.1	59.0	61.3	72.9	79.0	77.9												
ary_Arab	34.9	34.9	34.9	55.1	55.1	49.4	49.8	46.2	48.0	61.6	58.7	60.8	53.9	59.2	60.5	72.6	78.6	76.8												
awa_Deva	34.4	34.5	33.0	48.1	46.5	41.8	41.8	43.4	40.9	59.9	63.0	59.5	57.4	53.8	58.7	58.7	70	77.6	76.5											
ayr_Arab	29.0	28.0	28.2	33.3	33.5	31.1	30.5	30.3	29.7	40.6	41.3	41.9	35.4	34.5	35.7	44.0	47.3	46.1												
azb_Arab	34.5	34.5	34.5	54.3	54.3	49.3	49.3	45.1	45.1	59.1	58.5	58.6	55.6	58.0	59.0	61.3	62.8	62.8												
azb_Arab	34.6	33.1	32.2	53.8	53.4	47.9	47.9	44.1	45.2	60.5	60.0	61.0	55.8	59.0	59.0	72.8	78.6	78.1												
bak_Cyril	32.8	32.8	32.8	49.3	48.8	42.7	41.5	40.5	39.3	54.9	55.8	57.0	49.4	51.2	53.1	71.0	76.2	75.0												
bak_Cyril	28.0	27.7	26.9	34.7	34.3	32.4	32.3	31.2	31.5	41.9	42.0	42.6	40.5	40.2	41.4	43.4	52.5	59.3	58.3											
bel_Cyril	34.1	34.5	33.7	53.2	52.7	49.1	47.2	45.9	43.6	53.8	57.7	58.0	51.9	57.3	58.2	73.0	78.6	78.2												
bel_Cyril	28.7	28.9	28.6	39.5	36.2	40.2	46.6	35.4	34.7	34.9	46.4	47.7	47.3	39.9	40.3	40.9	56.0	58.0	55.4											
ben_Beng	32.8	33.6	32.2	48.2	48.5	46.1	47.0	44.3	43.3	62.2	58.1	54.7	51.7	55.7	54.3	72.8	78.0	77.1												
ben_Beng	34.9	34.9	34.9	54.7	54.7	48.1	48.1	45.3	45.3	62.8	60.6	58.1	54.7	58.1	58.1	72.8	78.0	77.1												
ben_Beng	35.0	34.9	34.9	54.7	54.7	48.1	48.1	45.3	45.3	62.8	60.6	58.1	54.7	58.1	58.1	72.8	78.0	77.1												
ben_Beng	34.9	34.9	34.9	54.7	54.7	48.1	48.1	45.3	45.3	62.8	60.6	58.1	54.7	58.1	58.1	72.8	78.0	77.1												
ben_Beng	34.9	34.9	34.9	54.7	54.7	48.1	48.1	45.3	45.3	62.8	60.6	58.1	54.7	58.1	58.1	72.8	78.0	77.1												
ben_Beng	34.9	34.9	34.9	54.7	54.7	48.1	48.1	45.3	45.3	62.8	60.6	58.1	54.7	58.1	58.1	72.8	78.0	77.1												
ben_Beng	34.9	34.9	34.9	54.7	54.7	48.1	48.1	45.3	45.3	62.8	60.6	58.1	54.7	58.1	58.1	72.8	78.0	77.1												
ben_Beng	34.9	34.9	34.9	54.7	54.7	48.1	48.1	45.3	45.3	62.8	60.6	58.1	54.7	58.1	58.1	72.8	78.0	77.1												
ben_Beng	34.9	34.9	34.9	54.7	54.7	48.1	48.1	45.3	45.3	62.8	60.6	58.1	54.7	58.1	58.1	72.8	78.0	77.1												
ben_Beng	34.9	34.9	34.9	54.7	54.7	48.1	48.1	45.3	45.3	62.8	60.6	58.1	54.7	58.1	58.1	72.8	78.0	77.1												
ben_Beng	34.9	34.9	34.9	54.7	54.7	48.1	48.1	45.3	45.3	62.8	60.6	58.1	54.7	58.1	58.1	72.8	78.0	77.1												
ben_Beng	34.9	34.9	34.9	54.7	54.7	48.1	48.1	45.3	45.3	62.8	60.6	58.1	54.7	58.1	58.1	72.8	78.0	77.1												
ben_Beng	34.9	34.9	34.9	54.7	54.7	48.1	48.1	45.3	45.3	62.8	60.6	58.1	54.7	58.1	58.1	72.8	78.0	77.1												
ben_Beng	34.9	34.9	34.9	54.7	54.7	48.1	48.1	45.3	45.3	62.8	60.6	58.1	54.7	58.1	58.1	72.8	78.0	77.1												
ben_Beng	34.9	34.9	34.9	54.7	54.7	48.1	48.1	45.3	45.3	62.8	60.6	58.1	54.7	58.1	58.1	72.8	78.0	77.1												
ben_Beng	34.9	34.9	34.9	54.7	54.7	48.1	48.1	45.3	45.3	62.8	60.6	58.1	54.7	58.1	58.1	72.8	78.0	77.1												
ben_Beng	34.9	34.9	34.9	54.7	54.7	48.1	48.1	45.3	45.3	62.8	60.6	58.1	54.7	58.1	58.1	72.8	78.0	77.1												
ben_Beng	34.9	34.9	34.9	54.7	54.7	48.1	48.1	45.3	45.3	62.8	60.6	58.1	54.7	58.1	58.1	72.8	78.0	77.1												
ben_Beng	34.9	34.9	34.9	54.7	54.7	48.1	48.1	45.3	45.3	62.8	60.6	58.1	54.7	58.1	58.1	72.8	78.0	77.1												
ben_Beng	34.9	34.9	34.9	54.7	54.7	48.1	48.1	45.3	45.3	62.8	60.6	58.1	54.7	58.1	58.1	72.8	78.0	77.1												
ben_Beng	34.9	34.9	34.9	54.7	54.7	48.1	48.1	45.3	45.3	62.8	60.6	58.1	54.7	58.1	58.1	72.8	78.0	77.1												
ben_Beng	34.9	34.9	34.9	54.7	54.7	48.1	48.1	45.3	45.3	62.8	60.6	58.1	54.7	58.1	58.1	72.8	78.0	77.1												
ben_Beng	34.9	34.9	34.9	54.7	54.7	48.1	48.1	45.3	45.3	62.8	60.6	58.1	54.7	58.1	58.1	72.8	78.0	77.1												
ben_Beng	34.9	34.9	34.9	54.7	54.7	48.1	48.1	45.3	45.3	62.8	60.6	58.1	54.7	58.1	58.1	72.8	78.0	77.1												
ben_Beng	34.9	34.9	34.9	54.7	54.7	48.1	48.1	45.3	45.3	62.8	60.6	58.1	54.7	58.1	58.1	72.8	78.0	77.1												
ben_Beng	34.9	34.9	34.9	54.7	54.7	48.1	48.1	45.3	45.3	62.8	60.6	58.1	54.7	58.1	58.1	72.8	78.0	77.1												
ben_Beng	34.9	34.9	34.9	54.7	54.7	48.1	48.1	45.3	45.3	62.8	60.6	58.1	54.7	58.1	58.1	72.8	78.0	77.1												
ben_Beng	34.9	34.9	34.9	54.7	54.7	48.1	48.1	45.3	45.3	62.8	60.6	58.1	54.7	58.1	58.1	72.8	78.0	77.1												
ben_Beng	34.9	34.9	34.9	54.7	54.7	48.1	48.1	45.3	45.3	62.8	60.6	58.1	54.7	58.1	58.1	72.8	78.0	77.1												
ben_Beng	34.9	34.9	34.9	54.7	54.7	48.1	48.1	45.3	45.3	62.8	60.6	58.1	54.7	58.1	58.1	72.8	78.0	77.1												
ben_Beng	34.9	34.9	34.9	54.7	54.7	48.1	48.1	45.3	45.3	62.8	60.6	58.1	54.7	58.1																

A.7.3 Topic-To-Sentence

Lang.	OwenVL-2 2B	QwenVL-2 7B	InternVL-2.5 4B	InternVL-2.5 8B	CentriQ-Qwen	4o-mini
avg	38.6	66.5	54.5	68.0	72.4	81.5
eng_Latin	56.7	85.7	81.4	87.0	85.4	90.0
ace_Arab	29.9	46.4	34.9	42.1	54.1	66.3
ace_Latin	39.6	68.7	52.0	71.8	72.9	80.6
acm_Arab	46.6	81.1	71.5	80.5	83.2	89.9
acc_Arab	47.5	81.6	71.3	80.6	84.0	90.0
aeb_Arab	45.7	79.6	68.4	72.2	82.1	85.9
afn_Latin	49.9	82.5	70.6	82.3	83.4	89.3
ajp_Arab	46.3	81.1	70.5	80.2	83.6	89.3
aka_Latin	32.3	53.2	44.4	59.6	63.5	79.9
als_Latin	41.2	77.6	57.9	76.1	80.0	89.4
amb_Ethi	26.0	38.1	37.5	31.2	68.5	81.6
apc_Arab	46.0	80.9	70.0	80.4	82.7	89.4
arb_Arab	48.1	82.7	73.1	81.3	84.9	89.9
arb_Latin	31.2	56.1	44.9	56.6	69.8	85.1
ars_Arab	48.0	82.4	72.7	81.2	84.3	90.0
arr_Arab	43.3	77.0	64.1	75.9	81.4	88.2
arz_Arab	46.6	81.2	71.5	80.7	83.1	89.3
asm_Beng	37.0	69.6	57.8	72.0	79.2	88.2
ast_Latin	50.1	82.1	71.6	82.2	83.5	89.4
awa_Deva	41.5	74.5	59.4	78.2	81.7	88.6
ayr_Latin	32.2	46.3	37.0	52.3	54.3	62.0
azb_Arab	38.7	69.5	46.0	68.0	77.8	83.2
azj_Latin	42.6	80.2	65.5	76.9	82.0	89.1
bak_Cyril	38.1	74.3	53.1	69.6	74.7	85.2
ban_Latin	30.9	46.6	39.4	53.3	58.8	58.5
bar_Latin	46.2	75.3	56.5	74.8	79.0	85.4
bel_Cyril	44.2	80.4	63.7	77.7	83.9	89.4
bem_Latin	31.1	52.9	42.3	59.3	59.9	72.6
ben_Beng	38.9	77.2	64.5	76.9	81.9	88.6
bho_Deva	39.1	69.8	56.4	76.8	80.3	88.0
bij_Arab	28.8	44.3	33.0	41.2	52.8	64.8
bjn_Latin	43.3	73.6	57.2	72.8	78.6	84.9
bob_Latin	25.3	28.4	25.6	35.2	40.9	49.9
bos_Latin	50.8	83.6	73.5	83.6	82.9	88.8
bug_Latin	39.9	64.8	50.4	68.8	69.3	77.2
bul_Cyril	47.1	83.3	73.4	83.5	83.9	89.3
cat_Latin	51.9	83.8	75.2	84.0	84.2	89.9
ceb_Latin	41.8	79.7	63.6	79.9	82.1	88.6
ces_Latin	50.8	82.8	75.3	85.3	84.6	90.1
cjk_Latin	31.0	50.6	40.6	55.2	57.0	60.8
cdo_Arab	30.3	61.3	43.1	49.2	58.2	82.3
chr_Latin	41.1	74.8	59.4	75.7	80.0	86.4
cym_Latin	33.6	72.8	49.1	72.8	75.9	88.5
dan_Latin	52.6	83.6	74.8	84.5	83.3	88.5
deu_Latin	52.9	85.2	77.4	85.5	84.7	89.7
dik_Latin	32.0	48.6	41.0	59.3	55.2	61.9
duy_Latin	33.4	51.3	40.4	58.6	59.0	63.0
dzo_Tibet	24.9	24.8	24.2	46.9	46.8	45.1
e11_Grek	77.7	60.3	79.2	81.6	90.1	90.1
epo_Arab	46.2	83.1	81.1	83.2	82.2	89.2
est_Latin	42.3	79.0	66.3	75.5	80.2	89.3
eus_Latin	40.2	71.2	55.6	74.5	77.5	88.4
ewe_Latin	29.3	45.1	37.2	53.5	50.5	57.2
fao_Latin	36.0	74.0	54.3	73.5	76.0	87.4
fij_Latin	31.6	51.8	38.8	57.5	56.7	77.6
fin_Latin	42.3	82.0	70.5	82.3	81.8	89.2
fon_Latin	29.6	42.6	36.5	50.9	49.2	55.2
fra_Latin	53.3	85.2	77.7	86.4	85.9	89.9
fur_Arab	47.6	76.4	61.4	77.5	79.9	86.6
fuv_Arab	33.0	52.0	42.6	54.7	56.9	68.5
gaz_Latin	30.3	45.8	38.9	48.4	53.1	81.7
gia_Latin	32.3	59.9	44.7	59.3	71.6	85.8
gle_Latin	32.1	66.5	47.6	62.2	73.8	87.7
glg_Latin	51.2	83.7	75.8	85.0	85.2	88.8
grn_Latin	40.2	64.1	51.7	69.6	54.5	81.3
guj_Gujr	34.0	72.8	56.5	74.3	79.0	88.7
hau_Arab	40.6	70.8	57.9	70.0	82.0	87.3
hau_Latin	32.6	52.1	40.0	53.8	66.5	87.3
heb_Hebr	46.0	82.7	69.1	76.9	82.7	88.8
hin_Deva	39.9	78.9	65.2	80.3	84.1	89.1
hne_Deva	41.1	70.6	55.6	76.5	80.3	88.3
hrv_Latin	51.5	84.0	73.2	84.0	83.1	89.7
hun_Latin	44.3	80.9	66.8	81.3	81.2	89.9
hye_Arm	28.0	69.7	40.5	38.4	64.2	88.4
ibo_Latin	29.6	49.4	41.1	57.0	71.1	85.8
ill_Latin	39.7	56.6	47.0	68.1	76.3	87.8
ind_Latin	53.4	84.2	75.0	84.7	84.4	89.3
isl_Latin	34.1	74.1	55.2	70.6	77.2	88.8
ita_Latin	53.6	84.5	75.5	86.6	86.0	89.3
jav_Latin	46.0	77.2	59.9	74.4	80.9	87.5
jpn_Jpan	52.0	81.9	76.0	86.8	84.7	90.0
kab_Latin	27.4	35.7	29.7	40.3	41.3	46.3
kac_Latin	30.8	47.0	37.7	56.3	52.6	57.3
kam_Latin	31.6	51.6	40.4	56.2	57.9	69.7
kal_Arab	30.9	69.3	52.7	68.1	77.5	87.9
kas_Arab	36.1	62.5	46.3	67.9	71.9	84.9
kas_Deva	34.6	55.2	43.5	63.3	67.4	80.8
kat_Georg	32.8	69.0	51.9	44.9	76.8	88.9
kaz_Cyril	39.2	76.3	58.9	71.7	81.1	89.0
kbp_Latin	30.4	45.3	38.1	54.2	51.4	57.1
kea_Latin	45.9	76.3	63.8	78.2	81.1	86.2
khk_Cyril	33.3	69.1	44.6	57.6	61.4	88.5
khm_Khmer	29.9	63.6	52.7	57.9	73.0	87.4
klg_Latin	32.2	52.5	45.6	60.0	64.9	70.3
kin_Latin	31.2	49.9	39.8	54.4	54.8	86.8
kir_Cyril	37.9	72.1	53.6	67.7	77.9	88.0
kmb_Latin	29.3	46.6	38.6	53.3	53.0	59.2
kmr_Latin	34.7	60.1	50.5	65.1	65.9	83.4
knc_Arab	28.0	34.2	29.0	37.3	36.9	46.1
knc_Latin	33.1	52.0	42.8	59.8	60.5	64.2
kon_Latin	33.6	59.5	46.0	63.5	65.4	73.9
kor_Hang	48.6	82.6	76.0	86.0	83.3	88.9

Lang.	QwenVL-2 2B	QwenVL-2 7B	InternVL-2.5 4B	InternVL-2.5 8B	CentriQ-Qwen	4o-mini
lao_Lao	28.2	53.1	47.6	57.8	74.6	80.1
lij_Latin	42.2	75.8	60.3	76.7	79.2	86.4
lim_Latin	42.5	77.8	61.9	75.7	80.0	85.1
lin_Latin	34.6	59.4	46.5	62.5	75.6	81.2
lit_Latin	44.6	81.3	63.7	77.3	80.3	89.6
lmo_Latin	44.1	77.2	61.3	72.6	75.4	83.4
ltg_Latin	38.2	70.0	54.5	72.6	75.4	83.4
luu_Latin	44.8	75.6	59.7	75.7	82.8	89.0
lua_Latin	34.3	56.0	45.2	61.9	62.1	69.1
lug_Latin	30.5	48.6	39.9	53.4	53.6	76.5
luo_Latin	31.7	50.0	39.9	55.3	56.7	63.1
lus_Latin	38.0	62.6	49.2	69.2	68.3	74.7
lvu_Latin	42.6	81.4	66.1	78.0	82.2	89.5
mag_Deva	40.6	71.7	56.5	77.0	80.7	87.9
mai_Deva	42.5	70.5	53.8	77.7	80.7	89.3
mai_Myan	27.8	47.4	41.0	67.2	78.7	85.5
nld_Latin	54.1	84.0	76.9	85.3	85.1	89.0
nno_Latin	50.5	83.3	72.8	83.7	83.5	88.4
nob_Latin	52.3	83.8	75.2	84.0	84.2	88.7
npi_Deva	42.4	73.8	59.0	73.8	79.1	88.4
npo_Nkoo	25.6	24.8	25.3	27.0	26.6	26.6
nso_Latin	31.9	51.8	41.0	56.9	62.6	80.3
nus_Latin	28.4	38.7	34.8	45.2	53.1	58.9
nya_Latin	31.8	43.8	41.3	62.4	65.2	85.2
oed_Latin	40.9	57.0	50.4	70.8	83.3	89.0
ori_Orya	27.5	62.4	54.2	71.5	75.9	86.4
pag_Latin	42.8	75.3	59.5	77.5	78.6	84.3
pan_Guru	34.0	68.4	57.4	67.4	68.8	74.4
pap_Latin	43.3	75.1	62.7	76.8	78.8	85.4
pbt_Arab	37.7	70.1	48.6	65.8	72.8	86.6
pes_Arab	44.5	82.0	67.5	82.5	82.7	90.0
pit_Latin	33.6	53.8	42.0	57.6	59.9	87.0
pol_Latin	51.1	73.4	57.0	73.5	84.7	89.3
por_Latin	53.9	85.5	75.5	86.4	85.4	89.7
sot_Latin	40.5	76.4	59.6	76.2	80.2	87.7
spa_Latin	52.7	84.8	76.6	85.9	85.3	89.2
srd_Latin	42.3	74.3	56.9	76.9	80.9	84.7
srp_Cyril	46.2	82.9	71.5	80.9	83.8	89.3
ssw_Latin	30.7	48.6	39.7	53.0	69.1	82.6
sun_Latin	45.2	77.0	63.7	75.8	79.9	88.1
swi_Latin	50.8	83.9	76.1	84.8	84.0	89.6
swi_Myan	35.5	57.2	48.6	68.8	70.8	80.7
szl_Latin	42.0	76.1	57.1	76.6	78.7	87.4
tam_Taml	28.3	65.9	50.4	67.4	77.2	87.5
taq_Tfng	31.9	52.7	41.8	57.7	58.2	66.0
tat_Cyril	38.1	73.5	53.9	69.4	75.0	88.8
tel_Telu	31.5	68.7	54.6	72.0	78.8	87.7
tgk_Cyril	33.2	68.3	50.6	61.8	74.4	88.0
tgl_Latin	43.2	80.4	69.1	83.8	82.1	89.4
tha_Latin	48.0	74.2	59.3	79.4	84.5	89.3
tkz_Latin	29.2	57.3	40.3	70.3	82.0	87.3
tgl_Ethi	25.2	57.4	37.5	75.7	82.6	86.4
tpi_Latin	41.7	75.4	57.5	77.7	82.6	87.7
ts						

A.7.4 Sentences-To-Images

Lang.	QwenVL-2B	QwenVL-7B	InternVL-2.5-4B	InternVL-2.5-8B	Centrio-Qwen	4o-mini	Lang.	InternVL-2.5-4B	InternVL-2.5-4B	QwenVL-2.5-2B	QwenVL-2.5-2B	Centrio-Qwen	4o-mini						
k	1	3	5	3	5	1	1	3	5	3	5	1	3	5					
Avg	35.1	48.8	52.8	45.8	48.4	32.2	33.4	36.6	45.8	51.7	59.2	1.5	69.2	80.2	83.0				
eng_Latn	41.9	43.3	43.4	70.7	70.4	32.7	44.5	43.0	66.2	59.0	57.3	35.3	36.1	35.6	77.5	86.4	89.1		
ace_Arab	29.2	30.4	30.6	36.7	38.8	41.8	26.3	28.2	30.9	29.7	31.8	30.7	26.8	26.5	26.4	53.7	70.0	73.6	
ace_Latn	36.5	38.0	38.0	52.2	55.8	56.9	31.2	36.9	38.2	45.5	54.1	54.6	28.5	29.4	29.4	66.5	80.9	84.9	
acm_Arab	41.4	46.4	43.9	44.6	66.1	62.0	37.6	38.7	39.9	54.7	60.4	59.9	30.1	32.6	33.3	77.2	86.4	88.6	
aceb_Arab	43.1	44.1	44.1	64.7	62.0	65.3	65.4	36.8	38.5	39.0	53.2	58.6	57.6	29.8	31.7	32.8	76.2	86.4	88.6
afr_Arab	40.3	41.6	42.7	64.4	64.2	63.2	36.2	40.3	40.2	54.3	60.8	60.6	31.0	32.4	32.0	76.1	86.3	88.6	
ajb_Arab	39.2	40.4	40.4	59.4	59.2	59.2	32.3	36.9	37.0	54.9	59.0	59.0	30.9	32.3	32.3	76.0	86.0	88.0	
akz_Latn	30.2	31.7	32.4	43.1	45.2	46.6	28.2	32.3	33.6	40.8	48.9	48.7	28.4	27.9	28.3	66.4	79.5	82.9	
als_Deva	36.8	40.5	42.4	35.3	61.1	61.1	37.1	38.4	41.8	54.3	60.9	59.9	30.4	32.3	32.1	76.8	87.4	89.2	
ahl_Ethi	24.8	25.0	25.5	30.2	31.1	32.0	26.1	27.8	28.0	26.8	27.6	27.7	70.1	70.1	80.6	82.9			
apis_Arab	40.4	40.0	40.0	64.2	64.2	64.2	39.0	40.0	40.0	55.2	59.8	59.9	30.4	32.8	33.5	77.7	86.8	88.7	
arb_Arab	44.3	44.8	44.3	66.5	67.6	67.5	65.7	39.2	39.2	40.0	55.2	59.8	59.9	30.4	32.8	33.5	77.7	86.8	88.7
arb_Latn	29.6	31.5	31.0	43.3	48.0	51.1	28.5	31.5	33.2	37.5	44.4	44.4	27.5	28.2	28.2	73.4	84.8	87.4	
ars_Arab	43.9	44.5	44.5	65.8	65.8	65.8	37.8	38.0	38.0	56.6	59.2	59.2	30.2	32.3	32.3	77.7	86.8	88.4	
ary_Arab	37.0	37.4	38.4	60.9	61.7	61.7	35.9	37.4	38.0	49.9	56.6	56.6	30.2	32.3	32.3	77.1	86.0	88.4	
arz_Arab	43.3	44.9	44.0	64.4	66.1	65.0	37.5	39.2	40.2	53.9	58.9	58.9	33.3	34.3	34.3	77.6	86.6	88.6	
asm_Beng	35.2	36.9	36.9	63.9	59.9	59.9	36.6	41.9	43.9	45.9	51.4	51.4	32.2	33.2	33.2	73.9	85.8	88.1	
ast_Arab	30.4	30.3	30.3	59.4	59.4	59.4	30.3	30.4	30.4	50.9	59.8	59.8	33.0	33.9	33.9	74.8	85.8	88.3	
ast_Deva	40.5	42.4	43.3	61.1	61.1	61.1	35.1	38.4	41.8	54.3	62.3	61.6	33.6	33.8	32.8	75.4	86.3	88.2	
ayr_Arab	28.4	30.4	30.3	58.4	58.4	58.4	29.3	27.5	27.5	45.5	52.7	51.1	26.7	27.1	26.4	69.7	86.0	87.8	
azb_Arab	37.6	40.0	38.7	53.2	58.1	60.2	31.1	32.6	35.7	49.3	49.3	49.9	29.4	30.9	30.9	71.0	83.8	86.6	
azj_Arab	39.4	41.3	41.3	64.5	64.5	64.5	30.6	31.1	31.1	49.3	52.7	52.7	26.4	27.1	27.1	73.0	86.0	88.0	
bak_Cyril	34.6	39.2	38.6	57.4	60.1	61.2	30.6	36.4	38.2	42.9	48.8	48.2	29.2	29.5	29.5	72.9	86.0	88.7	
ban_Latn	27.7	30.3	30.7	32.4	37.3	37.3	25.8	29.4	31.2	34.1	34.8	30.3	40.3	26.6	27.0	27.0	64.2	63.3	
baq_Arab	41.6	41.6	41.6	62.1	62.1	62.1	30.4	30.4	30.4	50.4	51.0	51.0	28.1	28.1	28.1	72.6	84.7	87.4	
baq_Latn	20.6	21.0	21.2	42.2	63.2	63.3	34.4	38.1	40.1	50.2	54.4	54.4	31.0	32.2	32.2	73.3	83.8	86.8	
ben_Latn	31.5	33.3	33.3	45.5	46.5	46.5	28.4	31.8	34.2	41.6	48.2	47.9	27.6	27.7	27.7	75.7	86.3	88.2	
benz_Arab	39.8	40.4	41.1	64.0	67.0	67.0	31.2	43.4	43.9	50.9	58.1	58.1	33.9	33.9	33.9	74.4	85.7	88.4	
bjb_Arab	37.8	38.4	38.4	59.0	59.0	59.0	30.6	30.6	30.6	52.6	52.6	52.6	26.2	26.2	26.2	71.0	81.6	84.0	
bjn_Arab	28.4	29.5	28.3	30.4	35.6	57.5	25.3	26.7	29.7	39.2	50.9	50.9	31.0	31.0	31.0	71.7	81.7	84.0	
bob_Tibet	23.6	24.7	24.7	32.9	24.0	24.0	22.8	23.4	24.9	33.3	33.5	33.5	30.6	30.6	30.6	70.3	81.7	84.5	
bug_Arab	39.4	40.4	40.4	64.7	64.7	64.7	30.4	30.4	30.4	50.6	50.6	50.6	29.2	29.2	29.2	72.9	85.6	88.4	
bug_Deva	35.8	38.1	38.5	51.4	53.3	53.8	31.1	35.4	36.7	44.5	52.5	53.4	28.2	29.2	29.2	72.9	85.2	88.4	
bug_Latn	29.7	32.3	33.1	37.3	39.2	39.2	27.3	30.5	32.2	37.4	45.1	44.7	26.7	27.7	27.7	71.1	81.6	88.4	
dyu_Latn	29.7	32.3	33.1	37.3	39.2	39.2	27.3	30.5	32.2	37.4	45.1	44.7	26.7	27.7	27.7	71.1	81.6	88.4	
dzo_Tibet	23.6	23.9	24.4	31.4	31.4	31.4	21.8	26.6	26.6	30.4	30.4	30.4	25.7	25.7	25.7	77.7	87.7	88.4	
edz_Arab	39.2	40.9	40.9	63.1	63.4	63.7	30.3	31.1	31.1	51.9	55.3	55.3	30.9	30.9	30.9	77.7	86.5	88.5	
ces_Latn	42.7	43.2	42.7	65.5	64.8	63.0	37.3	39.7	41.1	60.4	65.4	64.6	31.1	32.8	34.2	77.9	86.6	88.6	
cjk_Arab	30.4	33.4	32.9	38.4	40.9	42.4	27.8	30.1	31.5	36.4	43.5	43.5	36.7	36.7	36.7	77.7	86.7	88.7	
crk_Arab	37.8	38.4	38.4	61.6	61.6	61.6	30.8	31.1	31.1	34.1	36.4	36.4	29.1	29.1	29.1	77.7	86.7	88.7	
crk_Deva	37.8	40.2	39.8	58.9	62.6	62.6	30.6	31.1	31.1	34.1	36.4	36.4	29.1	29.1	29.1	77.7	86.7	88.7	
cym_Deva	31.9	34.6	34.5	55.4	58.1	59.3	29.2	34.0	35.2	46.3	53.2	53.2	30.3	30.5	30.5	76.6	86.7	88.8	
dan_Latn	41.9	43.1	43.2	66.2	66.5	67.0	38.3	40.5	40.5	58.1	63.1	63.1	30.2	32.9	32.9	77.7	86.8	88.4	
deu_Arab	39.2	40.4	40.4	61.3	61.3	61.3	30.3	31.3	31.3	46.1	52.3	52.3	29.1	29.1	29.1	77.7	86.8	88.4	
dik_Latn	29.2	31.9	31.9	39.3	43.3	45.6	28.2	31.9	32.9	38.4	40.7	40.7	26.3	26.3	26.3	77.7	86.8	88.4	
dik_Arab	39.2	37.0	35.2	35.1	35.2	30.1	30.1	35.2	40.7	40.3	40.3	26.3	26.3	26.3	77.7	86.8	88.4		
don_Arab	41.2	42.4	42.4	62.6	62.6	62.6	30.4	32.4	32.4	40.3	46.5	46.5	27.7	27.7	27.7	77.7	86.8	88.4	
dot_Arab	39.1	41.3	41.3	62.9	62.9	62.9	33.9	37.4	37.4	50.8	58.8	58.8	31.0	32.1	32.1	77.7	86.8	88.4	
gaz_Latn	27.0	30.0	29.7	33.8	36.0	37.3	25.6	28.9	30.7	31.9	34.1	34.1	26.1	26.6	26.6	77.7	86.8	88.4	
glia_Latn	29.6	31.5	30.3	44.6	46.2	47.2	28.3	30.7	30.7	37.0	40.4	40.4	27.6	27.6	27.6	77.7	86.8	88.4	
gle_Latn	30.9	31.4	31.4	47.1	47.1	47.1	20.3	20.3	20.3	44.6	44.6	44.6	27.4	27.4	27.4	77.7	86.8	88.4	
gle_Deva	32.3	32.3	32.3	42.9	42.9	42.9	20.8	20.8	20.8	44.6	44.6	44.6	27.4	27.4	27.4	77.7	86.8	88.4	
grn_Deva	36.4	38.4	38.5	55.4	58.7	58.7	31.6	36.6	38.4	44.6	55.3	55.3	26.9	26.9	26.9	77.7	86.8	88.4	
guj_Guru	34.5	35.8	34.3	63.0	63.0	63.0	30.4	31.1	31.1	41.4	41.4	41.4	27.9	27.9	27.9	77.7	86.8	88.4	
hau_Latn	29.2	31.5	30.8	38.0	40.6	42.6	28.4	31.9	32.9	39.4	41.6	41.6	27.1	27.1	27.1	77.7	86.8	88.4	
heb_Hebr	41.6	43.6	43.4	65.9	65.9	63.3	37.1	39.5	39.5	50.4	56.8	56.8	30.3	32.4	32.4	77.7	86.8	88.4	
hei_Arab	39.0	40.8	40.8	60.0	60.0	60.0	30.4	30.4	30.4	46.0	50.0	50.0	28.8	28.8	28.8	77.7	86.8	88.4	
hrv_Latn	42.8	43.2	44.2	66.9	65.7	63.5	37.5	39.9	39.6	50.8	52.5	52.5	30.3	32.2	32.2	77.7	86.8	88.4	
hrv_Arab	38.1	40.4	41.1	64.7	64.7	64.7	30.4	32.1	32.1										

