

Breaking Language Barriers in Visual Language Models via Multilingual Textual Regularization

Anonymous ACL submission

Abstract

Rapid advancements in Visual Language Models (VLMs) have transformed multimodal understanding but are often constrained by generating English responses regardless of the input language. This phenomenon has been termed as Image-induced Fidelity Loss (IFL) and stems from limited multimodal multilingual training data. To address this, we propose a continuous multilingual integration strategy that injects text-only multilingual data during visual instruction tuning, preserving the language model’s original multilingual capabilities. Extensive evaluations demonstrate that our approach significantly improves linguistic fidelity across languages without degradation in visual performance. We also explore model merging, which improves language fidelity but comes at the cost of visual performance. In contrast, our core method achieves robust multilingual alignment without trade-offs, offering a scalable and effective path to mitigating IFL for global VLM adoption.

1 Introduction

Large Language Models (LLMs) have significantly advanced multimodal understanding, leading to the rise of VLMs, which integrate vision encoders into LLM backbones. A widely adopted paradigm is the LLaVA-style architecture (Liu et al., 2023b, 2024a), where a decoder-only LLM is coupled with a vision encoder and an adapter module to align visual representations with textual embeddings.

Despite their success, VLMs exhibit a strong bias toward English due to the predominance of monolingual vision-language training data. Consequently, they often generate English responses regardless of the input language, a phenomenon termed Image-induced Fidelity Loss (or IFL) (Hinck et al., 2024). This issue stems from limitations in the underlying LLM rather than the visual representations.

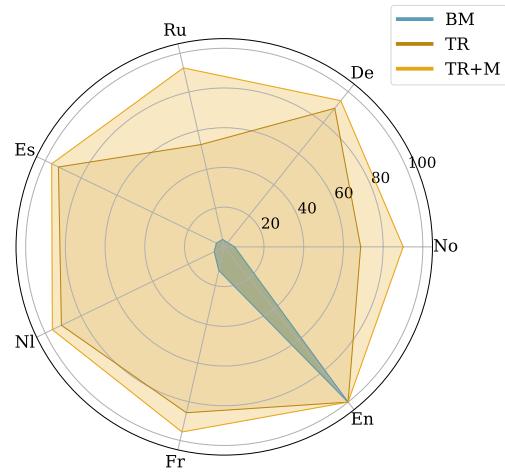


Figure 1: Language Fidelity (LF) accuracy on Crossmodal-3600. (BM: Base Model, TR: model trained with multilingual Textual Regularization, TR+M: TR and merging the final model with the original LLM Backbone)

Ensuring multilingual capability in VLMs is essential for their adoption across diverse linguistic communities, as reliance on English-centric outputs risk erasing cultural and linguistic nuances. Prior work (Qiu et al., 2022; Li et al., 2023b) has explored dataset translation, but this approach incurs high computational costs and introduces translation errors, especially in images with language-dependent elements.

In this paper, we propose an alternative solution by integrating multilingual text-only data during the visual instruction tuning process. Additionally, we explore model merging, combining the visually fine-tuned model with the original multilingual backbone LLM to further preserve linguistic fidelity. As shown in Figure 1, our method effectively prevents the model from defaulting to English in non-English queries. To the best of our knowledge, no previous work has achieved full multilingual competence in VLMs through such a simple and scalable approach.

063 Our contributions are as follows:

- 064 • We systematically demonstrate that integrating
065 multilingual text-only data during training
066 significantly reduces IFL bias in LLaVA-style
067 VLMs while maintaining core capabilities.
- 068 • We conduct an extensive analysis on the optimal
069 proportion of text-only data required for
070 effective multilingual adaptation.
- 071 • We explore a model merging strategy, combining
072 the visually fine-tuned model with the original
073 multilingual backbone LLM, and assess
074 its impact on preserving linguistic fidelity.

075 Our findings suggest that we can develop high-
076 quality multilingual VLMs that maintain strong
077 performance across multiple languages in a simple
078 and scalable way. By avoiding the need to
079 translate or construct multimodal datasets for each
080 language, our approach lowers the entry barrier for
081 multilingual VLM development. This makes it especially
082 attractive for low-resource settings, where
083 monolingual text is often available but collecting
084 vision-language data is costly or impractical.

085 2 Related Work

086 2.1 Multimodal Large Language Models

087 VLMs typically integrate an image encoder, usually CLIP (Radford et al., 2021; Dosovitskiy et al.,
088 2021), with an LLM backbone. Various strategies exist for combining these components. The predominant approach follows a decoder-only architecture, as seen in the LLaVA series, where an adapter module projects visual representations into the textual embedding space. Other methods include cross-attention mechanisms (Grattafiori et al., 2024), and some models, like NVLM (Dai et al., 2024), adopt a hybrid strategy combining both approaches.

099 LLaVA-style models tend to default to English
100 due to the scarcity of multimodal training data in
101 other languages (Hinck et al., 2024). This issue
102 arises because the LLM’s parameters are updated
103 for a distinct task, which can disrupt its original
104 language capabilities. Llama 3 (Grattafiori et al.,
105 2024) takes a different approach by freezing the
106 LLM during training, which helps preserve its pre-
107 trained abilities while incorporating visual informa-
108 tion. However, freezing the LLM also limits the

109 model’s capacity to learn new visual tasks, creating
110 a trade-off between language preservation and
111 multimodal learning.

112 2.2 Multilingual Multimodal Learning

113 A widely adopted approach to improving multilingual-
114 ity in VLMs is translating existing multimodal
115 datasets. Several works (Song et al., 2024; Hu
116 et al., 2024) have analyzed this strategy and pro-
117 posed methods to enhance its effectiveness. Several
118 models, such as PALI (Chen et al., 2023), PALI-
119 X (Chen et al., 2024d), mBLIP (Geigle et al., 2024),
120 PAPO (Maaz et al., 2024) and Pangea (Yue et al.,
121 2025), have pursued this approach. However, this
122 strategy presents challenges, including computa-
123 tional overhead, translation inconsistencies, and
124 the loss of cultural context in visual-text pairs.

125 Moreover, recent research (Aggarwal et al.,
126 2024) suggests that continual fine-tuning can harm
127 an LLM’s performance. When a model undergoes
128 two consecutive fine-tuning phases with differ-
129 ing task distributions, its ability to perform earlier tasks
130 deteriorates. This raises concerns that direct fine-
131 tuning solely on translated multimodal data may
132 degrade the LLM’s original capabilities.

133 2.3 Catastrophic Forgetting Prevention

134 In the context of LLMs, the problem of maintaining
135 performance across tasks while integrating new in-
136 formation is known as lifelong learning. This field
137 focuses on a system’s ability to acquire, integrate,
138 and retain knowledge without catastrophically for-
139 getting previous information. Visual Instruction
140 Tuning is a case of lifelong learning, and it faces
141 the same challenges. One known mitigation strat-
142 egy is episodic or experience replay (Zheng et al.,
143 2025), which helps prevent catastrophic forgetting
144 by reintroducing previously learned information.

145 Several studies (Liu et al., 2022; Ibrahim et al.,
146 2024) have explored ways to incorporate pretrain-
147 ing data during fine-tuning. Bethune et al. (2025)
148 further analyze the impact of this approach and
149 suggest that even a small amount of pretraining
150 data can help retain previously learned knowledge,
151 reducing the risk of performance degradation.

152 In the case of VLMs, NVLM (Dai et al., 2024)
153 and InternVL 2.5 (Chen et al., 2024e) demon-
154 strate that incorporating high-quality text-only data dur-
155 ing Visual Instruction Tuning, not only improves
156 the overall text-generation capabilities, but also
157 multimodal performance. Our approach builds
158 upon these findings by integrating multilingual text-

159 only data throughout VLM training to mitigate IFL,
160 without requiring extensive multimodal multilingual
161 data collection.

162 2.4 Model Merging

163 Model merging is a technique that involves combining
164 two or more pre-trained models to create
165 a new model that leverages the strengths of each.
166 By merging a fine-tuned model with its original
167 backbone, this process preserves the model’s prior
168 capabilities while incorporating additional refine-
169 ments from further training. This strategy has been
170 applied in various contexts, such as language trans-
171 fer, where [Alexandrov et al. \(2024\)](#) demonstrate
172 that model merging facilitates fine-tuning for new
173 linguistic capabilities without compromising the
174 performance of the original LLM.

175 Building on this insight, we explore model merg-
176 ing as a means of preserving the multilingual com-
177 petencies of a VLM during the visual fine-tuning
178 process. We adopt the same model merging strat-
179 egy as Aya Vision ([Dash et al., 2025](#)), which has
180 shown strong empirical results, and combine it with
181 our multilingual textual regularization strategy.

182 3 Experimental Setup

183 3.1 Data

184 Our training framework combines multimodal
185 visual-language data from LLaVA-OneVision ([Li](#)
186 et al., 2025) with multilingual text-only instruc-
187 tion data from the Salamandra family of mod-
188 els ([Gonzalez-Agirre et al., 2025](#)). This hybrid ap-
189 proach ensures robust visual understanding while
190 addressing IFL through explicit multilingual text
191 supervision. All datasets are documented in Ap-
192 pendix A.

193 **Visual Data** We employ LLaVA-OneVision’s
194 English-only visual pipeline, consisting of a to-
195 tal number of 9,286,732, which is divided into two
196 main groups:

- 197 • **General and Detailed Image Captions:** This
198 dataset comprises both basic and highly de-
199 tailed image captions. The basic captions
200 align the visual embedding space with the
201 LLM’s embedding space, while the detailed
202 captions refine the mapping between the two
203 providing a high-quality understanding of the
204 images. This group comprises 4.4M unique
205 instances.

206 • Task-Specific, Multi-Image, and Video

207 **Data:** This dataset is used to instruct the
208 aligned model on specific tasks, including
209 Optical Character Recognition (OCR), info-
210 graphic understanding, and math & reasoning.
211 Additionally, multi-image and video data are
212 incorporated to enhance the model’s ability
213 to interpret diverse visual inputs. This group
214 comprises 4.9M unique instances.

215 Multilingual Text-Only Data

216 To further enhance the model’s multilingual profi-
217 ciency, we incorporate 315,496 text-only samples drawn
218 from 11 diverse datasets covering domains such as general
219 language tasks, multilingual instructions, conversa-
220 tional QA, and code annotations. These sources
221 include human-annotated datasets ([Rajani et al.,](#)
222 [2023](#)), multilingual instruction collections ([Singh](#)
223 et al., 2024; [Costa-jussà et al., 2024](#)) and conversa-
224 tional data ([Conover et al., 2023](#); [Köpf et al.,](#)
225 [2023](#)).

226 Notably, the text-only samples cover 21 of the
227 35 languages used in training Salamandra, ensuring
228 extensive linguistic representation. A significant
229 portion of this dataset is machine-translation data.

230 Although most of the text-only data is in English,
231 matching the language of the visual data, its inclu-
232 sion remains important. This alignment reinforces
233 the model’s linguistic foundation and facilitates the
234 integration of multilingual supervision, ultimately
235 ensuring balanced performance across modalities.
236 The final distribution of languages in the text-only
237 data, complementing LLaVA-OneVision’s training
238 data, is shown in Figure 2.

239 **Evaluation Data** For evaluation, we use
240 both monolingual and multilingual multimodal
241 datasets. To assess visual performance, we
242 include AI2D ([Kembhavi et al., 2016](#)), which
243 tests understanding of diagram-based questions;
244 RealWorldQA¹, a real-world image dataset
245 with open-ended and multiple-choice questions;
246 MMMU ([Yue et al., 2024](#)), a diverse multimodal
247 reasoning benchmark; and MMStar ([Chen et al.,](#)
248 [2024b](#)), which aggregates vision-language tasks for
249 broad multimodal evaluation.

250 The first two benchmarks primarily assess the
251 exact match accuracy, quantifying the proportion of
252 responses that exactly match the predefined ground
253 truth. These targets are typically short-form text

¹<https://huggingface.co/datasets/visheratin/realworldqa>

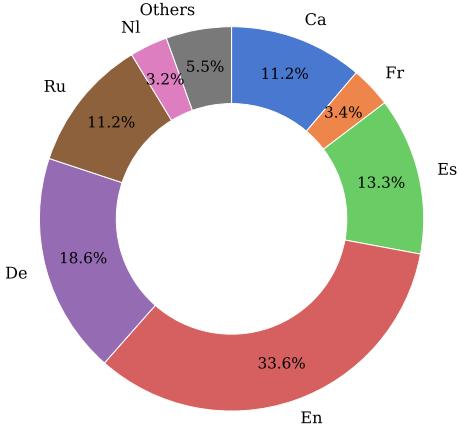


Figure 2: Distribution of the multilingual text-only data used for Textual Regularization. Languages with a volume smaller than 3% are grouped under *Others*, which collectively account for 5.5% of the data. The most frequent languages in this group are Portuguese (2.1%), Italian (0.7%), Polish (0.47%), Swedish (0.42%), Irish (0.39%), Lithuanian (0.29%), Galician (0.22%), Greek (0.20%), and Ukrainian (0.17%).

or multiple-choice answers. On the other hand, MMMU and MMStar are classification tasks that are measured using accuracy.

For multilingual multimodal performance, we have selected Crossmodal-3600 (Thapliyal et al., 2022), a geographically diverse multilingual multimodal dataset for image captioning. It is particularly well-suited for our experiments as it covers the highest number of overlapping languages with Salamandra while allowing for image captioning with multiple reference targets per instance. The dataset comprises approximately 3,600 samples in 36 languages, from which we evaluate on German (De), Russian (Ru), Spanish (Es), Dutch (Du), French (Fr), and English (En). Appendix C discusses the prompts used for caption generation during evaluation and the rationale behind their selection.

3.2 Model Framework

Our implementation also follows the LLaVA-OneVision framework², adapted for enhanced multilingual capabilities. The architecture consists of three main components: SigLIP-SO400M³ (Zhai et al., 2023) as the visual encoder, a 2-layer MLP projection module (with GELU activation functions (Hendrycks and Gimpel, 2016)), and

Salamandra-7b-instruct⁴ (Gonzalez-Agirre et al., 2025) as the backbone LLM. Salamandra has been chosen for its high multilinguality, as it has been trained with 35 European languages.

We also adopt LLaVA-OneVision’s curriculum learning training strategy, which progresses through four distinct stages:

- **Stage 1 (Language-Image Alignment):** In this initial phase, only the MLP projector is trained, while both the visual encoder and LLM remain frozen. General image captions are employed to establish basic cross-modal connections.
- **Stage 1.5 (Full Model Training):** At this stage, all model components are unfrozen to enable end-to-end training. A high-quality set of detailed image captions is used in conjunction with an increased image resolution to enhance visual detail processing.
- **Stage 2 (Single-Image Instruction Tuning):** Once the model has achieved a deep understanding of images, it is fine-tuned for a diverse set of visual tasks. The image resolution is further increased to support fine-grained visual analysis.
- **Stage 2.5 (Multi-Image and Video Training):** In the final stage, multi-image and video data are incorporated to enable reasoning across multiple visual inputs. Additionally, single-image data from the previous stage is also utilized.

The key innovation in our approach lies in the strategic injection of multilingual text-only data throughout these training stages detailed in §3.4. After the Visual Instruction Tuning, the model is merged with the baseline LLM weights using linear interpolation.

3.3 Metrics

To evaluate language fidelity and consistency, we employ a common metric established in prior multilingual multimodal evaluation work (Hinck et al., 2024; Schneider and Sitaram, 2024):

Language Fidelity We use GlotLID (Kargaran et al., 2023) to obtain the accuracy of whether the language of the generated captions over

²<https://github.com/LLaVA-VL/LLaVA-NeXT>

³<https://huggingface.co/google/siglip-so400m-patch14-384>

⁴<https://huggingface.co/BSC-LT/salamandra-7b-instruct>

Crossmodal-3600 images is the same as the user prompt⁵. We named this metric LF, and we observed that, in many cases, it considered as correct samples that had single words in English, or with minor code-switching errors. To address this issue, we extend this metric (LF+) by using Llama-3.1-8B-Instruct⁶ (Grattafiori et al., 2024) as an LLM-as-a-judge, evaluating if the samples already classified by GlotLID are entirely in the same language or not. Nevertheless, due to a majority voting strategy in its implementation (see Appendix B), the LLM-as-a-judge does not work perfectly, as it sometimes misclassifies correct samples as non-consistent in language. For this reason, this metric can be interpreted as a statistical lower bound of language fidelity.

Visual Performance To evaluate visual performance, we use the English-only multimodal benchmarks detailed in Section 3.1: AI2D, RealWorldQA, MMMU, and MMStar. AI2D and RealWorldQA are evaluated using exact match accuracy, measuring the proportion of responses identical to ground-truth answers (typically short text or multiple-choice). MMMU and MMStar, however, are treated as classification tasks and evaluated via accuracy.

To evaluate multilingual multimodal performance, we also use the same approach used in Hinck et al. (2024); Schneider and Sitaram (2024), and evaluate the captioning quality across different languages with chrF++ (Popović, 2016, 2017) over Crossmodal-3600 samples.

Further discussion on metric selection can be found in Appendix D.

3.4 Experiments

We focus on testing different text-only integration strategies, analyzing the influence of data quantity, examining generalization capabilities on languages not contained during textual regularization, and assessing the effect of model merging on overall performance. To be able to quantify the results obtained with these experiments, we also trained a baseline model (BM) by only conducting the Visual Instruction Tuning, without textual regularization.

⁵Crossmodal-3600 does not include a predefined reference generation prompt. For completeness, we present and explain the employed prompt in Appendix C.

⁶<https://huggingface.co/meta-llama/Llama-3-1-8B-Instruct>

Multilingual Data Integration Strategies We explore three distinct strategies for incorporating the text-only multilingual data (315,496 instances) during the visual instruction tuning process:

- **Textual Regularization across Three Stages (TR-3S)**: Multilingual text data was distributed proportionally across the final three training stages (1.5, 2, and 2.5).
- **Textual Regularization across Two Stages (TR-2S)**: Multilingual text data was integrated proportionally only in the last two stages (2 and 2.5).
- **Textual Regularization at a Single Stage (TR-1S)**: Multilingual text data was added exclusively during the final stage (2.5).

Multilingual Generalization Capabilities To investigate whether regularization with multilingual text data extends to languages not explicitly seen during training, we train a variant of the TR-3S model where German was excluded from the multilingual text dataset. This experiment allows us to evaluate the model’s generalization ability to new languages.

Influence of Data Balance We vary the proportion of multilingual text data used for regularization alongside visual–text pairs. Starting from the text-only, multilingual instruction-tuning subset of 315k text-only instances, we inject this data following the TR-3S approach. Across these three stages, the total number of visual–text pairs is $\approx 9M$, of which the 315k text-only examples account for approximately 3.7% of the entire set. We denote $2x$, $0.5x$, and $0.25x$ as using twice, half, and one quarter of the amount of text-only data, respectively, where x corresponds to the original dataset used in the TR-3S configuration.

Model Merging To explore the potential of further enhancing the multilingual capabilities of our best-performing model (TR-3S), we apply model merging. To do so, we perform a linear interpolation between the weights of the visually instructed model with those of the backbone LLM, maintaining the encoder and MLP layers. This allows us to evaluate whether model merging could combine the model’s visual understanding capabilities with the language fidelity of the original model.

As explained in §3.3, we evaluate multimodal performance on a suite of English benchmarks and

417 extend it to multiple languages evaluating chrF++
418 on Crossmodal-3600. Moreover, we use LF to
419 assess IFL and further analyze its bounds with LF+.

420 3.5 Implementation Details

421 Our experiments were conducted on custom
422 NVIDIA H100 GPUs, each with 64GB of memory.
423 We trained each model for 6 days in a distributed
424 setup with 8 nodes, each containing 4 GPUs, to-
425 taling 32 GPUs per experiment. As we trained 8
426 models (excluding the merged model, which did
427 not require separate training), the total compute
428 usage amounted to 36,864 GPU hours.

429 For evaluation, we assessed 9 models across 6
430 languages, with each requiring one node for 24
431 hours, resulting in 5,184 GPU hours.

432 The training hyperparameters were largely based
433 on those used in LLaVA-OneVision and Salamand-
434 dra’s Instruction Tuning, ensuring consistency with
435 prior work. Further details on the training process,
436 including specific hyperparameters and configura-
437 tions, can be found in Appendix E.

438 4 Results

439 This section presents the outcomes of our exper-
440 imental investigation into the effectiveness of in-
441 corporating multilingual text-only data during the
442 visual instruction tuning process for reducing IFL
443 in VLMs.

444 4.1 Quantifying the Baseline English Bias

445 To better understand the starting point of our in-
446 vestigation, we first evaluated the baseline model
447 (BM), trained exclusively on English visual instruc-
448 tion data. As anticipated, this model exhibits a
449 pronounced English-centric behavior, responding
450 predominantly in English even when prompted in
451 other languages. This confirms the strong presence
452 of IFL and underscores the necessity of multilin-
453 gual regularization.

454 As shown in Table 1, the model demonstrates
455 a very limited capacity to generate non-English
456 responses. Languages such as German and Spanish,
457 for example, show particularly low consistency,
458 often defaulting back to English. This behavior
459 reveals how the training process strongly anchors
460 the model to English due to the lack of multilingual
461 signals.

462 Interpreting these results in context, the baseline
463 model’s bias highlights a fundamental limitation of
464 current VLM training pipelines, where even models

465 based on multilingual backbones revert to English
466 if not explicitly trained with multilingual supervi-
467 sion.

468 4.2 Impact of Multilingual Text-Only Data 469 Integration Strategies

470 The results presented in Table 1 clearly demon-
471 strate that integrating multilingual text-only data substi-
472 tantly mitigates English bias across all evaluated
473 strategies. Notably, proportional integration across
474 the final three training stages (TR-3S) consistently
475 achieves superior LF scores for most non-English
476 languages. This suggests that continuous expo-
477 sure to multilingual text throughout training stages
478 is most effective in maintaining linguistic fidelity.
479 The strategy of introducing multilingual data ex-
480clusively in the final stage (TR-1S) yields the least
481 improvement, indicating that delaying multilingual
482 exposure is insufficient to counteract the English
483 bias ingrained during earlier training phases. The
484 stronger performance of TR-3S can be attributed
485 to its role as a continual regularizer. By more ex-
486 tensively interleaving multilingual text-only data,
487 the model consistently reinforces its multilingual
488 representations, thus more effectively preserving
489 previously acquired language capabilities and re-
490 ducing IFL.

Lang.	BM	TR-3S	TR-2S	TR-1S
De	2.7	88.7	81.3	24.5
Es	4.4	92.9	65.4	38.4
Fr	12.2	85.7	74.9	29.9
Nl	5.6	91.8	91.3	49.2
Ru	3.8	52.9	24.8	50.9
En	100.0	100.0	100.0	100.0

491 Table 1: LF accuracy for different integration strategies.
492 The best results are shown in bold.

493 4.3 Evaluating Multilingual Generalization 494 Capabilities

495 We obtained a LF scores of 5.4% for German in
496 this scenario. While these results are slightly above
497 the English-biased baseline (2.7% LF), the perfor-
498 mance remains very limited. This suggests that
499 the multilingual regularization approach, in the ab-
500 sence of explicit exposure to the target language,
501 does not meaningfully help mitigate IFL. In other
502 words, the model struggles to generalize to unseen
503 languages, and explicit inclusion during training

502 appears necessary for achieving satisfactory multi-
 503 lingual fidelity.

504 4.4 Analyzing the Influence of Data Quantity

505 The LF score for these variations are presented on
 506 Table 2.

Lang.	$0.25x$	$0.5x$	x	$2x$
De	85.0	88.9	88.7	73.3
Es	91.9	92.4	92.9	76.4
Fr	88.6	83.5	85.7	69.5
Nl	69.8	96.0	91.8	93.5
Ru	91.5	81.9	52.9	93.3

Table 2: LF accuracy under varying quantities of text-only to visual data. All this variants are trained with the TR-3S approach. The best results are shown in bold.

507 The results demonstrate a complex relationship
 508 between the text-only to visual data quantity and
 509 language fidelity, making straightforward inter-
 510 pretation challenging. Increasing the quantity by
 511 changing the configuration from $0.25x$ to $0.5x$ gen-
 512 erally improves LF scores for most languages, sug-
 513 gesting a positive impact of increased text-only data
 514 within this range. However, further increasing the
 515 relative quantity of text-only data switching to the
 516 $2x$ configuration ($\approx 7\%$ of the total data) does not
 517 consistently yield better results and, in some cases,
 518 significantly reduces performance, particularly for
 519 German, Spanish, and French.

520 Notably, none of the tested variations drastically
 521 degrade LF across all languages compared to the
 522 baseline x configuration (3.7% of text-only data
 523 from the total amount). This indicates that while
 524 the optimal text-only data ratio requires careful con-
 525 sideration, moderate variations around the original
 526 amount do not necessarily lead to a substantial loss
 527 in language fidelity.

528 4.5 Evaluating the Effect of Model Merging

	TR-3S (%)	TR-3S + M (%)
De	88.7	94.1
Es	92.9	96.4
Fr	85.7	95.5
Nl	91.8	96.1
Ru	52.9	92.4

Table 3: LF score on multiple languages before (TR-3S) and after (TR-3S + M) model merging. The best results are shown in bold.

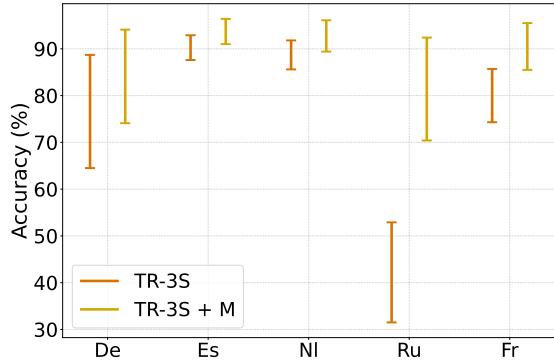


Figure 3: Interval Plot contrasting LF (upper bars) vs. LF+ (lower bars) across languages of our best-performing models.

The results presented in Table 3 demonstrate a clear positive impact of model merging on language fidelity (LF). On average, LF improved by over 12.5 points across the five languages (5.75 points removing Russian). This improvement is consistent across all evaluated languages, indicating a robust effect of the merging strategy in enhancing multilingual fidelity, rather than a language-specific anomaly. Importantly, this analysis focuses solely on language fidelity. The impact of model merging on other VLM capabilities will be discussed in §4.7. These findings support the hypothesis that model merging can be a valuable tool in mitigating IFL.

4.6 Bounding IFL

Figure 3 represents the upper (LF) and lower bounds (LF+) of IFL of our best-performing models. As it can be seen, in most languages we can observe a span of approximately 10% or less, except for De and Ru, which is around 20%. We attribute this difference to the lower performance of the LLM-as-a-judge on these languages (detailed in Appendix B). Nonetheless, it is important to note that this is not a statistical representation, so the interval width is not significant in terms of performance.

The results confirm the validity of our method, as all the lower bounds (excluding TR-3S on Russian) surpass the 65% accuracy.

4.7 Impact on General Tasks Performance

Our analysis shown in Table 4 reveals that the models trained with our proposed regularization techniques, specifically those utilizing proportional multilingual text-only data integration, generally maintain or slightly improve performance on stan-

Model	AI2D _{EM}	RealWorldQA _{EM}	MMMU _{Acc} (val)	MMStar _{Acc} (avg)
BM	73.96	56.99	34.22	47.33
TR-3S	75.39	54.25	33.56	48.87
TR-3S + M	57.19	52.03	34.11	42.25

Table 4: Performance on general VLM benchmarks (only in English). All scores are reported on a 0–100 scale. The best results are shown in bold.

dard VLM benchmarks compared to the baseline English-centric model (BM). For instance, the TR-3S model, which incorporates text-only data across three training stages, exhibits an increase in AI2D and MMStar scores. This demonstrates that our method effectively mitigates IFL without sacrificing the model’s core visual-language understanding capabilities. The strategic injection of multilingual text-only data appears to reinforce the LLM’s inherent multilingual abilities without disrupting its ability to process and understand visual information.

The evaluation using the chrF++ metric, which measures the quality of text generation by comparing character n-grams, further supports the effectiveness of our multilingual regularization techniques. As shown in Table 5, the TR-3S model demonstrates improved chrF++ scores across all non-English languages compared to the baseline (BM). For instance, German improves from 15.0 to 20.4, and Spanish from 19.1 to 23.7. This indicates that the model not only maintains language fidelity but also generates more accurate and coherent text in multilingual settings.

Model	De	Ru	Es	Nl	Fr	En
BM	15.0	9.9	19.1	16.2	18.1	27.5
TR-3S	20.4	12.5	23.7	22.0	22.8	28.2
TR-3S + M	16.1	10.5	21.6	14.7	18.5	25.5

Table 5: Performance on Crossmodal-3600 by language (chrF++). The best results are shown in bold.

However, a notable observation is the performance degradation observed in the merged model. Despite achieving substantial improvements in multilingual fidelity, the TR-3S M model shows a significant decrease in performance on benchmarks such as AI2D. This decline suggests a potential trade-off between enhanced multilingual capabilities and general task performance when employing model merging techniques. We hypothesize that the merging process, while beneficial for consolidating multilingual knowledge, may introduce con-

flicts or misalignments in the model’s learned visual representations. We further investigated alternative merging methods, including spherical linear interpolation (slerp) and both asymmetric weightings that favor the original backbone (75–25) and the visually instructed model (25–75). These variants, detailed in Appendix F, confirm the trade-offs between language fidelity and multimodal performance, without revealing a universally superior configuration.

Examples of the generation with the TR-3S model can be found in Appendix G.

5 Conclusion

We addressed the challenge of Image-induced Fidelity Loss in VLMs, where models trained on predominantly English data tend to default to English responses. Our approach integrates multilingual text-only data into the visual instruction tuning process, preserving the multilingual abilities of the underlying language model.

Experiments show that proportional multilingual integration (TR-3S) significantly reduces English bias while maintaining core multimodal capabilities. We also analyzed data quantity effects, finding that moderate variations in text-to-visual data ratios do not compromise fidelity, though explicit inclusion of target languages remains necessary. Additionally, we introduced a model merging strategy that further improves language fidelity, albeit with some trade-offs in general task performance, highlighting the need for balance in practical applications.

Overall, our findings demonstrate that multilingual textual regularization is a simple and scalable solution to enhance VLM multilingual competence without large multimodal multilingual datasets. This paves the way for future research on optimizing data integration and refining model merging techniques to balance fidelity and overall performance.

639 Limitations

640 Language Coverage

641 While our approach improves multilingual align-
642 ment through text-only supervision, the language
643 coverage remains predominantly European. This
644 raises concerns about the model’s ability to gen-
645 eralize to typologically diverse languages, partic-
646 ularly those with non-Latin scripts (e.g., Arabic,
647 Hindi, Chinese). Future work should explore the
648 integration of a wider array of language families
649 and scripts to validate and expand the method’s
650 applicability.

651 Metric Reliability

652 The fidelity metric (GlotLID and LF+) relies on
653 automatic tools and heuristic judgments, including
654 LLM-as-a-judge assessments that exhibit sensitiv-
655 ity to code-switching and short prompts. Despite
656 efforts to address false positives, such metrics are
657 not infallible and may fail to fully capture semantic
658 fidelity across languages.

659 Ethical Considerations

660 Our research tackles key ethical issues related to
661 multilingual representation and inclusivity in visual
662 language models. Enhancing multilingual capa-
663 bilities promotes accessibility and fairness across
664 diverse linguistic communities.

665 However, relying on machine-translated datasets
666 may introduce biases or cultural inaccuracies. En-
667 suring responsible translation and ongoing refine-
668 ment is crucial.

669 Real-world deployment also demands cultural
670 sensitivity, especially in sectors like education,
671 health, or governance. We emphasize the need
672 for transparency, continuous monitoring, and col-
673 laboration with diverse communities to ensure re-
674 sponsible development and use.

675 References

676 Manoj Acharya, Kushal Kafle, and Christopher Kanan.
677 2019. [Tallyqa: Answering complex counting ques-](#)
678 [tions](#). *Proceedings of the AAAI Conference on Artif-*
679 *cial Intelligence*, 33(01):8076–8084.

680 Divyanshu Aggarwal, Sankarshan Damle, Navin Goyal,
681 Satya Lokam, and Sunayana Sitaram. 2024. [Towards](#)
682 [exploring continual fine-tuning for enhancing lan-](#)
683 [guage ability in large language model](#). In *NeurIPS*
684 *2024 Workshop on Fine-Tuning in Modern Machine*
685 *Learning: Principles and Scalability*.

686 Anton Alexandrov, Veselin Raychev, Mark Niklas
687 Mueller, Ce Zhang, Martin Vechev, and Kristina
688 Toutanova. 2024. [Mitigating catastrophic forgetting](#)
689 [in language transfer via model merging](#). In *Find-*
690 *ings of the Association for Computational Linguistics:*
691 *EMNLP 2024*, pages 17167–17186, Miami, Florida,
692 USA. Association for Computational Linguistics.

693 Duarte Miguel Alves, José Pombal, Nuno M Guerreiro,
694 Pedro Henrique Martins, João Alves, Amin Farajian,
695 Ben Peters, Ricardo Rei, Patrick Fernandes, Sweta
696 Agrawal, Pierre Colombo, José G. C. de Souza, and
697 Andre Martins. 2024. [Tower: An open multilingual](#)
698 [large language model for translation-related tasks](#). In
699 *First Conference on Language Modeling*.

700 Aida Amini, Saadia Gabriel, Shanchuan Lin, Rik
701 Koncel-Kedziorski, Yejin Choi, and Hannaneh Ha-
702 jishirzi. 2019. [MathQA: Towards interpretable math](#)
703 [word problem solving with operation-based for-](#)
704 [malisms](#). In *Proceedings of the 2019 Conference*
705 *of the North American Chapter of the Association for*
706 *Computational Linguistics: Human Language Tech-*
707 *nologies, Volume 1 (Long and Short Papers)*, pages
708 2357–2367, Minneapolis, Minnesota. Association for
709 Computational Linguistics.

710 Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Mar-
711 garet Mitchell, Dhruv Batra, C. Lawrence Zitnick,
712 and Devi Parikh. 2015. [VQA: Visual Question](#)
713 [Answering](#). In *2015 IEEE International Conference on*
714 *Computer Vision (ICCV)*, pages 2425–2433.

715 Daichi Azuma, Taiki Miyanishi, Shuhei Kurita, and
716 Motoaki Kawanabe. 2022. [ScanQA: 3D Question](#)
717 [Answering for Spatial Scene Understanding](#). In *2022*
718 *IEEE/CVF Conference on Computer Vision and Pat-*
719 *tern Recognition (CVPR)*, pages 19107–19117.

720 Haoping Bai, Shancong Mou, Tatiana Likhomanenko,
721 Ramazan Gokberk Cinbis, Oncel Tuzel, Ping Huang,
722 Jiulong Shan, Jianjun Shi, and Meng Cao. 2023. [Vi-](#)
723 [sion datasets: A benchmark for vision-based indus-](#)
724 [trial inspection](#). *arXiv preprint arXiv:2306.07890*.

725 Louis Bethune, David Grangier, Dan Busbridge,
726 Eleonora Gualdoni, Marco Cuturi, and Pierre Ablin.
727 2025. [Scaling laws for forgetting during finetun-](#)
728 [ing with pretraining data injection](#). *arXiv preprint*
729 *arXiv:2502.06042*.

730 Ali Furkan Biten, Rubèn Tito, Andrés Mafla, Lluis
731 Gomez, Marçal Rusiñol, C.V. Jawahar, Ernest Val-
732 veny, and Dimosthenis Karatzas. 2019. [Scene text](#)
733 [visual question answering](#). In *2019 IEEE/CVF Inter-*
734 *national Conference on Computer Vision (ICCV)*,
735 pages 4290–4300.

736 Holger Caesar, Varun Bankiti, Alex H. Lang, Sourabh
737 Vora, Venice Erin Liang, Qiang Xu, Anush Krishnan,
738 Yu Pan, Giancarlo Baldan, and Oscar Beijbom. 2020.
739 [nuscenes: A multimodal dataset for autonomous driv-](#)
740 [ing](#). In *2020 IEEE/CVF Conference on Computer*
741 *Vision and Pattern Recognition (CVPR)*, pages 11618–
742 11628.

743	Jimmy Carter. 2024. Textocr-gpt4v. https://huggingface.co/datasets/jimmycarter/textocr-gpt4v .	798
744		799
745		800
746	Shuaichen Chang, David Palzer, Jialin Li, Eric Fosler-Lussier, and Ningchuan Xiao. 2022a. MapQA: A Dataset for Question Answering on Choropleth Maps. In <i>NeurIPS 2022 First Table Representation Workshop</i> .	801
747		802
748		803
749		804
750		805
751	Yingshan Chang, Guihong Cao, Mridu Narang, Jianfeng Gao, Hisami Suzuki, and Yonatan Bisk. 2022b. WebQA: Multihop and Multimodal QA. In <i>2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pages 16474–16483.	806
752		807
753		808
754		809
755		810
756	Guiming Hardy Chen, Shunian Chen, Ruifei Zhang, Junying Chen, Xiangbo Wu, Zhiyi Zhang, Zhihong Chen, Jianquan Li, Xiang Wan, and Benyou Wang. 2024a. ALLaVA: Harnessing GPT4V-synthesized Data for A Lite Vision-Language Model. <i>Preprint</i> , arXiv:2402.11684.	811
757		812
758		813
759		814
760		815
761		
762	Jiaqi Chen, Tong Li, Jinghui Qin, Pan Lu, Liang Lin, Chongyu Chen, and Xiaodan Liang. 2022. UniGeo: Unifying geometry logical reasoning via reformulating mathematical expression. In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 3313–3323, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	816
763		817
764		818
765		819
766		820
767		821
768		822
769		823
770	Jiaqi Chen, Jianheng Tang, Jinghui Qin, Xiaodan Liang, Lingbo Liu, Eric Xing, and Liang Lin. 2021. GeoQA: A geometric question answering benchmark towards multimodal numerical reasoning. In <i>Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021</i> , pages 513–523, Online. Association for Computational Linguistics.	824
771		825
772		826
773		827
774		828
775		
776		
777	Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Conghui He, Jiaqi Wang, Feng Zhao, and Dahua Lin. 2025. ShareGPT4V: Improving Large Multi-modal Models with Better Captions. In <i>Computer Vision – ECCV 2024</i> , pages 370–387, Cham. Springer Nature Switzerland.	829
778		830
779		831
780		832
781		833
782		834
783	Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Zehui Chen, Haodong Duan, Jiaqi Wang, Yu Qiao, Dahua Lin, and Feng Zhao. 2024b. Are We on the Right Way for Evaluating Large Vision-Language Models? In <i>The Thirty-eighth Annual Conference on Neural Information Processing Systems</i> .	835
784		836
785		837
786		838
787		839
788		840
789		841
790	Lin Chen, Xilin Wei, Jinsong Li, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Zehui Chen, Haodong Duan, Bin Lin, Zhenyu Tang, Li Yuan, Yu Qiao, Dahua Lin, Feng Zhao, and Jiaqi Wang. 2024c. ShareGPT4Video: Improving Video Understanding and Generation with Better Captions. In <i>Advances in Neural Information Processing Systems</i> , volume 37, pages 19472–19495. Curran Associates, Inc.	842
791		
792		
793		
794		
795		
796		
797		
Xi Chen, Josip Djolonga, Piotr Padlewski, Basil Mustafa, Soravit Changpinyo, Jalin Wu, Carlos Riquelme Ruiz, Sebastian Goodman, Xiao Wang, Yi Tay, Siamak Shakeri, Mostafa Dehghani, Daniel Salz, Mario Lucic, Michael Tschannen, Arsha Nagrani, Hexiang Hu, Mandar Joshi, Bo Pang, Ceslee Montgomery, Paulina Pietrzyk, Marvin Ritter, AJ Piergiovanni, Matthias Minderer, Filip Pavetic, Austin Waters, Gang Li, Ibrahim Alabdulmohsin, Lucas Beyer, Julien Amelot, Kenton Lee, Andreas Peter Steiner, Yang Li, Daniel Keysers, Anurag Arnab, Yuanzhong Xu, Keran Rong, Alexander Kolesnikov, Mojtaba Seyedhosseini, Anelia Angelova, Xiaohua Zhai, Neil Houlsby, and Radu Soricut. 2024d. On Scaling Up a Multilingual Vision and Language Model. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pages 14432–14444.	808	
		809
		810
		811
		812
		813
		814
		815
Zi Chen, Xiao Wang, Soravit Changpinyo, AJ Piergiovanni, Piotr Padlewski, Daniel Salz, Sebastian Goodman, Adam Grycner, Basil Mustafa, Lucas Beyer, Alexander Kolesnikov, Joan Puigcerver, Nan Ding, Keran Rong, Hassan Akbari, Gaurav Mishra, Linting Xue, Ashish V Thapliyal, James Bradbury, Weicheng Kuo, Mojtaba Seyedhosseini, Chao Jia, Burcu Karagol Ayan, Carlos Riquelme Ruiz, Andreas Peter Steiner, Anelia Angelova, Xiaohua Zhai, Neil Houlsby, and Radu Soricut. 2023. PaLI: A jointly-scaled multilingual language-image model. In <i>The Eleventh International Conference on Learning Representations</i> .	816	
		817
		818
		819
		820
		821
		822
		823
		824
		825
		826
		827
		828
Zhe Chen, Weiyun Wang, Yue Cao, Yangzhou Liu, Zhangwei Gao, Erfei Cui, Jinguo Zhu, Shenglong Ye, Hao Tian, Zhaoyang Liu, et al. 2024e. Expanding performance boundaries of open-source multimodal models with model, data, and test-time scaling. <i>arXiv preprint arXiv:2412.05271</i> .	829	
		830
		831
		832
		833
		834
Zhoujun Cheng, Haoyu Dong, Zhiruo Wang, Ran Jia, Jiaqi Guo, Yan Gao, Shi Han, Jian-Guang Lou, and Dongmei Zhang. 2022. HiTab: A hierarchical table dataset for question answering and natural language generation. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1094–1110, Dublin, Ireland. Association for Computational Linguistics.	835	
		836
		837
		838
		839
		840
		841
		842
Mike Conover, Matt Hayes, Ankit Mathur, Jianwei Xie, Jun Wan, Sam Shah, Ali Ghodsi, Patrick Wendell, Matei Zaharia, and Reynold Xin. 2023. Free dolly: Introducing the world’s first truly open instruction-tuned llm.	843	
		844
		845
		846
		847
Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barraud, 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	848	
		849
		850
		851
		852
		853
		854
		855
		856
		857

858	Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, Jeff Wang, and NLLB Team. 2024. Scaling neural machine translation to 200 languages . <i>Nature</i> , 630(8018):841–846.	Research (ALVR), pages 7–25, Bangkok, Thailand. Association for Computational Linguistics.	915 916	
863	Erfei Cui, Yinan He, Zheng Ma, Zhe Chen, Hao Tian, Weiyun Wang, Kunchang Li, Yi Wang, Wenhui Wang, Xizhou Zhu, Lewei Lu, Tong Lu, Yali Wang, Limin Wang, Yu Qiao, and Jifeng Dai. 2024. Sharegpt-4o: Comprehensive multimodal annotations with gpt-4o .	Aitor Gonzalez-Agirre, Montserrat Marimon, Carlos Rodriguez-Penagos, Javier Aula-Blasco, Irene Baucells, Carme Armentano-Oller, Jorge Palomar-Giner, Baybars Kulebi, and Marta Villegas. 2024. Building a data infrastructure for a mid-resource language: The case of Catalan . In <i>Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)</i> , pages 2556–2566, Torino, Italia. ELRA and ICCL.	917 918 919 920 921 922 923 924 925 926	
868	Angela Dai, Angel X. Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. 2017. Scannet: Richly-annotated 3d reconstructions of indoor scenes . In <i>2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pages 2432–2443.	Aitor Gonzalez-Agirre, Marc Pàmies, Joan Llop, Irene Baucells, Severino Da Dalt, Daniel Tamayo, José Javier Saiz, Ferran Espuña, Jaume Prats, Javier Aula-Blasco, and others. 2025. Salamandra technical report . <i>arXiv preprint arXiv:2502.08489</i> .	927 928 929 930 931	
874	Wenliang Dai, Nayeon Lee, Boxin Wang, Zhuolin Yang, Zihan Liu, Jon Barker, Tuomas Rintamaki, Mohammad Shoeybi, Bryan Catanzaro, and Wei Ping. 2024. NVLM: Open Frontier-Class Multimodal LLMs . <i>arXiv preprint arXiv:2409.11402</i> .	Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. 2024. The llama 3 herd of models . <i>arXiv preprint arXiv:2407.21783</i> .	932 933 934 935 936	
879	Saurabh Dash, Yiyang Nan, Arash Ahmadian, and John Dang. 2025. A deepdive into aya vision: Advancing the frontier of multilingual multimodality . Hugging Face Blog.	Kristen Grauman, Andrew Westbury, Eugene Byrne, Vincent Cartillier, Zachary Chavis, Antonino Furnari, Rohit Girdhar, Jackson Hamburger, Hao Jiang, Devansh Kukreja, Miao Liu, Xingyu Liu, Miguel Martin, Tushar Nagarajan, Ilija Radosavovic, Santhosh Kumar Ramakrishnan, Fiona Ryan, Jayant Sharma, Michael Wray, Mengmeng Xu, Eric Zhongcong Xu, Chen Zhao, Siddhant Bansal, Dhruv Batra, Sean Crane, Tien Do, Morrie Doulaty, Akshay Erappalli, Christoph Feichtenhofer, Adriano Fragnomeni, Qichen Fu, Abraham Gebreselasie, Cristina Gonzalez, James Hillis, Xuhua Huang, Yifei Huang, Wenqi Jia, Leslie Khoo, Jachym Kolar, Satwik Kottur, Anurag Kumar, Federico Landini, Chao Li, Yanghao Li, Zhenqiang Li, Karttikeya Mangalam, Raghava Modhuguri, Jonathan Munro, Tullie Murrell, Takumi Nishiyasu, Will Price, Paola Ruiz Puentes, Merey Ramazanova, Leda Sari, Kiran Somasundaram, Audrey Southerland, Yusuke Sugano, Ruijie Tao, Minh Vo, Yuchen Wang, Xindi Wu, Takuma Yagi, Ziwei Zhao, Yunyi Zhu, Pablo Arbelaez, David Crandall, Dima Damen, Giovanni Maria Farinella, Christian Fuegen, Bernard Ghanem, Vamsi Krishna Ithapu, C. V. Jawahar, Hanbyul Joo, Kris Kitani, Haizhou Li, Richard Newcombe, Aude Oliva, Hyun Soo Park, James M. Rehg, Yoichi Sato, Jianbo Shi, Mike Zheng Shou, Antonio Torralba, Lorenzo Torresani, Mingfei Yan, and Jitendra Malik. 2024. Ego4D: Around the World in 3,000 Hours of Egocentric Video . <i>IEEE Transactions on Pattern Analysis & Machine Intelligence</i> , 01:1–32.	937 938 939 940 941 942 943 944 945 946 947 948 949 950 951 952 953 954 955 956 957 958 959 960 961 962 963 964 965 966 967	
883	Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. 2021. An image is worth 16x16 words: Transformers for image recognition at scale . In <i>International Conference on Learning Representations</i> .	Stephanie Fu, Netanel Tamir, Shobhita Sundaram, Lucy Chai, Richard Zhang, Tali Dekel, and Phillip Isola. 2023. Dreamsim: Learning new dimensions of human visual similarity using synthetic data . <i>Advances in Neural Information Processing Systems</i> , 36:50742–50768.	Tanmay Gupta, Dustin Schwenk, Ali Farhadi, Derek Hoiem, and Aniruddha Kembhavi. 2018. Imagine this! scripts to compositions to videos . In <i>Proceedings of the European conference on computer vision (ECCV)</i> , pages 598–613.	968 969 970 971 972
888	Maxwell Forbes, Christine Kaeser-Chen, Piyush Sharma, and Serge Belongie. 2019. Neural naturalist: Generating fine-grained image comparisons . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 708–717, Hong Kong, China. Association for Computational Linguistics.	Jiahui Gao, Renjie Pi, Jipeng Zhang, Jiacheng Ye, Wan-jun Zhong, Yufei Wang, Lanqing Hong, Jianhua Han, Hang Xu, Zhenguo Li, and Lingpeng Kong. 2023. G-llava: Solving geometric problem with multi-modal large language model . <i>CoRR</i> , abs/2312.11370.	Danna Gurari, Qing Li, Abigale J. Stangl, Anhong Guo, Chi Lin, Kristen Grauman, Jiebo Luo, and Jeffrey P.	973 974
906	Gregor Geigle, Abhay Jain, Radu Timofte, and Goran Glavaš. 2024. mBLIP: Efficient bootstrapping of multilingual vision-LLMs . In <i>Proceedings of the 3rd Workshop on Advances in Language and Vision</i>			

975	Bigham. 2018. VizWiz Grand Challenge: Answering Visual Questions from Blind People . In <i>2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pages 3608–3617, Los Alamitos, CA, USA. IEEE Computer Society.	1032
976		1033
977		1034
978		1035
979		1036
980		1037
981	Dan Hendrycks and Kevin Gimpel. 2016. Bridging non-linearities and stochastic regularizers with gaussian error linear units . <i>CoRR</i> , abs/1606.08415.	1038
982		1039
983	Musashi Hinck, Carolin Holtermann, Matthew Lyle Olson, Florian Schneider, Sungduk Yu, Anahita Bhawaniwalla, Anne Lauscher, Shao-Yen Tseng, and Vasudev Lal. 2024. Why do LLaVA vision-language models reply to images in English? In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 13402–13421, Miami, Florida, USA. Association for Computational Linguistics.	1040
984		1041
985		1042
986		1043
987		
988	Jinyi Hu, Yuan Yao, Chongyi Wang, SHAN WANG, Yinxu Pan, Qianyu Chen, Tianyu Yu, Hanghao Wu, Yue Zhao, Haoye Zhang, Xu Han, Yankai Lin, Jiao Xue, dahai li, Zhiyuan Liu, and Maosong Sun. 2024. Large multilingual models pivot zero-shot multi-modal learning across languages . In <i>The Twelfth International Conference on Learning Representations</i> .	1044
989		1045
990		1046
991		1047
992		
993	Ting-Hao Kenneth Huang, Francis Ferraro, Nasrin Mostafazadeh, Ishan Misra, Aishwarya Agrawal, Jacob Devlin, Ross Girshick, Xiaodong He, Pushmeet Kohli, Dhruv Batra, C. Lawrence Zitnick, Devi Parikh, Lucy Vanderwende, Michel Galley, and Margaret Mitchell. 2016. Visual storytelling . In <i>Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 1233–1239, San Diego, California. Association for Computational Linguistics.	1048
994		1049
995		1050
996		1051
997		1052
998		1053
999		1054
1000		
1001	Kushal Kafle, Brian Price, Scott Cohen, and Christopher Kanan. 2018. Dvqa: Understanding data visualizations via question answering . In <i>2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 5648–5656.	1055
1002		1056
1003		1057
1004		1058
1005		1059
1006		
1007	Samira Ebrahimi Kahou, Vincent Michalski, Adam Atkinson, Ákos Kádár, Adam Trischler, and Yoshua Bengio. 2017. Figureqa: An annotated figure dataset for visual reasoning . <i>arXiv preprint arXiv:1710.07300</i> .	1060
1008		1061
1009		1062
1010		1063
1011		1064
1012		
1013	Amir Hossein Kargaran, Ayyoob Imani, François Yvon, and Hinrich Schuetze. 2023. GlotLID: Language identification for low-resource languages . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 6155–6218, Singapore. Association for Computational Linguistics.	1065
1014		1066
1015		1067
1016		1068
1017		1069
1018		1070
1019	Mehran Kazemi, Hamidreza Alvari, Ankit Anand, Jialin Wu, Xi Chen, and Radu Soricut. 2024. Geomverse: A systematic evaluation of large models for geometric reasoning . In <i>AI for Math Workshop@ ICML</i> .	1071
1020		1072
1021		1073
1022		1074
1023		
1024	Aniruddha Kembhavi, Mike Salvato, Eric Kolve, Minjoon Seo, Hannaneh Hajishirzi, and Ali Farhadi. 2016. A diagram is worth a dozen images . In <i>Computer Vision – ECCV 2016</i> , pages 235–251, Cham. Springer International Publishing.	1075
1025		1076
1026		1077
1027		1078
1028		1079
1029		
1030	Aniruddha Kembhavi, Minjoon Seo, Dustin Schwenk, Jonghyun Choi, Ali Farhadi, and Hannaneh Hajishirzi. 2017. Are you smarter than a sixth grader? textbook question answering for multimodal machine comprehension . In <i>2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pages 5376–5384.	1080
1031		1081
1032		1082
1033		1083
1034		1084
1035		1085
1036		1086
1037		

1087	Douwe Kiela, Hamed Firooz, Aravind Mohan, Vedanuj Goswami, Amanpreet Singh, Pratik Ringshia, and Davide Testuggine. 2020. The hateful memes challenge: Detecting hate speech in multimodal memes. <i>Advances in neural information processing systems</i> , 33:2611–2624.	1144
1088		1145
1089		1146
1090		1147
1091		1148
1092		1149
1093	Geewook Kim, Teakgyu Hong, Moonbin Yim, JeongYeon Nam, Jinyoung Park, JinYeong Yim, Won-seok Hwang, Sangdoo Yun, Dongyoon Han, and Seunghyun Park. 2022. Ocr-free document understanding transformer. In <i>European Conference on Computer Vision (ECCV)</i> .	1150
1094		1151
1095		1152
1096		1153
1097		1154
1098		1155
1099	Andreas Köpf, Yannic Kilcher, Dimitri Von Rütte, Sotiris Anagnostidis, Zhi Rui Tam, Keith Stevens, Abdullah Barhoum, Duc Nguyen, Oliver Stanley, Richárd Nagyfi, et al. 2023. Openassistant conversations-democratizing large language model alignment. <i>Advances in Neural Information Processing Systems</i> , 36:47669–47681.	1156
1100		1157
1101		
1102		
1103		
1104		
1105		
1106	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. <i>International Journal of Computer Vision</i> , 123(1):32–73.	1158
1107		1159
1108		1160
1109		1161
1110		1162
1111		1163
1112		1164
1113	Benno Krojer, Vaibhav Adlakha, Vibhav Vineet, Yash Goyal, Edoardo Ponti, and Siva Reddy. 2022. Image retrieval from contextual descriptions. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 3426–3440, Dublin, Ireland. Association for Computational Linguistics.	1165
1114		1166
1115		1167
1116		1168
1117		
1118		
1119		
1120	Jason J. Lau, Soumya Gayen, Asma Ben Abacha, and Dina Demner-Fushman. 2018. A dataset of clinically generated visual questions and answers about radiology images. <i>Scientific Data</i> , 5(1):180251.	1169
1121		1170
1122		1171
1123		1172
1124		1173
1125	Hugo Laurençon, Léo Tronchon, and Victor Sanh. 2024. Unlocking the conversion of web screenshots into html code with the websight dataset. <i>Preprint</i> , arXiv:2403.09029.	1174
1126		1175
1127		1176
1128	Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Peiyuan Zhang, Yanwei Li, Ziwei Liu, and Chunyuan Li. 2025. LLaVA-onevision: Easy visual task transfer. <i>Transactions on Machine Learning Research</i> .	1177
1129		1178
1130		1179
1131		1180
1132		1181
1133	Juncheng Li, Kaihang Pan, Zhiqi Ge, Minghe Gao, Wei Ji, Wenqiao Zhang, Tat-Seng Chua, Siliang Tang, Hanwang Zhang, and Yuetong Zhuang. 2023a. Fine-tuning multimodal llms to follow zero-shot demonstrative instructions. <i>arXiv preprint arXiv:2308.04152</i> .	1182
1134		1183
1135		1184
1136		1185
1137		1186
1138		1187
1139	Lei Li, Yuwei Yin, Shicheng Li, Liang Chen, Peiyi Wang, Shuhuai Ren, Mukai Li, Yazheng Yang, Jingjing Xu, Xu Sun, et al. 2023b. M³it: A large-scale dataset towards multi-modal multilingual instruction tuning. <i>arXiv preprint arXiv:2306.04387</i> .	1188
1140		1189
1141		1190
1142		1191
1143		1192
1144	Yitong Li, Zhe Gan, Yelong Shen, Jingjing Liu, Yu Cheng, Yuexin Wu, Lawrence Carin, David Carlson, and Jianfeng Gao. 2019. Storygan: A sequential conditional gan for story visualization. In <i>2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pages 6322–6331.	1193
1145		1194
1146		1195
1147		1196
1148		1197
1149		1200

1201	Multimodal reasoning via thought chains for science question answering. In <i>Advances in Neural Information Processing Systems</i> , volume 35, pages 2507–2521. Curran Associates, Inc.	Jason Obeid and Enamul Hoque. 2020. <i>Chart-to-text: Generating natural language descriptions for charts by adapting the transformer model</i> . In <i>Proceedings of the 13th International Conference on Natural Language Generation</i> , pages 138–147, Dublin, Ireland. Association for Computational Linguistics.	1258
1202			1259
1203			1260
1204			1261
1205	Pan Lu, Liang Qiu, Kai-Wei Chang, Ying Nian Wu, Song-Chun Zhu, Tamay Rajpurohit, Peter Clark, and Ashwin Kalyan. 2023. Dynamic prompt learning via policy gradient for semi-structured mathematical reasoning. In <i>International Conference on Learning Representations (ICLR)</i> .	Dong Huk Park, Trevor Darrell, and Anna Rohrbach. 2019. Robust change captioning. In <i>2019 IEEE/CVF International Conference on Computer Vision (ICCV)</i> , pages 4623–4632.	1262
1206			1263
1207			1264
1208			1265
1209			1266
1210			1267
1211	Pan Lu, Liang Qiu, Jiaqi Chen, Tony Xia, Yizhou Zhao, Wei Zhang, Zhou Yu, Xiaodan Liang, and Song-Chun Zhu. 2021b. Iconqa: A new benchmark for abstract diagram understanding and visual language reasoning. In <i>Advances in Neural Information Processing Systems (NeurIPS)</i> .	Renjie Pi, Jianshu Zhang, Jipeng Zhang, Rui Pan, Zhekai Chen, and Tong Zhang. 2024. Image textualization: An automatic framework for creating accurate and detailed image descriptions. <i>arXiv preprint arXiv:2406.07502</i> .	1268
1212			1269
1213			1270
1214			1271
1215			1272
1216			
1217	Muhammad Maaz, Hanoona Rasheed, Abdelrahman Shaker, Salman Khan, Hisham Cholakal, Rao M Anwer, Tim Baldwin, Michael Felsberg, and Fahad S Khan. 2024. Palo: A polyglot large multimodal model for 5b people. <i>arXiv preprint arXiv:2402.14818</i> .	Maja Popović. 2016. chrF deconstructed: beta parameters and n-gram weights. In <i>Proceedings of the First Conference on Machine Translation: Volume 2, Shared Task Papers</i> , pages 499–504, Berlin, Germany. Association for Computational Linguistics.	1273
1218			1274
1219			1275
1220			1276
1221			1277
1222			
1223	Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. 2019. OK-VQA: A Visual Question Answering Benchmark Requiring External Knowledge. In <i>Proceedings of the IEEE/cvpr conference on computer vision and pattern recognition</i> , pages 3195–3204.	Maja Popović. 2017. chrF++: words helping character n-grams. In <i>Proceedings of the Second Conference on Machine Translation</i> , pages 612–618, Copenhagen, Denmark. Association for Computational Linguistics.	1278
1224			1279
1225			1280
1226			1281
1227			1282
1228			
1229	U.-V. Marti and H. Bunke. 2002. The iam-database: an english sentence database for offline handwriting recognition. <i>International Journal on Document Analysis and Recognition</i> , 5(1):39–46.	Chen Qiu, Dan Oneață, Emanuele Bugliarello, Stella Frank, and Desmond Elliott. 2022. Multilingual multimodal learning with machine translated text. In <i>Findings of the Association for Computational Linguistics: EMNLP 2022</i> , pages 4178–4193, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	1283
1230			1284
1231			1285
1232			1286
1233	Ahmed Masry, Xuan Long Do, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. 2022. ChartQA: A benchmark for question answering about charts with visual and logical reasoning. In <i>Findings of the Association for Computational Linguistics: ACL 2022</i> , pages 2263–2279, Dublin, Ireland. Association for Computational Linguistics.	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastri, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In <i>International conference on machine learning</i> , pages 8748–8763. PMLR.	1287
1234			1288
1235			1289
1236			
1237			
1238			
1239			
1240	Minesh Mathew, Viraj Bagal, Rubèn Tito, Dimosthenis Karatzas, Ernest Valveny, and C. V. Jawahar. 2022. InfographicVQA. In <i>2022 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)</i> , pages 2582–2591.	Nazneen Rajani, Lewis Tunstall, Edward Beeching, Nathan Lambert, Alexander M. Rush, and Thomas Wolf. 2023. No robots. https://huggingface.co/datasets/HuggingFace4/no_robots .	1296
1241			1297
1242			1298
1243			1299
1244			
1245	Minesh Mathew, Dimosthenis Karatzas, and C. V. Jawahar. 2021. Docvqa: A dataset for vqa on document images. In <i>2021 IEEE Winter Conference on Applications of Computer Vision (WACV)</i> , pages 2199–2208.	Hareesh Ravi, Kushal Kafle, Scott Cohen, Jonathan Brandt, and Mubbasis Kapadia. 2021. Aesop: Abstract encoding of stories, objects, and pictures. In <i>2021 IEEE/CVF International Conference on Computer Vision (ICCV)</i> , pages 2032–2043.	1300
1246			1301
1247			1302
1248			1303
1249	Anand Mishra, Karteek Alahari, and CV Jawahar. 2012. Scene text recognition using higher order language priors. In <i>BMVC-British machine vision conference</i> . BMVA.	Florian Schneider and Sunayana Sitaram. 2024. M5 – a diverse benchmark to assess the performance of large multilingual models across multilingual and multicultural vision-language tasks. In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 4309–4345, Miami, Florida, USA. Association for Computational Linguistics.	1304
1250			1305
1251			1306
1252			1307
1253			1308
1254			1309
1255			1310
1256			1311
1257	Anand Mishra, Shashank Shekhar, Ajeet Kumar Singh, and Anirban Chakraborty. 2019. Ocr-vqa: Visual question answering by reading text in images. In <i>2019 International Conference on Document Analysis and Recognition (ICDAR)</i> , pages 947–952.	Dustin Schwenk, Apoorv Khandelwal, Christopher Clark, Kenneth Marino, and Roozbeh Mottaghi. 2022.	1312
			1313

1314	A-okvqa: A benchmark for visual question answering using world knowledge. In <i>European conference on computer vision</i> , pages 146–162. Springer.	1370
1315		1371
1316		1372
1317		1373
1318		1374
1319		
1320	Minjoon Seo, Hannaneh Hajishirzi, Ali Farhadi, Oren Etzioni, and Clint Malcolm. 2015. Solving geometry problems: Combining text and diagram interpretation. In <i>Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing</i> , pages 1466–1476, Lisbon, Portugal. Association for Computational Linguistics.	
1321		
1322		
1323		
1324	Ken Shoemake. 1985. Animating rotation with quaternion curves. In <i>Proceedings of the 12th annual conference on Computer graphics and interactive techniques</i> , pages 245–254. ACM.	
1325		
1326		
1327		
1328		
1329	Mohit Shridhar, Jesse Thomason, Daniel Gordon, Yonatan Bisk, Winson Han, Roozbeh Mottaghi, Luke Zettlemoyer, and Dieter Fox. 2020. Alfred: A benchmark for interpreting grounded instructions for everyday tasks. In <i>2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pages 10737–10746.	
1330		
1331		
1332		
1333		
1334		
1335	Oleksii Sidorov, Ronghang Hu, Marcus Rohrbach, and Amanpreet Singh. 2020. Textcaps: a dataset for image captioning with reading comprehension. In <i>Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part II 16</i> , pages 742–758. Springer.	
1336		
1337		
1338		
1339		
1340		
1341	Gunnar A. Sigurdsson, Gülcin Varol, Xiaolong Wang, Ali Farhadi, Ivan Laptev, and Abhinav Gupta. 2016. Hollywood in homes: Crowdsourcing data collection for activity understanding. In <i>Computer Vision – ECCV 2016</i> , pages 510–526, Cham. Springer International Publishing.	
1342		
1343		
1344		
1345		
1346		
1347	Shivalika Singh, Freddie Vargas, Daniel D’souza, Börje Karlsson, Abinaya Mahendiran, Wei-Yin Ko, Herumb Shandilya, Jay Patel, Deividas Mataciunas, Laura O’Mahony, Mike Zhang, Ramith Hetiarachchi, Joseph Wilson, Marina Machado, Luisa Moura, Dominik Krzeminski, Hakimeh Fadaee, Irem Ergun, Ifeoma Okoh, Aisha Alaagib, Oshan Mudannayake, Zaid Alyafeai, Vu Chien, Sebastian Ruder, Surya Guthikonda, Emad Alghamdi, Sebastian Gehrmann, Niklas Muennighoff, Max Bartolo, Julia Kreutzer, Ahmet Üstün, Marzieh Fadaee, and Sara Hooker. 2024. Aya dataset: An open-access collection for multilingual instruction tuning. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 11521–11567, Bangkok, Thailand. Association for Computational Linguistics.	
1348		
1349		
1350		
1351		
1352		
1353		
1354		
1355		
1356		
1357		
1358		
1359		
1360		
1361		
1362		
1363		
1364	Yueqi Song, Simran Khanuja, and Graham Neubig. 2024. What is missing in multilingual visual reasoning and how to fix it. <i>arXiv preprint arXiv:2403.01404</i> .	
1365		
1366		
1367		
1368	Alane Suhr, Stephanie Zhou, Ally Zhang, Iris Zhang, Huajun Bai, and Yoav Artzi. 2019. A corpus for reasoning about natural language grounded in photographs. In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 6418–6428, Florence, Italy. Association for Computational Linguistics.	
1369		
1370		
1371		
1372		
1373		
1374		
1375	Hao Tan, Franck Dernoncourt, Zhe Lin, Trung Bui, and Mohit Bansal. 2019. Expressing visual relationships via language. In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 1873–1883, Florence, Italy. Association for Computational Linguistics.	
1376		
1377		
1378		
1379		
1380		
1381	Ryota Tanaka, Kyosuke Nishida, and Sen Yoshida. 2021. VisualMRC: Machine Reading Comprehension on Document Images. <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , 35(15):13878–13888.	
1382		
1383		
1384		
1385	Benny Tang, Angie Boggust, and Arvind Satyanarayan. 2023. VisText: A benchmark for semantically rich chart captioning. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 7268–7298, Toronto, Canada. Association for Computational Linguistics.	
1386		
1387		
1388		
1389		
1390		
1391		
1392	Ashish V. Thapliyal, Jordi Pont Tuset, Xi Chen, and Radu Soricut. 2022. Crossmodal-3600: A massively multilingual multimodal evaluation dataset. In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 715–729, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	
1393		
1394		
1395		
1396		
1397		
1398		
1399	Peter Tong, Ellis Brown, Penghao Wu, Sanghyun Woo, Adithya Jairam Vedagiri IYER, Sai Charitha Akula, Shusheng Yang, Jihan Yang, Manoj Middepogu, Ziteng Wang, et al. 2024. Cambrian-1: A fully open, vision-centric exploration of multimodal llms. <i>Advances in Neural Information Processing Systems</i> , 37:87310–87356.	
1400		
1401		
1402		
1403		
1404		
1405		
1406	Bryan Wang, Gang Li, Xin Zhou, Zhourong Chen, Tovi Grossman, and Yang Li. 2021. Screen2words: Automatic mobile ui summarization with multimodal learning. In <i>The 34th Annual ACM Symposium on User Interface Software and Technology</i> , pages 498–510.	
1407		
1408		
1409		
1410		
1411		
1412	Haoning Wu, Hanwei Zhu, Zicheng Zhang, Erli Zhang, Chaofeng Chen, Liang Liao, Chunyi Li, Annan Wang, Wenxiu Sun, Qiong Yan, et al. 2024. Towards open-ended visual quality comparison. In <i>European Conference on Computer Vision</i> , pages 360–377. Springer.	
1413		
1414		
1415		
1416		
1417		
1418	Junbin Xiao, Xindi Shang, Angela Yao, and Tat-Seng Chua. 2021. NExT-QA: Next Phase of Question-Answering to Explaining Temporal Actions . In <i>2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pages 9772–9781, Los Alamitos, CA, USA. IEEE Computer Society.	
1419		
1420		
1421		
1422		
1423		
1424	Zhangchen Xu, Fengqing Jiang, Luyao Niu, Yun-tian Deng, Radha Poovendran, Yejin Choi, and Bill Yuchen Lin. 2024a. Magpie: Alignment data	
1425		
1426		

1427		synthesis from scratch by prompting aligned llms with nothing. <i>arXiv preprint arXiv:2406.08464</i> .	1483
1428			1484
1429	Zhiyang Xu, Chao Feng, Rulin Shao, Trevor Ashby, Ying Shen, Di Jin, Yu Cheng, Qifan Wang, and Lifu Huang. 2024b. Vision-Flan: Scaling Human-Labeled Tasks in Visual Instruction Tuning . In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 15271–15342, Bangkok, Thailand. Association for Computational Linguistics.	1485	
1430			1486
1431			1487
1432			1488
1433			1489
1434			1490
1435			1491
1436	Semih Yagcioglu, Aykut Erdem, Erkut Erdem, and Nazli Ikizler-Cinbis. 2018. RecipeQA: A challenge dataset for multimodal comprehension of cooking recipes . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 1358–1368, Brussels, Belgium. Association for Computational Linguistics.	1492	
1437			1493
1438			1494
1439			1495
1440			1496
1441			1497
1442			1498
1443	Jiabo Ye, Anwen Hu, Haiyang Xu, Qinghao Ye, Ming Yan, Guohai Xu, Chenliang Li, Junfeng Tian, Qi Qian, Ji Zhang, Qin Jin, Liang He, Xin Lin, and Fei Huang. 2023. UReader: Universal OCR-free visually-situated language understanding with multimodal large language model . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 2841–2858, Singapore. Association for Computational Linguistics.	1499	
1444			1500
1445			1501
1446			1502
1447			1503
1448			1504
1449			1505
1450			1506
1451			1507
1452	Licheng Yu, Patrick Poirson, Shan Yang, Alexander C Berg, and Tamara L Berg. 2016. Modeling context in referring expressions . In <i>Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part II 14</i> , pages 69–85. Springer.	1508	
1453			1509
1454			1510
1455			1511
1456			1512
1457			1513
1458	Zhou Yu, Dejing Xu, Jun Yu, Ting Yu, Zhou Zhao, Yuet-ing Zhuang, and Dacheng Tao. 2019. Activitynet-qa: A dataset for understanding complex web videos via question answering . In <i>AAAI</i> , pages 9127–9134.	1514	
1459			1515
1460			1516
1461			1517
1462	Ye Yuan, Xiao Liu, Wondimu Dikubab, Hui Liu, Zhilong Ji, Zhongqin Wu, and Xiang Bai. 2022. Syntax-aware network for handwritten mathematical expression recognition . In <i>2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pages 4543–4552.	1518	
1463			1519
1464			1520
1465			1521
1466			1522
1467			1523
1468	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, Wenhai Huang, Huan Sun, Yu Su, and Wenhua Chen. 2024. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi . In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pages 9556–9567.	1524	
1469			1525
1470			1526
1471			1527
1472			1528
1473			1529
1474			1530
1475			1531
1476			1532
1477			1533
1478			1534
1479	Xiang Yue, Yueqi Song, Akari Asai, Seungone Kim, Jean de Dieu Nyandwi, Simran Khanuja, Anjali Kantharuban, Lintang Sutawika, Sathyaranayanan Ramamoorthy, and Graham Neubig. 2025. Pangea: A	1535	
1480			1536
1481			1537
1482			1538

A Data Sources

This section showcases the visual datasets (Table 10) and the text-only datasets (Table 11) used throughout this work.

1537 B Language Consistency Evaluation via 1538 LLM-as-a-judge

1539 We have observed cases of code-switched genera-
1540 tions, and in some cases, even if most of the sen-
1541 tence is generated in the target language, a few
1542 words may still appear in English. The primary
1543 goal of the textual LLM-as-a-judge evaluation is to
1544 address GlotLID’s limitation to classify these cases
1545 as incorrect.

1546 To effectively evaluate each sentence, we have
1547 designed a prompt (see Figure 4) that instructs the
1548 evaluator to perform multiple tasks beyond basic
1549 language identification, enabling the computation
1550 of a language consistency score:

- 1551 1. Guess the language of the sentence (to com-
1552 pare with GlotLID, even if we will use
1553 GlotLID’s outputs).
- 1554 2. Assign a language consistency score (between
1555 0 and 1).
- 1556 3. Determine whether the sentence is fully in
1557 the target language (a boolean value, where
1558 False indicates that at least one word appears
1559 in another language).
- 1560 4. Generate a summary explaining the decisions
1561 made by the model.

1562 This method allows us to evaluate the language
1563 consistency of our model at the word level from
1564 different perspectives, both through a numerical
1565 score and a boolean indicator.

1566 For additional robustness, we have performed
1567 the text evaluations using three different generation
1568 configurations (defined in Table 6) and then applied
1569 a majority voting.

Conf.	Temperature	Top_p	Max new tokens
A	0.6	0.7	50
B	0.8	0.6	50
C	1.0	0.5	50

1570 Table 6: Generation parameter settings for LLM-as-a-
1571 Judge evaluation.

1572 As shown in Figure 4, a structured response for-
1573 mat was explicitly requested to ensure that each
1574 field could be reliably extracted from the output.
1575 The responses that did not conform to the expected
1576 format were replaced with "N/A".

1577 To compute the final scores, we averaged the re-
1578 sults across the three different configurations. For

1579 numerical scores, when a value was missing, we
1580 computed the average using the available values.
1581 For boolean scores, we applied a majority voting
1582 approach. In cases where one score was missing
1583 and the remaining two were True and False, we
1584 defaulted to False as the final instance score. This
1585 ensures that our results provide a lower bound, mak-
1586 ing the evaluation more conservative and reliable.
1587 All the results are shown in Table 7.

1588 B.1 Judging the Judge: Evaluation of the 1589 LLM-as-a-judge

1590 Even though we have implemented measures for
1591 robustness (such as using different generation con-
1592 figurations for the LLM-as-a-judge evaluators and
1593 requiring evaluators to provide summaries to justify
1594 their scores) we aim to further ensure the fairness
1595 of the provided scores by evaluating them against
1596 the reference captions in the dataset.

1597 To assess the reliability of the language consis-
1598 tency evaluator model, we scored the reference
1599 captions from the evaluation dataset (crossmodal-
1600 3600) using the same evaluation process applied
1601 to the captions generated by our models. We have
1602 evaluated all the available references (up to 3 per
1603 image) and only used a single generation configura-
1604 tion (configuration B in Table 6).

1605 As shown in Table 7, most of the obtained
1606 scores are above 90%, demonstrating the ef-
1607 fectiveness of the chosen model as an evalua-
1608 tor. However, some minor errors are present,
1609 which can be attributed to multilingual limita-
1610 tions. Llama-3.1-8B-Instruct officially sup-
1611 ports seven languages in addition to English
1612 (French, German, Hindi, Italian, Portuguese, Span-
1613 ish, and Thai). While this allows it to handle most
1614 European languages, it is expected that the model
1615 may occasionally struggle with languages outside
1616 its primary training set, leading to some misclassi-
1617 fications.

1618 In terms of language consistency, we have dis-
1619 carded selected the binary score due to its higher
1620 scores. The evaluator classifies the samples as cor-
1621 rect more than the 95% of times in the majority
1622 of languages. The lower performance in German
1623 (91.7%) can be attributed to the fact that it is a
1624 Germanic language that shares a large amount of
1625 words with English, what may induce classification
1626 errors. In the case of Russian (83.3%), the lower
1627 performance may be explained by the limited sup-
1628 port for languages using the Cyrillic alphabet in the
1629 LLM.

Model	Lang.	GL	LLM-L	LLM-S	LLM-B	Model	Lang.	GL	LLM-L	LLM-S	LLM-B
Normal Models											
BM	De	2.7	2.4	85.8	73.5	TR-1S	De	24.5	18.3	84.0	75.9
	En	100.0	98.4	99.0	99.3		En	100.0	98.1	98.7	99.1
	Es	4.4	4.7	91.5	87.2		Es	38.4	37.9	94.4	95.4
	Fr	12.2	10.8	90.3	88.1		Fr	29.9	26.8	91.2	91.1
	Nl	5.6	6.02	93.8	92.5		Nl	49.2	48.1	93.4	93.1
	No	5.2	3.84	92.0	89.6		No	46.8	30.3	92.8	91.5
	Ru	3.8	2.75	82.2	58.5		Ru	50.9	28.4	82.2	62.9
TR-2S	De	81.3	54.4	81.9	73.2	TR-3S	De	88.7	61.3	82.8	72.7
	En	100.0	98.3	98.9	99.4		En	100.0	98.7	99.1	99.5
	Es	65.4	61.8	94.1	93.7		Es	92.9	86.8	94.6	94.3
	Fr	74.9	63.8	90.5	89.5		Fr	85.7	71.4	89.6	86.7
	Nl	91.3	84.9	93.4	94.3		Nl	91.8	86.7	92.9	93.2
	No	52.0	40.4	93.7	93.3		No	69.5	49.5	92.6	91.0
	Ru	24.8	13.4	81.6	63.1		Ru	52.9	31.7	83.1	59.6
TR-3S-0.25x	De	85.0	58.4	82.1	75.2	TR-3S-0.5x	De	88.9	58.9	81.2	72.8
	En	100.0	98.1	98.8	99.1		En	100.0	98.4	98.7	99.4
	Es	91.9	85.9	94.1	94.8		Es	92.4	85.9	94.0	93.7
	Fr	88.6	76.3	90.7	89.8		Fr	83.5	71.2	90.0	88.8
	Nl	69.8	65.3	93.8	95.3		Nl	96.0	90.6	94.1	96.0
	No	53.2	41.2	94.0	94.6		No	82.7	62.6	93.5	93.6
	Ru	91.5	48.5	81.6	63.2		Ru	81.9	43.3	81.1	62.7
TR-3S-2x	De	73.3	48.6	81.0	72.8						
	En	100.0	98.7	98.7	99.4						
	Es	76.4	71.7	94.3	95.2						
	Fr	69.5	59.1	90.0	88.6						
	Nl	93.5	88.3	93.8	95.9						
	No	67.1	53.7	93.9	94.0						
	Ru	93.3	48.9	81.3	64.2						
Merged Models (+M)											
BM+M	De	15.0	12.8	84.1	76.6	TR-1S+M	De	76.5	58.3	84.5	77.4
	En	100.0	98.1	98.9	99.5		En	100.0	98.8	98.8	99.7
	Es	37.0	36.9	93.4	91.5		Es	86.3	81.9	93.6	93.7
	Fr	47.4	40.0	89.8	86.8		Fr	87.9	74.3	90.2	89.4
	Nl	9.3	10.0	93.1	92.5		Nl	63.7	65.3	92.0	91.4
	No	14.7	10.9	91.8	90.4		No	58.1	47.3	89.5	87.6
	Ru	15.4	10.8	82.2	60.8		Ru	77.7	56.1	85.2	73.1
TR-2S+M	De	95.1	66.5	84.1	77.1	TR-3S+M	De	94.1	69.6	84.4	78.7
	En	100.0	97.8	98.8	99.2		En	100.0	98.0	98.9	99.2
	Es	97.4	89.1	94.8	95.3		Es	96.4	88.9	94.5	94.4
	Fr	97.7	83.1	91.8	90.1		Fr	95.5	80.8	91.1	89.4
	Nl	95.6	88.8	92.7	93.1		Nl	96.1	90.5	92.3	93.0
	No	97.6	72.7	90.1	85.7		No	90.1	70.9	91.7	90.7
	Ru	96.7	60.8	85.2	74.0		Ru	92.4	65.7	86.7	76.2
Crossmodal-3600 reference samples (Evaluating the LLM-as-a-judge)											
references	De	100.0	83.3	85.4	91.7						
	En	99.7	99.2	97.7	99.5						
	Es	99.6	97.2	95.2	99.0						
	Fr	99.9	90.5	90.8	95.2						
	Nl	99.4	95.8	94.5	98.7						
	No	99.1	91.4	93.5	97.9						
	Ru	99.9	76.0	79.8	83.3						

Table 7: Model comparison showing GlotLID detection percentages for the target language (GL) and scores obtained using LLM-as-a-judge. LLM-L represents the target language detection, LLM-S indicates the numerical language consistency score, and LLM-B denotes the binary language consistency score. The GlotLID+LLM score is calculated as the product of the GlotLID score and LLM-B: GlotLID+LLM = GlotLID × LLM-B. For each group except the LLM-as-a-judge evaluator evaluation, the highest scores are marked in bold.

Language Consistency Evaluation Prompt Template

Analyze the following text and determine the language it is written in.

- Identify the most likely language.
- Ensure the probability score is a single value, not a range or estimate.
- Determine a language consistency score between 0 and 1.0, where 1.0 means the text is entirely in one language, and 0.0 means it is completely incomprehensible.
- Lower the score proportionally if foreign words are present, but do not assign 0.0 unless the text is nonsensical.
- The language score must be a single number between 0 and 1.0.
- Indicate whether the text is completely written in the identified language (True or False).
- In both language consistency metrics, do not penalize for proper nouns, brand names, or commonly used foreign terms (e.g., 'software', 'email') that do not alter the overall language structure.
- Avoid unnecessary explanations. Summarize the feedback (reason of the mark) in at most 30 words.

Use the exact format below:

- Language: [language_guess]
- Language Score: [single value between 0 and 1.0]
- Fully in Language: [True/False]
- Summary: [Concise explanation (max 30 words)]

Keep your answer short and concise. The sentence to analyze is the following:

<CAPTION GENERATED BY VISUAL SALAMANDRA>

Figure 4: Prompt used to evaluate language consistency via LLM-as-a-judge. The evaluator model assesses the language fidelity of the caption generated by the VLM using multiple criteria. Note that this evaluation focuses solely on language fidelity, not the overall quality of the caption.

C Caption Generation Prompt

The Crossmodal-3600 dataset does not specify an explicit prompt for caption generation. However, in their work they provide instructions for generating captions, which we used as a guideline. Based on these instructions, we formulated a simplified captioning approach. The prompts used for generating image captions in our evaluation are presented in Figure 5. The same prompt was applied consistently across all selected languages.

D Metrics Discussion

In this study, we chose not to use teacher forcing loss or perplexity as evaluation metrics due to their inherent limitations in interpretability and comparative analysis across models.

We selected chrF++ as our primary evaluation metric rather than BLEU or ROUGE due to its suitability for multilingual assessments. BLEU relies heavily on exact n-gram matching, often pe-

nalizing legitimate linguistic variations common in multilingual contexts, while ROUGE primarily measures recall and is optimized for summarization tasks, making it suboptimal for assessing generative multilingual output quality. In contrast, chrF++ evaluates based on character-level n-gram overlaps, accommodating linguistic diversity and morphological richness across multiple languages, thus providing a more robust and linguistically sensitive assessment for multilingual visual language models.

Additionally, during the study, we employed VLM-as-a-judge to evaluate the quality of multilingual generations. However, we found that chrF++ effectively addressed the limitations related to multilingual performance evaluation inherent in other metrics, thereby serving as a comprehensive solution for our assessment needs.

		Stage-1	Stage-1.5	Stage-2	OneVision
Vision	Resolution	384	AnyRes Max 5	AnyRes Max 9	AnyRes Max 9
Data	# Tokens	729	Max 729 × 5	Max 729 × 10	Max 729 × 10
Model	Dataset	Single-Image 558K	Single-Image 3.8M	Single-Image 3.1M	Single/Multi-Image, Video 1.6M
Trainable	Projector 20.0M	Full Model 8.2B	Full Model 8.2B	Full Model 8.2B	Full Model 8.2B
Training	Batch Size	128	64	64	64
	LR: ψ_{vision}	1×10^{-3}	2×10^{-6}	2×10^{-6}	2×10^{-6}
	LR: $\{\theta_{proj}, \phi_{LLM}\}$	1×10^{-3}	1×10^{-5}	1×10^{-5}	1×10^{-5}
	Epoch	1	1	1	1
	Warmup Ratio	0.03	0.03	0.03	0.03
	LR Scheduler	Cosine	Cosine	Cosine	Cosine
	Grad. Accum.	1	2	2	2

Table 8: Detailed configuration for each training stage of the LLaVA-OneVision model. For a detailed explanation of AnyRes Max, refer to (Li et al., 2025). Anyres Max 5: $384 \times \{2 \times 2, 1 \times \{2, 3\}, 2, 3 \times 1\}$. AnyRes Max 9: $384 \times \{1 \times 1\}, \dots, \{6 \times 6\}$.

E Training Hyperparameters

The training hyperparameters used during the training of the models evaluated in this work are detailed in Table 8.

F Alternative Merging Strategies

To better understand the trade-offs involved in model merging, we conducted a series of additional experiments comparing different interpolation methods and weight ratios. In particular, we investigated:

Linear Interpolation (lerp) This method interpolates model weights using the standard formula $w = (1 - \alpha)w_1 + \alpha w_2$, where w_1 and w_2 are the weights of the visually instructed and backbone models, respectively, and α is the interpolation ratio.

Spherical Linear Interpolation (slerp) Unlike lerp, slerp (Shoemake, 1985) interpolates weights along a great arc on the hypersphere, preserving the norm and relative directionality. It is computed as:

$$\text{slerp}(w_1, w_2, \alpha) = \frac{\sin((1 - \alpha)\theta)}{\sin(\theta)} w_1 + \frac{\sin(\alpha\theta)}{\sin(\theta)} w_2 \quad (1)$$

where θ is the angle between the two weight vectors. This method can yield smoother transitions in weight space, especially when the models differ significantly.

We evaluated merged models using both interpolation methods under three weighting scenarios:

- **50–50**, giving equal weight to the visually instructed and original backbone models.

- **75–25**, favoring the original backbone to preserve pretrained language capabilities.

- **25–75**, prioritizing the visually instructed model to reinforce vision-language alignment.

The results are presented in Table 9, showing both language fidelity and downstream task performance.

Discussion. As shown, asymmetric merges favoring the backbone (e.g., lerp_075 and slerp_075) achieve near-perfect language fidelity but show weaker performance in multimodal benchmarks. Conversely, merges favoring the visually instructed model (e.g., lerp_025 and slerp_025) lead to substantially improved task performance, but at the cost of lower fidelity in certain languages. The slerp_050 model—corresponding to our main TR-3S+M—offers a more balanced trade-off.

Overall, no single merging configuration yields a clearly optimal trade-off. The best strategy depends on the intended use case: 75% visually instructed weights are preferable when multilingual fidelity is critical, while 25% weights better support general multimodal performance.

G Generation Examples

Figures 6–9 present examples generated using the TR-3S-0.05x model across various languages and diverse tasks.

Table 9: Merged model results: language fidelity (left) and task performance (right). lerp_050 and slerp_050 correspond to 50–50 merges using linear and spherical interpolation, respectively; slerp_050 corresponds to our main model TR-3S+M. lerp_075 and slerp_075 are asymmetric 75–25 merges favoring the base text model. lerp_025 and slerp_025 invert this ratio to prioritize the visually instructed model.

Model	DE	ES	FR	NL	RU	Avg.	AI2D	MMMU	MMStar	RWQA	Avg.
lerp_075	98.75	99.89	99.92	99.64	98.66	99.37	41.84	27.00	32.56	37.39	34.70
slerp_075	99.08	99.86	99.89	99.64	99.02	99.50	42.16	27.33	32.66	36.99	34.79
lerp_050	92.83	96.42	95.83	94.78	91.11	94.19	57.03	34.44	42.06	52.03	46.39
slerp_050	94.10	96.40	95.50	96.10	92.40	94.90	57.19	34.11	42.25	52.16	46.43
lerp_025	89.32	93.06	91.06	92.83	63.56	85.97	72.51	34.89	47.13	55.95	52.62
slerp_025	90.11	93.53	90.81	94.44	66.30	87.44	72.38	34.89	47.27	56.08	52.66

Story Generation In Figure 6, the model is prompted to generate a story from an image. It accurately reads text within the image to identify characters and establish the setting, demonstrating its ability to craft diverse narratives consistently across different languages—even when the text is in English.

Image Description Figure 7 presents a brief image description task. Although the image shows a salamander perched on a person’s hand, some language outputs mistakenly label it as an insect or a predator. Despite these inaccuracies, the descriptions remain largely appropriate.

OCR and Translation Figure 8 showcases a task combining OCR with translation. The model extracts text from an image and then translates it into a target language. This two-step process: OCR followed by translation, highlights the model’s ability to merge visual analysis with its linguistic capabilities. Minor errors do occur, particularly in languages not extensively represented during training, resulting in slightly erroneous translations or defaulting to English.

Multi-Image Reasoning Finally, Figure 9 illustrates a multi-image scenario where the model must comprehend the content of several images and reason to provide an appropriate answer. This example further confirms the model’s effectiveness in real-world applications.

In general, these examples demonstrate how the VLMs instructed via our approach perform optimally across a wide range of tasks, especially for languages where text-only data was incorporated during the visual instruction process.

G.1 Code Switching in Caption Generation

Figure 10 presents examples of code switching observed during caption generation for the

Crossmodal-3600 dataset. We could identify two primary patterns emerge:

- **Independent Words:** Certain technical or less common words are generated in English.
- **Language Alternation:** In some cases, once a word is switched to English, all subsequent words continue in English.

Caption Generation Prompts

[English]

Give me a brief summary of the following image, without too many details. The description should be general and have a maximum of 10 words. To do this, identify the most relevant object or person in the image, the main relationship between the highlighted objects, the most important activity represented, the most outstanding attributes of the main object or person, and the context in which the scene takes place. Then, synthesize everything into a single descriptive and concise sentence, without including additional text.

[Spanish]

Dame un breve resumen de la siguiente imagen, sin dar muchos detalles. La descripción debe ser general y tener un máximo de 10 palabras. Para ello, identifica el objeto o persona más relevante en la imagen, la relación principal entre los objetos destacados, la actividad más importante representada, los atributos más sobresalientes del objeto o persona principal y el contexto en el que ocurre la escena. Luego, sintetiza todo en una sola frase descriptiva y concisa, sin incluir texto adicional.

[French]

Donne-moi un bref résumé de l'image suivante, sans trop de détails. La description doit être générale et contenir un maximum de 10 mots. Pour cela, identifie l'objet ou la personne la plus importante dans l'image, la relation principale entre les objets mis en avant, l'activité la plus significative représentée, les attributs les plus marquants de l'objet ou de la personne principale et le contexte dans lequel la scène se déroule. Ensuite, synthétise tout en une seule phrase descriptive et concise, sans ajouter de texte supplémentaire.

[German]

Gib mir eine kurze Zusammenfassung des folgenden Bildes, ohne zu viele Details. Die Beschreibung sollte allgemein sein und maximal 10 Wörter umfassen. Identifiziere dazu das relevanteste Objekt oder die wichtigste Person im Bild, die Hauptbeziehung zwischen den hervorgehobenen Objekten, die wichtigste dargestellte Aktivität, die auffälligsten Merkmale des Hauptobjekts oder der Hauptperson und den Kontext, in dem die Szene stattfindet. Fasse dann alles in einem einzigen prägnanten und beschreibenden Satz zusammen, ohne zusätzlichen Text hinzuzufügen.

[Italian]

Dammi un breve riassunto della seguente immagine, senza troppi dettagli. La descrizione deve essere generale e avere un massimo di 10 parole. Per farlo, identifica l'oggetto o la persona più rilevante nell'immagine, la relazione principale tra gli oggetti evidenziati, l'attività più importante rappresentata, gli attributi più evidenti dell'oggetto o della persona principale e il contesto in cui si svolge la scena. Quindi, sintetizza tutto in un'unica frase descrittiva e concisa, senza includere testo aggiuntivo.

(Continues on next page...)

[Dutch]

Geef me een korte samenvatting van de volgende afbeelding, zonder te veel details. De beschrijving moet algemeen zijn en maximaal 10 woorden bevatten. Identificeer hiervoor het meest relevante object of de belangrijkste persoon in de afbeelding, de hoofdrelatie tussen de uitgelichte objecten, de belangrijkste weergegeven activiteit, de meest opvallende kenmerken van het hoofdobject of de belangrijkste persoon en de context waarin de scène zich afspeelt. Vat vervolgens alles samen in één beschrijvende en beknopte zin, zonder extra tekst toe te voegen.

[Portuguese]

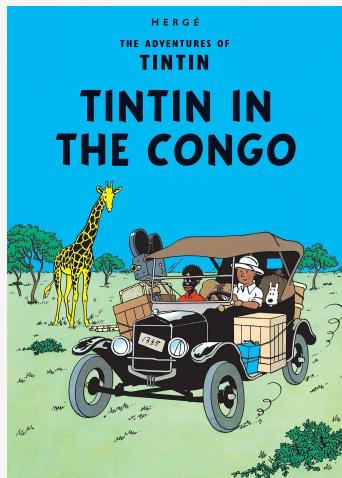
Dá-me um breve resumo da seguinte imagem, sem muitos detalhes. A descrição deve ser geral e ter um máximo de 10 palavras. Para isso, identifica o objeto ou a pessoa mais relevante na imagem, a principal relação entre os objetos destacados, a atividade mais importante representada, os atributos mais marcantes do objeto ou da pessoa principal e o contexto em que ocorre a cena. Depois, sintetiza tudo numa única frase descritiva e concisa, sem incluir texto adicional.

[Russian]

Дайте мне краткое описание следующего изображения без лишних деталей. Описание должно быть общим и содержать не более 10 слов. Для этого определите наиболее важный объект или человека на изображении, основную связь между выделенными объектами, главное изображенное действие, самые заметные характеристики главного объекта или человека и контекст, в котором происходит сцена. Затем объедините все это в одно описательное и краткое предложение, не добавляя лишнего текста.

Figure 5: Prompts used to evaluate via LLM-as-a-judge the language consistency of the caption provided by the model.

Story Generation



[English]

User: Create a short story using the following image as source.

VLM: Once upon a time, in the heart of the African savannah, there was a young boy named Tintin. Tintin was an adventurous soul, always eager to explore new places and meet new people. One day, he heard about a mysterious place called the Congo, a land of exotic animals and uncharted territories. Tintin decided to embark on an adventure to discover the secrets of the Congo. As he traveled through the dense jungle, he encountered a variety of animals, including a majestic giraffe and a wise elephant. He also met a group of explorers who were searching for a hidden treasure. Tintin joined their team and together they set out on a thrilling journey through the jungle.

They faced many challenges and obstacles, but Tintin's curiosity and bravery never faltered. Finally, they reached the heart of the Congo and discovered a treasure—a collection of rare artifacts from different cultures, each with its own story to tell.

Tintin and his friends returned home with a newfound appreciation for the wonders of the world and the importance of exploration and discovery. And so, the legend of Tintin in the Congo began.

[Spanish]

User: Crea una historia corta usando la siguiente imagen como fuente.

VLM: Había una vez un aventurero llamado Tintín que viajaba por el Congo en un coche antiguo. Lo acompañaban un guía y un perro, ambos llamados Milú. Tintín y sus amigos se encontraron con un grupo de gorilas y decidieron ayudarlos a escapar de los cazadores furtivos. Juntos, salvaron a los gorilas y regresaron a la ciudad, convirtiéndose en héroes en el Congo y continuando sus aventuras por el mundo.

[English] VLM: Once upon a time there was an adventurer named Tintin who traveled through the Congo in an antique car. He was accompanied by a guide and a dog, both named Snowy. Tintin and his friends came across a group of gorillas and decided to help them escape from poachers. Together, they saved the gorillas and returned to the city, becoming heroes in the Congo and continuing their adventures around the world.

(Continues on next page...)

[Catalan]

User: Crea una història curta utilitzant la següent imatge com a font.

VLM: En Tintín i el seu gos Milú es troben en un safari a l'Àfrica. El cotxe d'en Tintín es treu i han de demanar ajuda a un grup de nadius. Aquests reparen el cotxe i, a canvi, en Tintín i Milú els ensenyen una pel·lícula que és tot un èxit. Els nadius, encantats, els demanen que tornin l'any següent per mostrar-ne una altra. Finalment, en Tintín i Milú decideixen quedar-se a l'Àfrica i muntar un cinema per als habitants, començant així una nova vida plena d'amistat i intercanvi cultural.

[English] VLM: Tintin and his dog Snowy are on a safari in Africa. Tintin's car breaks down and they have to ask a group of natives for help. They repair the car and, in return, Tintin and Snowy show them a film that is a huge success. The natives, delighted, ask them to come back the following year to show another one. Finally, Tintin and Snowy decide to stay in Africa and set up a cinema for the inhabitants, thus beginning a new life full of friendship and cultural exchange.

[French]

User: Créez une histoire courte en utilisant l'image suivante comme source.

VLM: Il était une fois, dans la savane africaine, un petit garçon nommé Tintin, passionné d'aventures et rêvant d'explorer le continent. Un jour, il partit à la recherche de son idole, le célèbre explorateur Tintin, et se lança dans une aventure à travers la savane, rencontrant girafes, lions et éléphants. Finalement, il rencontra son idole dans le Congo, et ensemble, ils découvrirent de nombreux secrets et trésors, réalisant ainsi le rêve d'explorer l'Afrique.

[English] VLM: Once upon a time, in the African savannah, a little boy named Tintin had a passion for adventure and dreamed of exploring the continent. One day, he set off in search of his idol, the famous explorer Tintin, and embarked on an adventure across the savannah, encountering giraffes, lions and elephants. Finally, he met his idol in the Congo, and together they discovered many secrets and treasures, fulfilling the dream of exploring Africa.

[German]

User: Erstelle eine kurze Geschichte, indem du das folgende Bild als Quelle verwendest.

VLM: Tintin und sein Hund Milú begeben sich auf ein Abenteuer im Kongo. Während sie mit einem alten Auto durch die Savanne fahren, begegnen sie einem riesigen Giraffen, das sie fasziniert. Nachdem sie weitere Tiere wie Elefanten und Löwen getroffen haben, erleben sie eine spannende Reise, bei der sie schließlich einen Schatz entdecken und stolz auf ihre erlebten Abenteuer zurückblicken.

[English] VLM: Tintin and his dog Milú go on an adventure in the Congo. While driving through the savannah in an old car, they come across a giant giraffe that fascinates them. After meeting other animals such as elephants and lions, they experience an exciting journey during which they finally discover a treasure and look back proudly on their adventures.

[Italian]

User: Crea una breve storia utilizzando la seguente immagine come fonte.

VLM: Tintin e il suo cane Milú sono in un'avventura in Africa. Durante il viaggio, mentre guidano il loro veicolo, notano una grande giraffa che si avvicina e decidono di fermarsi per osservarla da vicino, vivendo un momento unico nel cuore della savana.

[English] VLM: Tintin and his dog Milú are on an adventure in Africa. On the way, while driving their vehicle, they notice a large giraffe approaching and decide to stop and observe it closely, experiencing a unique moment in the heart of the savannah.

(Continues on next page...)

[Dutch]

User: Maak een kort verhaal met behulp van de volgende afbeelding als bron.

VLM: Tintin in de Congo is een stripverhaal van Hergé waarin Tintin als journalist naar het Congo reist, vergezeld door zijn beste vriend Milú, een geredde chimpansee. Ze beleven tal van avonturen in de jungle, ontdekken de schoonheid van de natuur en maken kennis met de lokale bevolking, wat leidt tot een onvergetelijke reis vol humor en spanning.

[English] VLM: Tintin in the Congo is a comic strip by Hergé in which Tintin travels to the Congo as a journalist, accompanied by his best friend Milú, a rescued chimpanzee. They have numerous adventures in the jungle, discover the beauty of nature and meet the local people, leading to an unforgettable journey full of humor and excitement.

[Portuguese]

User: Crie uma curta história usando a seguinte imagem como fonte.

VLM: Era uma vez um explorador chamado Tintin, que se encontrava numa aventura no Congo. Ao volante de um carro vintage, acompanhado pelo seu amigo e animal de estimação, Milú, Tintin estava em busca de um tesouro escondido. No caminho, tiveram de improvisar para contornar um enorme elefante que bloqueava a passagem, mas conseguiram seguir em frente e, no final, encontraram o tesouro que buscavam.

[English] VLM: Once upon a time there was an explorer called Tintin, who was on an adventure in the Congo. At the wheel of a vintage car, accompanied by his friend and pet, Milú, Tintin was looking for hidden treasure. I don't walk, you had to improvise to get around a huge elephant that was blocking the passage, but we managed to keep going and, finally, we found the treasure we were looking for.

[Galician]

User: Crea unha historia curta utilizando a seguinte imaxe como fonte.

VLM: Había unha vez un explorador chamado Tintin no Congo. Durante a súa viaxe, atopouse con animais exóticos como un gorila, un elefante e un león, pero o que más chamou a súa atención foi un gran xaguar no medio do bosque. Cansado, o xaguar foi axudado por Tintin a chegar a un río para beber, e xuntos descubriron unha cova misteriosa con un tesouro oculto que decidiron compartir.

[English] VLM: Once upon a time there was an explorer named Tintin in the Congo. During his journey, he encountered exotic animals such as a gorilla, an elephant and a lion, but what caught his attention the most was a large jaguar in the middle of the forest. Tired, the jaguar was helped by Tintin to reach a river to drink, and together they discovered a mysterious cave with a hidden treasure that they decided to share.

[Polish]

User: Stwórz krótką historię, używając poniższego obrazu jako źródła.

VLM: Tytuł: "Tintin w Kong" Autor: Hergé Data wydania: 1931

[English] VLM: Title: "Tintin in the Congo" Author: Hergé Publication date: 1931

[Swedish]

User: Skapa en kort berättelse med hjälp av följande bild som källa.

VLM: Tintin och Milú är på en safari i Afrika. De sitter i en old-fashioned bil och ser på en giraffe som är i vägen. Tintin och Milú tar sig över till den andra sidan av vägen och ser på en elefant som är i vägen. De tar sig över till den andra sidan av vägen och ser på en elefant som är i vägen.

[English] VLM: Tintin and Milú are on a safari in Africa. They sit in an old-fashioned car and look at a giraffe that is in the road. Tintin and Milú cross to the other side of the road and look at an elephant in the road. They cross to the other side of the road and look at an elephant in the road.

(Continues on next page...)

[Russian]

User: Создайте короткий рассказ, используя следующее изображение в качестве источника.

VLM: Однажды, в далкой Африке, в самом сердце Конго, в 1920-х годах, произошло уникальное приключение, которое стало известно как "Тинтин в Конго". Тинтин, знаменитый репортер и путешественник, был отправлен в Африку, чтобы рассказать о жизни и культуре этого загадочного континента. Он и его верные спутники, Микель и Нельсон, сели в свой старенький, но надежный автомобиль и отправились в путешествие. Они пробирались через дикие джунгли, где их окружали дикие животные и таинственные племена. Тинтин и его команда были не только reporters, но и смелые путешественники, которые хотели узнать больше о мире, который был для них загадкой. Они встретили множество людей, которые были для них новыми и интересными, и узнали о их культуре и традициях. Они также узнали о том, что Конго был не только домом для диких животных, но и для людей, которые жили в гармонии с природой. В конце концов, Тинтин и его команда достигли своей цели и рассказали о своем путешествии в Европе, где их история вызывала большой интерес. "Тинтин в Конго" стал одним из самых известных и любимых приключений Тинтина, и он всегда будет помнить об этом

[English] VLM: Once upon a time, in faraway Africa, in the heart of the Congo, in the 1920s, a unique adventure took place that became known as “Tintin in the Congo.” Tintin, a famous reporter and traveler, was sent to Africa to report on the life and culture of this mysterious continent. He and his faithful companions, Mikel and Nelson, got into their old but reliable car and set off on their journey. They made their way through a wild jungle where they were surrounded by wild animals and mysterious tribes. Tintin and his team were not only reporters, but also brave travelers who wanted to learn more about a world that was a mystery to them. They met many people who were new and interesting to them and learned about their culture and traditions. They also learned that the Congo was not only home to wild animals, but also to people who lived in harmony with nature. In the end, Tintin and his team reached their goal and told their journey in Europe, where their story generated a lot of interest. “Tintin in the Congo” became one of Tintin’s most famous and beloved adventures, and he will always be remembered for it

[Norwegian]

User: Lag en kort historie ved å bruke følgende bilde som kilde.

VLM: Tintin og Milú er på en safari i Afrika. De er i en bil og ser på en giraffe.

[English] VLM: Tintin and Milú are on a safari in Africa. They are in a car looking at a giraffe.

[Hungarian]

User: Írj egy rövid történetet a következő kép felhasználásával forrásként.

VLM: Tintin a Kongóba indul, hogy megmentse a királyt és a királynőt, akiket a gorilla rabolt el. A királynő és a király a gorilla foglyul ejtette, és a Kongóba menekült. Tintin és a barátai elindulnak, hogy megmentsék a királyt és a királynőt.

[English] VLM: Tintin goes to the Congo to rescue the king and queen, who have been kidnapped by the gorilla. The king and queen have been captured by the gorilla and have fled to the Congo. Tintin and his friends set out to rescue the king and queen.

Figure 6: Multilingual generation examples with Text Regularization and merged with the original backbone LLM.

Image Description



[English]

User: Briefly describe the following image, without many details and in a few words.

VLM: A person's hand holding a black and yellow striped salamander.

[Spanish]

User: Describe brevemente la siguiente imagen, sin muchos detalles y en pocas palabras.

VLM: Una mano sosteniendo a un insecto, específicamente un insecto de color negro y amarillo.

[English] VLM: A hand holding an insect, specifically a black and yellow insect.

[Catalan]

User: Descriu breument la següent imatge, sense molts detalls i en poques paraules.

VLM: A la mà d'una persona hi ha una salamandra de color negre i groc.

[English] VLM: In a person's hand there is a black and yellow salamander.

[French]

User: Décrivez brièvement l'image suivante, sans trop de détails et en quelques mots.

VLM: Un petit animal noir et blanc est dans la main d'une personne.

[English] VLM: A small black and white animal is in a person's hand.

[German]

User: Beschreiben Sie das folgende Bild kurz, ohne viele Details und in wenigen Worten.

VLM: In der Hand einer Person liegt ein Frillengecko. Der Gecko hat eine schwarze und gelbe markierte Haut. Die Person hat einen Arm mit einem Armband.

[English] VLM: A frilled gecko lies in the hand of a person. The gecko has black and yellow marked skin. The person has an arm with a bracelet.

[Italian]

User: Descrivi brevemente la seguente immagine, senza molti dettagli e in poche parole.

VLM: Un insetto in una mano.

[English] VLM: An insect in a hand.

[Dutch]

User: Beschrijf kort de volgende afbeelding, zonder veel details en in een paar woorden.

VLM: De hand van een mens met een salamander in het palm.

[English] VLM: The hand of a human with a salamander in the palm.

(Continues on next page...)

[Portuguese]

User: Descreva resumidamente a imagem seguinte, sem muitos pormenores e em poucas palavras.

VLM: Uma pessoa está a segurar um sapo em sua mão. O sapo é de cor preta e amarela. O sapo está vivo.

[English] VLM: A person is holding a frog in his hand. The frog is black and yellow in color. The frog is alive.

[Galician]

User: Describe brevemente a seguinte imaxe, sen moitos detalles e en poucas palabras.

VLM: A man sostén unha larva de salamandra. A larva é de cor negra e amarela. A man que sostén a larva é branca.

[English] VLM: A man holds a yellow and black lobster nail.

[Russian]

User: Кратко опишите следующее изображение, без подробностей и в нескольких словах.

VLM: В руке человека, который держит змею.

[English] VLM: In the hand of a man holding a snake.

[Polish]

User: Opisz krótko poniższy obraz, bez wielu szczegółów i w kilku słowach.

VLM: Oto krótki opis tego obrazu: "Czarna i żółta salamandra w dłoni." To jest fotografia salamandry, która jest trzymana w ręce. Salamandra ma czarne i żółte plamy na swym ciele.

[English] VLM: Here is a short description of this image: 'A black and yellow salamander in hand.' This is a photograph of a salamander that is being held in hand. The salamander has black and yellow spots on its body.

[Swedish]

User: Beskriv kort nedanstående bild, utan många detaljer och med några få ord.

VLM: En hand med en salamander i den.

[English] VLM: A hand with a salamander in it.

[Norwegian]

User: Beskriv kort det følgende bildet, uten mange detaljer og med noen få ord.

VLM: En hand holding en salamander.

[English] VLM: A hand holding a salamander.

[Hungarian]

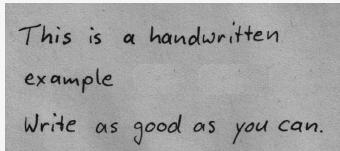
User: Röviden írja le a következő képet, sok részlet nélkül, néhány szóban.

VLM: A képen egy ember tartja a kezében egy fekete és sárga ragadozót.

[English] VLM: In the picture, a person holds in their hand a black and yellow predator.

Figure 7: Multilingual generation examples with the model trained with Text Regularization and merged with the original backbone LLM. Code switched words are underlined.

OCR & Translation



This is a handwritten example
example
Write as good as you can.

[English]

User: Read the text that appears in the image and translate it into English.

VLM: This is a handwritten example. Write as good as you can.

[Spanish]

User: Read the text that appears in the image and translate it into Spanish.

VLM: Este es un ejemplo de escritura a mano. Escribe tan bien como puedas.

[\[English\]](#) VLM: This is a handwritten example. Write as good as you can.

[Catalan]

User: Read the text that appears in the image and translate it into Catalan.

VLM: Aquest és un exemple de text escrit a mà. Escriu tan bé com puguis.

[\[English\]](#) VLM: This is a handwritten example. Write as good as you can.

[French]

User: Read the text that appears in the image and translate it into French.

VLM: C'est un exemple de lettre écrite à la main. Ecrivez comme vous pouvez.

[\[English\]](#) VLM: This is an example of a handwritten letter. Write as you can.

[German]

User: Read the text that appears in the image and translate it into German.

VLM: Dies ist ein Beispiel für einen handgeschriebenen Text. Schreiben Sie so gut wie Sie können.

[\[English\]](#) VLM: This is a handwritten example. Write as good as you can.

[Italian]

User: Read the text that appears in the image and translate it into Italian.

VLM: Questo è un esempio di scrittura a mano. Scrivete bene come voi potete.

[\[English\]](#) VLM: This is a handwritten example. Write as good as you can.

[Dutch]

User: Read the text that appears in the image and translate it into Dutch.

VLM: Dit is een handgeschreven voorbeeld. Schrijf zo goed als je kunt.

[\[English\]](#) VLM: This is a handwritten example. Write as good as you can.

[Portuguese]

User: Read the text that appears in the image and translate it into Portuguese.

VLM: Este é um exemplo de um texto escrito à mão. Escreva como bom como você pode.

[\[English\]](#) VLM: This is a handwritten example. Write as good as you can.

[Galician]

User: Read the text that appears in the image and translate it into Galician.

VLM: Isto é un exemplo de escrita a man. Escribe ben como se pode.

[\[English\]](#) VLM: This is a handwritten example. Write as good as you can.

(Continues on next page...)

[Russian]

User: Read the text that appears in the image and translate it into Russian.

VLM: Это пример написанного от руки текста. Напишите как можно лучше, чем вы можете.

[English] *VLM:* This is an example of a handwritten text. Write as best as you can.

[Polish]

User: Read the text that appears in the image and translate it into Polish.

VLM: Oto tekst w języku polskim: “To jest przykładowy handwriting example. Napisz tak, jak potrafisz”

[English] *VLM:* Here is the text in Polish: “This is a handwriting example. Write as much as you can”

[Swedish]

User: Read the text that appears in the image and translate it into Swedish.

VLM: This is a handwritten example. Write as good as you can.

[English] *VLM:* This is a handwritten example. Write as good as you can.

[Norwegian]

User: Read the text that appears in the image and translate it into Norwegian.

VLM: This is a handwritten example. Write as good as you can.

[English] *VLM:* This is a handwritten example. Write as good as you can.

[Hungarian]

User: Read the text that appears in the image and translate it into Hungarian.

VLM: Ez egy kézzel írott példa. Írjál úgy, mint te tudsz.

[English] *VLM:* This is a handwritten example. Write the way you know how.

Figure 8: Multilingual generation examples with the model trained with Text Regularization and merged with the original backbone LLM. Words in English whenever it is not the target are underlined.

Multi-Image Reasoning



[English]

User: Which dish from the ones shown in these images should I eat if I am lactose intolerant?

VLM: If you are lactose intolerant, you should eat the dish on the left, which is a stir fry of chicken and vegetables. This dish does not contain any dairy products, making it a good option for those who are lactose intolerant.

[Spanish]

User: ¿Qué plato de los que aparecen en estas imágenes debería de comer si soy intolerante a la lactosa?

VLM: Si eres intolerante a la lactosa, no deberías comer el plato de pizza.

[English] VLM: If you are lactose intolerant, you should not eat the pizza dish.

[Catalan]

User: Quin plat de les imatges donades hauria de menjar si sóc intolerant a la lactosa?

VLM: Si ets intolerant a la lactosa, hauries de menjar el plat de la primera imatge, ja que no conté productes lactis.

[English] VLM: If you are lactose intolerant, you should eat the dish in the first image, as it does not contain dairy products.

[French]

User: Quel plat parmi ceux qui apparaissent dans ces images devrais-je manger si je suis intolérant au lactose ?

VLM: Si vous êtes intolérant au lactose, je vous conseillerais de manger le plat de poulet grillé et de légumes.

[English] VLM: If you're lactose intolerant, I'd recommend eating the grilled chicken and vegetable dish.

[Italian]

User: Quale piatto tra quelli mostrati in queste immagini dovrei mangiare se sono intollerante al lattosio?

VLM: Se sei intollerante al lattosio, dovresti mangiare il primo piatto, il stir fry di pollo e verdura.

[English] VLM: If you are lactose intolerant, you should eat the first course, chicken and vegetable stir fry.

(Continues on next page...)

[German]

User: Welches Gericht aus den in diesen Bildern gezeigten sollte ich essen, wenn ich laktoseintolerant bin?

VLM: Wenn Sie laktoseintolerant sind, sollten Sie das erste Bild, das Chicken and Vegetables, essen.

[English] VLM: If you are lactose intolerant, you should eat the first picture, the Chicken and Vegetables.

[Dutch]

User: Welk gerecht van de getoonde gerechten in deze afbeeldingen zou ik moeten eten als ik lactose-intolerant ben?

VLM: Als je lactose-intolerant bent, zou je het gerecht in de eerste afbeelding moeten eten. Het is een salade met een dressing die geen lactose bevat. De salade is gemaakt van grilled chicken, rode en groeneppers, en is niet gebakken of gefrituurd, wat betekent dat het geen lactose bevat. De salade is dus een lactose-vrije optie.

[English] VLM: If you are lactose intolerant, you should eat the dish in the first picture. It is a salad with a dressing that does not contain lactose. The salad is made with grilled chicken, red and green peppers, and is not baked or fried, which means that it does not contain lactose. The salad is therefore a lactose-free option.

[Portuguese]

User: Qual prato entre os mostrados nestas imagens eu deveria comer se for intolerante à lactose?

VLM: Se você é intolerante à lactose, o prato que você deve evitar é o pizza. A pizza é feita com queijo, que é uma fonte de lactose, e pode causar sintomas como dor abdominal, diarreia, náuseas e vômitos. Se você tem intolerância à lactose, é importante evitar alimentos que contenham lactose para evitar esses sintomas.

[English] VLM: If you are lactose intolerant, the food you should avoid is pizza. Pizza is made with cheese, which is a source of lactose, and can cause symptoms such as abdominal pain, diarrhea, nausea, and vomiting. If you are lactose intolerant, it is important to avoid foods that contain lactose to avoid these symptoms.

[Galician]

User: Que prato das imaxes dadas debería comer se son intolerante á lactosa?

VLM: Se é intolerante á lactosa, debería comer o prato da primeira imaxe, que é unha ensalada de pollo e verduras.

[English] VLM: If you are lactose intolerant, you should eat the dish in the first picture, which is a chicken and salad.

[Polish]

User: Jakie danie z pokazanych na tych obrazkach powiniennem zjeść, jeśli jestem nietolerancyjny na laktozę?

VLM: Jeżeli jesteś nietolerancki na laktozę, to nie powinieneś jeść pizzy. Pizza jest gotowa z mozzarellą, a mozzarella jest produkowana z mleka. Jeżeli wolisz jeść potrawy bez laktozy, to możesz zjeść kurczaka z warzywami. Kurczak jest bezlaktozny, a warzywa są bezlaktozne.

[English] VLM: If you are lactose intolerant, you should not eat pizza. Pizza is made with mozzarella, and mozzarella is made from milk. If you prefer to eat lactose-free food, you can eat chicken with vegetables. Chicken is lactose-free, and vegetables are lactose-free.

(Continues on next page...)

[Russian]

User: Какое блюдо из показанных на этих изображениях мне следует съесть, если я непереносим к лактозе?

VLM: Если вы не переносите лактозу, то вам следует выбрать блюдо, которое не содержит молочных продуктов. На изображении 1 показаны кусочки курицы с овощами, которые не содержат лактозы. На изображении 2 показана пицца с сыром и другими toppings, которые могут содержать лактозу. Однако, если вы хотите попробовать что-то, что не содержит лактозы, то вам следует обратиться к другим блюдам, которые не содержат молочные продукты.

[English] VLM: If you are lactose intolerant, you should choose a dish that does not contain dairy products. Picture 1 shows chicken pieces with vegetables, which do not contain lactose. Picture 2 shows a pizza with cheese and other toppings, which may contain lactose. However, if you want to try something that does not contain lactose, you should turn to other dishes that do not contain dairy products.

[Swedish]

User: Vilken rätt av de som visas i dessa bilder bör jag äta om jag är laktosintolerant?

VLM: Om du är laktosintolerant, skulle jag rekommendera att du äter den första bilden, som är en stir-fry med kyckling, rödbetor, och peppar. Det har ingen laktos i sig och passar bra för dem som är laktosintoleranta.

[English] VLM: If you are lactose intolerant, I would recommend eating the first picture, which is a stir-fry with chicken, beetroot, and pepper. It has no lactose in it and is suitable for those who are lactose intolerant.

[Norwegian]

User: Hvilket måltid av de som vises på disse bildene bør jeg spise hvis jeg er laktoseintolerant?

VLM: For en person som er laktoseintolerant, er det best å spise en av de to alternativene i disse bildene:

1. **Chicken and Vegetable Stir-Fry:** - Chicken er en bra kilde til protein, men det er ikke laktose.
- Vegetable Stir-Fry er full av fiber, vitaminer, og mineraler, og har ingen laktose.
2. **Pepperoni and Sausage Pizza:** - Pepperoni og sausage er ikke laktoseintolerant, men de har ikke så store mengder laktose som fullpizza.

Det er viktig å merke seg at disse alternativene ikke er full substitut for en fullpizza, men de er bedre enn fullpizza for en laktoseintolerant person.

[English] VLM: or a person who is lactose intolerant, it is best to eat one of the two options in these pictures:

1. **Chicken and Vegetable Stir-Fry:** - Chicken is a good source of protein, but it's not lactose. - Vegetable Stir-Fry is full of fiber, vitamins, and minerals, and has no lactose.
2. **Pepperoni and Sausage Pizza:** - Pepperoni and Sausage are not lactose intolerant, but they don't have as high amounts of lactose as full pizzas.

It's important to note that these options are not a full substitute for a full pizza, but they are better than full pizza for a lactose intolerant person.

[Hungarian]

User: Melyik ételt kellene enned a képeken látható ételek közül, ha laktózérzékeny vagyok?

VLM: Ha laktózérzékeny vagy, akkor az ételt, amely nem tartalmaz laktózt, kellene enned. A laktózmentes alternatíváként a csirkét és a zöldségeket ajánlanám. A csirkét és a zöldségeket laktózmentesen is elkészítheted, és nem tartalmaznak laktózt.

[English] VLM: If you are lactose intolerant, you should eat the food that does not contain lactose. As a lactose-free alternative, I would recommend chicken and vegetables. You can also make chicken and vegetables lactose-free and they do not contain lactose. You can combine chicken and vegetables to get a delicious lactose-free meal.

Figure 9: Multilingual generation examples with the model trained with Text Regularization and merged with the original backbone LLM.

Generation Codeswitching Examples

Target Language: Spanish

Dos tacos con salsa y cilantro sobre una tortilla, **served in a paper wrapper and a lime wedge, with a white plate with a blue floral pattern underneath.**



Target Language: Dutch

Het is een scenische boulevard met hotels aan de waterkant, een restaurant met **outdoor seating**, en een **sailboat** aan de promenade..



Target Language: Russian

«Postnord» — это название компании, которая занимается логистикой и дистрибуцией товаров. **Company logo features bold, blue text on a wooden background, which gives a rustic and modern feel.**



Target Language: German

In einem urbanen **Setting** steht ein orange VW Golf auf einem **cobblestone**. Die license plate reads **HH 7293**; **and the car has a small antenna on the roof. A person stands behind the car, with a backpack on their back.**



Figure 10: Code switching examples observed when generating the captions of the Crossmodal-3600 dataset images using the Tr-3S-0.05x model. These examples showcase the need of adding a complementary evaluation apart from GlotLid's language detection in order to check language consistency. The words in English are shown in bold letters.

Visual Data			
Dataset	Field	Stage	Citation
LLaVA Pretrain LCS-558K	Image Captions	1	Liu et al. (2023b)
BLIP558K	Detailed Description	1.5	Liu et al. (2024b)
CC3M	Detailed Description	1.5	Liu et al. (2024b)
COCO118K	Detailed Description	1.5	Liu et al. (2024b)
Evol Instruct	Math/Reasoning	1.5	Chen et al. (2024a)
UReader	OCR	1.5	Ye et al. (2023)
SynthDOG	Language	1.5	Kim et al. (2022)
AI2D	Infographics	2/2.5	Kembhavi et al. (2016)
Allava Instruct	General	2/2.5	Chen et al. (2024a)
AOKVQA	General	2/2.5	Schwenk et al. (2022)
Cambrian (filtered)	General	2/2.5	Tong et al. (2024)
Chart2Text	Infographics	2/2.5	Obeid and Hoque (2020)
ChartQA	Infographics	2/2.5	Masry et al. (2022)
ChromeWriting	OCR	2/2.5	-
CLEVR	General	2/2.5	Johnson et al. (2017)
CLEVR-Math	Math/Reasoning	2/2.5	Johnson et al. (2017)
COCO Caption	General	2/2.5	Lin et al. (2014)
Diagram Image2Text	Infographics	2/2.5	-
DocVQA	Infographics	2/2.5	Mathew et al. (2021)
DVQA	Infographics	2/2.5	Kafle et al. (2018)
FigureQA	Infographics	2/2.5	Kahou et al. (2017)
GQA	Math/Reasoning	2/2.5	Hudson and Manning (2019)
Geo170K Align	Math/Reasoning	2/2.5	Gao et al. (2023)
Geo170K QA	Math/Reasoning	2/2.5	Gao et al. (2023)
Geo3K	Math/Reasoning	2/2.5	-
Geometry3K	Math/Reasoning	2/2.5	Lu et al. (2021a)
GeoMVerse	Math/Reasoning	2/2.5	Kazemi et al. (2024)
GeoQA+	Math/Reasoning	2/2.5	Chen et al. (2021)
GEOS	Math/Reasoning	2/2.5	Seo et al. (2015)
Hateful Memes	General	2/2.5	Kiela et al. (2020)
HiTab	Infographics	2/2.5	Cheng et al. (2022)
HME100K	OCR	2/2.5	Yuan et al. (2022)
IAM	OCR	2/2.5	Marti and Bunke (2002)
IconQA	General	2/2.5	Lu et al. (2021b)
IIIT5K	OCR	2/2.5	Mishra et al. (2012)
Infographic VQA	Infographics	2/2.5	Mathew et al. (2022)
InterGPS	General	2/2.5	Lu et al. (2021a)
Image Textualization	General	2/2.5	Pi et al. (2024)
K12 Printing	OCR	2/2.5	-
LLaVA-158K	General	2/2.5	Liu et al. (2023b)
LLaVA-Wild (train)	General	2/2.5	Liu et al. (2023b)
LLaVAR	General	2/2.5	Zhang et al. (2023b)
LRV-Chart	Infographics	2/2.5	Liu et al. (2023a)
LRV-Normal	Math/Reasoning	2/2.5	Liu et al. (2023a)
Magpie Pro	Language	2/2.5	Xu et al. (2024a)

(Continues on next page...)

Dataset	Field	Stage	Citation
MapQA	Math/Reasoning	2/2.5	Chang et al. (2022a)
MathQA	Math/Reasoning	2/2.5	Amini et al. (2019)
MAVIS	Math/Reasoning	2/2.5	Zhang et al. (2024)
OKVQA	General	2/2.5	Marino et al. (2019)
OCR-VQA	OCR	2/2.5	Mishra et al. (2019)
RAVEN	Math/Reasoning	2/2.5	Zhang et al. (2019)
RefCOCO	General	2/2.5	Yu et al. (2016)
Rendered Text	OCR	2/2.5	-
RoBUT	Infographics	2/2.5	Zhao et al. (2023)
ScienceQA	General	2/2.5	Lu et al. (2022)
Screen2Words	Infographics	2/2.5	Wang et al. (2021)
ShareGPT4O	General	2/2.5	Cui et al. (2024)
ShareGPT4V	General	2/2.5	Chen et al. (2025)
ST-VQA	General	2/2.5	Biten et al. (2019)
Super-CLEVR	Math/Reasoning	2/2.5	Li et al. (2023c)
TabMWP	Math/Reasoning	2/2.5	Lu et al. (2023)
TallyQA	General	2/2.5	Acharya et al. (2019)
TextCaps	OCR	2/2.5	Sidorov et al. (2020)
TextOCR-GPT4	OCR	2/2.5	Carter (2024)
TQA	Infographics	2/2.5	Kembhavi et al. (2017)
UniGeo	Math/Reasoning	2/2.5	Chen et al. (2022)
Ureader	Infographics	2/2.5	Ye et al. (2023)
Vision FLAN	General	2/2.5	Xu et al. (2024b)
Visual7W	General	2/2.5	Zhu et al. (2016)
Visual Genome	Math/Reasoning	2/2.5	Krishna et al. (2017)
VisText	General	2/2.5	Tang et al. (2023)
VisualMRC	Infographics	2/2.5	Tanaka et al. (2021)
VizWiz	General	2/2.5	Gurari et al. (2018)
VQARAD	General	2/2.5	Lau et al. (2018)
VQAv2	General	2/2.5	Antol et al. (2015)
VSR	General	2/2.5	Liu et al. (2023b)
WebSight	General	2/2.5	Laurençon et al. (2024)
Spot-the-Diff	Multi-Image	2.5	Jhamtani and Berg-Kirkpatrick (2018)
Birds-to-Words	Multi-Image	2.5	Forbes et al. (2019)
CLEVR-Change	Multi-Image	2.5	Park et al. (2019)
HQ-Edit-Diff	Multi-Image	2.5	Hui et al. (2024)
MagicBrush-Diff	Multi-Image	2.5	Zhang et al. (2023a)
IEdit	Multi-Image	2.5	Tan et al. (2019)
AESOP	Multi-Image	2.5	Ravi et al. (2021)
FlintstonesSV	Multi-Image	2.5	Gupta et al. (2018)
PororoSV	Multi-Image	2.5	Li et al. (2019)
VIST	Multi-Image	2.5	Huang et al. (2016)
WebQA	Multi-Image	2.5	Chang et al. (2022b)
TQA (MI)	Multi-Image	2.5	Kembhavi et al. (2017)
OCR-VQA (MI)	Multi-Image	2.5	Mishra et al. (2019)
DocVQA (MI)	Multi-Image	2.5	Mathew et al. (2021)

(Continues on next page...)

Dataset	Field	Stage	Citation
MIT-StateCoherence	Multi-Image	2.5	Isola et al. (2015)
MIT-PropertyCoherence	Multi-Image	2.5	Isola et al. (2015)
RecipeQA ImageCoherence	Multi-Image	2.5	Yagcioglu et al. (2018)
VISION	Multi-Image	2.5	Bai et al. (2023)
Multi-VQA	Multi-Image	2.5	Li et al. (2023a)
IconQA	Multi-Image	2.5	Lu et al. (2021b)
Co-Instruct	Multi-Image	2.5	Wu et al. (2024)
DreamSim	Multi-Image	2.5	Fu et al. (2023)
ImageCoDe	Multi-Image	2.5	Krojer et al. (2022)
nuScenes	Multi-Image	2.5	Caesar et al. (2020)
ScanQA	Multi-Image	2.5	Azuma et al. (2022)
ALFRED	Multi-Image	2.5	Shridhar et al. (2020)
ContrastCaption	Multi-Image	2.5	Jiang et al. (2024)
VizWiz (MI)	Multi-Image	2.5	Gurari et al. (2018)
ScanNet	Multi-Image	2.5	Dai et al. (2017)
COMICS Dialogue	Multi-Image	2.5	Iyyer et al. (2017)
NLVR2	Multi-Image	2.5	Suhr et al. (2019)
NExT-QA	Video	2.5	Xiao et al. (2021)
Ego-4D	Video	2.5	Grauman et al. (2024)
YouCook2	Video	2.5	Zhou et al. (2018)
ActivityNet	Video	2.5	Yu et al. (2019)
Charades	Video	2.5	Sigurdsson et al. (2016)
ShareGPT4Video	Video	2.5	Chen et al. (2024c)

Table 10: English only visual datasets used throughout this work. The same data as proposed in LLaVA-OneVision has been used.

Text-Only Data			
Dataset	Field	Stage	Citation
Aya Dataset	General	1.5/2/2.5	Singh et al. (2024)
CoqCat	Conversation QA	1.5/2/2.5	Gonzalez-Agirre et al. (2024)
Databricks Dolly 15k	General	1.5/2/2.5	Conover et al. (2023)
Databricks Dolly 3k CA	General	1.5/2/2.5	-
FLORES-200 (Instructions)	Translations	1.5/2/2.5	Costa-jussà et al. (2024)
MentorCA	General	1.5/2/2.5	-
No Robots	General	1.5/2/2.5	Rajani et al. (2023)
OASST	General	1.5/2/2.5	Köpf et al. (2023)
OASST-CA	General	1.5/2/2.5	-
RAG Multilingual	General	1.5/2/2.5	-
Tower-Blocks-v0.1	Text-Insight	1.5/2/2.5	Alves et al. (2024)

Table 11: Multilingual text-only datasets added throughout the visual instruction process.