

Multilingual evaluation of image captioning: How far can we get with CLIP models?

Anonymous ACL submission

Abstract

The evaluation of image captions, looking at both linguistic fluency and semantic correspondence to visual contents, has witnessed a significant effort. Still, despite advancements such as the CLIPScore metric, multilingual captioning evaluation has remained relatively unexplored. This work assesses the use of CLIPScore in multilingual captioning, evaluating different models in a variety of settings. To address the lack of multilingual test data, we consider two different strategies: (1) using machine-translated datasets with human judgements, and (2) re-purposing multilingual datasets that target inference and reasoning. Our results show that multilingual CLIP models can perform on par with their English-centric counterparts on English benchmarks while allowing for multilingual assessments. Performance increases with model finetuning and according to model size. Larger models, trained with more data, attained similar performance to more advanced methods that extended the original CLIPScore. Tests with machine-translated data show that multilingual CLIPScore can also maintain a high correlation with human judgements across different languages, and additional tests with natively multilingual and multicultural data further attest to the high-quality assessments.

1 Introduction

Computer-generated image captions are nowadays commonly used as descriptive annotations. The image captioning task has been extensively studied, including in multilingual settings, with many recent approaches combining established vision encoders with large language model decoders (Ramos et al., 2023b,a; Yang et al., 2023; Geigle et al., 2023; Ramos et al., 2024). The automatic evaluation of captions, accounting for linguistic and visual contents, has also witnessed a significant effort. Approaches such as CLIPScore (Hessel et al., 2021) have been proposed to evaluate captions through cosine similarity between image and text embeddings,

leveraging large-scale pre-trained vision and language models and achieving high correlations with human judgments. Still, despite the many recent advancements, most approaches are English-centric, while multilingual image captioning evaluation has remained relatively unexplored.

This work explores the use of CLIPScore in multilingual captioning, evaluating different CLIP models under various settings. Given the lack of available benchmarks for the evaluation of multilingual captioning metrics, we propose two different evaluation strategies: (1) using Machine Translation (MT) to obtain multilingual data from English-centric benchmarks, and (2) re-purposing multilingual benchmarks originally designed for the evaluation of the semantic inference and reasoning capabilities of vision-language models.

Through extensive experiments, we show that multilingual CLIP models achieve comparable or even better performance on English benchmarks while allowing for multilingual assessments. We also propose a multilingual fine-tuning strategy for CLIPScore, that allows to account for linguistic and cultural diversity while learning from human judgements, resulting in further performance improvements. Performance generally increases according to model size. Larger models, trained on more data, attained similar or even better performance to methods that extended the original CLIPScore (Sarto et al., 2023; Kim et al., 2022; Hu et al., 2023; Kim et al., 2023; Narins et al., 2024; Wada et al., 2024). Tests with machine-translated data show that multilingual CLIPScore can also maintain a high correlation with human judgements across different languages, and additional tests with natively multilingual and multicultural data further attest to high-quality of assessments across languages.

Our primary contributions include ① a comprehensive evaluation of existing models on English-centric and multilingual benchmarks, assessing correlation to human judgements; ② an extension of

043
044
045
046
047
048
049
050
051
052
053
054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083

English-centric benchmarks to multiple languages, incorporating human evaluations and diverse linguistic phenomena while preserving the original benchmarks' quality; ③ an adaptation of existing multilingual and multicultural datasets for captioning evaluation; ④ a finetuning strategy that accounts for linguistic and cultural diversity as well as alignment with human judgements, leading to multilingual CLIPScore models that outperform previous work across several benchmarks¹.

2 Related Work

Conventional image captioning evaluation has relied on reference-based assessments, where machine-generated captions are compared against human-generated ones (i.e., the references). Frequently used metrics such as BLEU or CIDEr (Vedantam et al., 2015) rely on lexical matches, and hence may fail to capture finer nuances and semantic overlaps in rich captions. A recent shift in the evaluation paradigm involves the use of learned vision-and-language models to enable evaluation through reference-free metrics.

The CLIPScore metric (Hessel et al., 2021) was one of the first proposals for evaluating image captions that departed from the traditional metrics. Grounded in a vision-and-language encoder, specifically the original Contrastive Language-Image Pre-training (CLIP) model (Radford et al., 2021), this strategy employs a modified cosine similarity between representations for the input image and the caption under evaluation. CLIPScore exhibits a high correlation with human judgments across various datasets, and despite being a reference-free metric, it even surpasses established reference-based metrics like BLEU and CIDEr. The authors also introduced a reference-augmented version named RefCLIPScore, which additionally uses the cosine similarity between candidate and reference captions to further improve the correlation with human assessments. CLIPScore and RefCLIPScore are currently the most widely used learned metrics for captioning evaluation. However, many studies in the area still only report results using traditional lexical-based metrics. Some previous studies have also proposed the combination of CLIPScore and CIDEr through a simple weighted average (Qiu et al., 2023), arguing that this can further boost the correlation with human assessments.

¹The code and adapted datasets supporting our evaluation experiments will be made publicly available upon acceptance.

PACScore (Sarto et al., 2023) extended CLIPScore by introducing a contrastive strategy that uses curated data to further finetune the CLIP projection layers for captioning evaluation. By augmenting with generated images and paraphrased texts, PACScore can achieve better correlations with human judgements across several datasets.

InfoMetIC also builds on CLIP, aiming to provide detailed and explainable feedback in the context of captioning evaluation (Hu et al., 2023), reporting incorrect words and unmentioned image regions at a fine-grained level, while also providing a text precision score, a vision recall score, and an overall quality score at a coarse-grained level. InfoMetIC was found to outperform CLIPScore even when the latter is finetuned on captioning data.

Mutual Information Divergence (MID) is instead a unified metric for multimodal generative models (Kim et al., 2022), targeting both text-to-image and image-to-text tasks. The metric quantifies the alignment between generated visual and textual features by employing the concept of Mutual Information (MI) from information theory.

Some recent studies have noted that metrics like CLIPScore can lack rating granularity, arguing the need for better benchmark datasets for assessing image captioning evaluation metrics (Ahmadi and Agrawal, 2023). Datasets like VICR (Narins et al., 2024) or Polaris (Wada et al., 2024) have recently been proposed, relying on more rigorous procedures for collecting human ratings, although still focusing only on the English language. Together with the proposal of the new datasets, the authors also showed that models trained specifically for image captioning evaluation can slightly outperform CLIPScore. We corroborate these findings, proposing a finetuned version of CLIPScore that outperforms other variants.

One previous study has specifically looked into multilingual image captioning evaluation (Kim et al., 2023), proposing a method based on finetuning the text encoder of CLIP with a language-agnostic method to distinguish the perturbed text from the original text. The authors have also developed a novel dataset to evaluate multilingual image captioning metrics, but unfortunately, this dataset is not yet publicly available.

3 Multilingual CLIPScore

The CLIPScore metric uses an adjusted cosine similarity to compare representations for the input im-

age and the caption being assessed, as originally described by Hessel et al. (2021). In our work, we adopted the original formulation. A more detailed explanation of the CLIPScore and RefCLIPScore metrics can be found in Appendix A.

To boost performance, we propose a strategy to finetune multilingual CLIP models in a setting that considers both linguistic and cultural diversity, while accounting for human preference alignment.

Two distinct datasets were used for finetuning. The first, CrossModal-3600, focuses on multilingual and multicultural imagery (Thapliyal et al., 2022), whereas the second, VICR, comprises English image-caption pairs that are evaluated by humans (Narins et al., 2024) and we machine translated to different languages following a strict translation scheme to help maintain quality. We also combined different losses, tailored to the specific characteristics of each dataset.

In more detail, to enhance the model’s ability to process multilingual and multicultural instances, we finetuned it on both datasets using the original CLIP contrastive loss, which can be formally described as follows:

$$\mathcal{L}_C = -\frac{1}{2N} \sum_{i=1}^N \left[\log \frac{e^{s_{i,i}/\tau}}{\sum_{j=1}^N e^{s_{i,j}/\tau}} + \log \frac{e^{s_{i,i}/\tau}}{\sum_{j=1}^N e^{s_{j,i}/\tau}} \right], \quad (1)$$

where N is the number of image-text pairs in a batch, $s_{i,j}$ is the similarity score between the i -th image and the j -th text description, and τ is a temperature parameter that scales the similarity scores and helps in controlling the concentration level of the distribution.

For the second dataset, to improve the alignment of CLIPScores with human ratings, we also consider a Pearson correlation loss:

$$\mathcal{L}_P = 1 - \frac{(x-\bar{x})^T (y-\bar{y})}{\| (x-\bar{x}) \| \| (y-\bar{y}) \|}, \quad (2)$$

where x is the vector of CLIPScores values, y is a vector with the human rating scores, and \bar{x} and \bar{y} are the respective average values.

Considering that both loss functions can benefit from larger batch sizes, we sample instances for training alternating between each task, without mixing instances from the different datasets in the same batch and applying the respective loss functions. We accumulate gradients for two steps before updating the network, effectively combining both loss effects while leveraging the benefits of larger batches, i.e., $\mathcal{L} \sim \mathcal{L}_C + \mathcal{L}_P$.

4 Experimental Evaluation

This section presents the datasets, the experimental setup, and the results for different CLIP models, considering English, multilingual, and multicultural scenarios for image captioning evaluation.

4.1 Datasets

For the assessment of human judgment correlations, experiments were conducted using the following well-established English-only datasets containing one or more human quality assessments for each image-caption pair.

- **Expert** (Flickr8K-Expert) contains 5,664 pairs (Hodosh et al., 2013).
- **Crowdfower** (Flickr8K-CF) contains 47,830 pairs (Hodosh et al., 2013).
- **Composite** contains 13,146 pairs (Aditya et al., 2015).
- **VICR** contains 10,175; 2,310; and 3,161 pairs respectively in training, validation and test splits (Narins et al., 2024).

To evaluate the robustness of our models to different linguistic phenomena, we used the VALSE (Parcalabescu et al., 2021) dataset, which contains 6,704 correct image-caption pairs, plus their respective foil caption versions.

Comparing multilingual versus monolingual models on English data provides a limited overview of model performance. However, high-quality multilingual resources with curated data featuring human assessments of caption quality are scarce (or even unavailable), hindering multilingual evaluation. To mitigate this limitation, we proposed a translation scheme leveraging large machine translation models (Fan et al., 2021; Alves et al., 2024; Liu et al., 2020), combined with language and translation quality estimation models (Rei et al., 2022), to automatically translate English captions for which we already have human assessments, preserving the highest translation quality possible. With high-quality translations, human judgments should be valid across the different target languages. This strategy is further detailed in Appendix C.

Our language selection is in line with recent machine translation studies (Alves et al., 2024), covering high-resource languages (i.e., English, French, German, Spanish, and Chinese) and also mid- (i.e., Portuguese, Italian, and Russian) to low-resource languages (i.e., Dutch and Korean). We machine-translated both the VICR and VALSE datasets into nine languages using this technique.

In addition, to further expand our multilingual evaluation, we used naturally multilingual and multicultural datasets, i.e., XVNLI (Bugliarello et al., 2022) and MaRVL (Liu et al., 2021), re-purposing them for the evaluation of image captioning metrics. The original datasets feature differences in terms of language composition, and we thus also expanded the XVNLI dataset by translating its data from English into the languages present in the MaRVL dataset. We chose not to further extend MaRVL to other languages because the dataset features image-caption pairs that focus on culturally specific concepts in association with the target languages.

4.2 Evaluation Metrics

We evaluate the different models using correlation with human judgements, and also with classification tasks. Regarding the correlation experiments, we measure performance using three different correlation coefficients, namely Spearman ρ and Kendall τ with variations b and c . The correlation metrics are formally defined in Appendix B.

For the multilingual/multicultural experiments, we measure performance under the assumption that a caption entailed by an image should reflect a higher CLIPScore than a contradiction/foil caption.

4.3 Experiments and Results

This section presents experimental results for the different models and evaluation datasets, establishing a comparison with previously reported results and contributing to the multi-linguistic exploration of existing models and datasets. We also performed a qualitative study focusing on image-caption pairs that feature concepts that could be associated with cultural bias, which is reported in Appendix F.

4.3.1 Correlation Assessment on English Data

Table 1 displays the correlation results between CLIPScore values and human ratings, across the four English datasets and considering existing and publicly available English-only and multilingual CLIP models. The results show that the CLIPScore estimates can improve significantly with larger CLIP models trained with more data.

Comparing results against the performance of the original CLIPScore computed with a ViT-B/32 vision encoder, we can observe substantial improvements. Apple’s ViT-H/14 model with 378px resolution achieves the highest correlation results among the English models, although the multilingual model of the same size from LAION is the

second best-performing alternative (even outperforming the English-only model on the Composite dataset). The largest multilingual CLIP model, trained on approximately 5 billion instances, also outperforms the similarly sized English-only model trained with approximately 2 billion instances, perhaps indicating that exposure to diverse language data can enable CLIPScore values to better correlate with human judgements. Overall, the findings support the argument that CLIPScore with a multilingual model can maintain, or even improve the performance, over an English-only model.

In turn, Table 2 compares recent studies proposing other metrics against the best performing English-only and multilingual CLIPScore models². We considered different evaluation settings, assessing results (a) without references, (b) using references, and (c) combining CIDEr with RefCLIPScore when using references (Qiu et al., 2023).

Results confirm that the original CLIPScore outperforms standard captioning metrics in all evaluated datasets, such as BLEU or CIDEr. The correlations consistently improve with RefCLIPScore, but the combination of RefCLIPScore with CIDEr only improved the smaller CLIP model (by previously published results from Qiu et al. (2023)), instead decreasing the performance in the cases involving the other CLIPScore variants.

The multilingual model competes head-to-head with the best English-only models, outperforming the previously published results for almost all the metrics, whether human references are considered or not. This competitive performance in both reference-free and reference-aided settings further solidifies the potential of multilingual CLIPScore.

Lastly, it is also interesting to note that our fine-tuned model outperformed the original multilingual LAION ViT-H model, and also Apple’s best model across all CLIPScore variants, on both the VICR and Expert datasets, even without using references. The version of our model that uses references achieved the best performance for these datasets, surpassing even specialized architectures such as VICR, InfoMetIC, and RefPACScore.

4.3.2 Correlation on Multilingual Data

We used the VICR multilingual variants to analyse the human human ratings, considering the best per-

²Note that we computed all the correlation scores for the CLIPScore variants and traditional image caption metrics based on lexical matches, whereas the results for other more recent metrics are taken from the corresponding publications.

	VICR			Expert			CrowdFlower			Composite		
	τ_b	τ_c	ρ	τ_b	τ_c	ρ	τ_b	τ_c	ρ	τ_b	τ_c	ρ
English												
openai/clip-vit-base-patch32	60.7	67.1	76.9	51.1	51.5	63.1	34.4	17.8	42.4	50.6	51.4	67.9
apple/DFN5B-CLIP-ViT-H-14-378	67.4	74.4	83.1	56.3	56.6	68.4	38.5	19.9	47.1	55.0	55.9	72.8
apple/DFN5B-CLIP-ViT-H-14	66.5	73.5	82.4	55.6	56.0	67.7	38.2	19.7	46.8	54.5	55.3	72.3
apple/DFN2B-CLIP-ViT-L-14	65.8	72.6	81.6	55.5	55.6	67.5	37.1	19.2	45.6	53.6	54.3	71.3
laion/CLIP-ViT-g-14-laion2B-s12B-b42K	66.4	73.0	82.0	55.3	55.1	66.9	37.1	19.1	45.5	54.9	55.7	72.7
laion/CLIP-ViT-H-14-laion2B-s32B-b79K	66.4	73.1	82.1	54.9	54.9	66.7	37.2	19.2	45.6	53.6	54.4	71.2
laion/CLIP-ViT-L-14-laion2B-s32B-b82K	65.9	72.6	81.7	54.4	54.5	66.3	36.7	18.9	45.1	53.6	54.4	71.3
BAAI/AltCLIP	65.1	71.9	81.1	54.1	54.4	66.2	36.2	18.7	44.6	53.8	54.6	71.6
openai/clip-vit-large-patch14-336	62.0	68.5	78.1	52.6	53.0	64.6	35.4	18.3	43.7	52.8	53.6	70.3
openai/clip-vit-large-patch14	62.3	68.9	78.5	52.6	53.0	64.6	35.2	18.2	43.4	52.5	53.3	70.0
Multilingual												
laion/CLIP-ViT-H-14-frozen-xml-roberta-large-laion5B-s13B-b90k	67.6	73.0	82.4	57.4	54.3	67.3	38.2	19.4	46.2	55.5	56.2	73.2
BAAI/AltCLIP-m18	66.6	73.4	82.3	54.8	55.0	66.8	37.1	19.2	45.6	55.2	56.0	73.0
M-CLIP/XLM-Roberta-Large-Vit-B-16Plus	64.6	71.3	80.6	54.6	54.8	66.7	36.0	18.6	44.3	52.1	52.9	69.6
BAAI/AltCLIP-m9	64.2	70.9	80.3	54.1	54.4	66.2	36.4	18.8	44.8	53.9	54.6	71.6
laion/CLIP-ViT-B-32-xml-roberta-base-laion5B-s13B-b90k	63.3	69.8	79.3	52.7	52.8	64.5	35.2	18.2	43.3	49.9	50.6	67.1
M-CLIP/XLM-Roberta-Large-Vit-L-14	62.2	68.7	78.4	53.0	53.4	65.0	35.4	18.3	43.7	52.9	53.6	70.4
M-CLIP/XLM-Roberta-Large-Vit-B-32	60.5	66.9	76.7	51.8	52.2	63.9	34.4	17.8	42.4	50.7	51.4	67.9
sentence-transformers/clip-ViT-B-32-multilingual-v1	60.3	66.7	76.6	51.5	51.8	63.6	33.3	17.2	41.1	48.9	49.6	65.9

Table 1: Correlation between CLIPScore values and human rankings, considering a set of different English (top rows) and multilingual (bottom rows) CLIP models.

	VICR			Expert			CrowdFlower			Composite		
	τ_b	τ_c	ρ	τ_b	τ_c	ρ	τ_b	τ_c	ρ	τ_b	τ_c	ρ
Related Work												
BLEU1	57.9	63.7	74.0	32.2	32.3	40.4	17.9	9.3	22.3	45.8	46.2	63.0
BLEU4	54.8	60.4	70.5	30.6	30.8	38.7	16.9	8.7	21.0	46.4	46.9	63.7
CIDEr	63.1	69.8	79.3	43.6	43.9	54.3	24.6	12.0	29.3	48.1	48.8	65.0
CLIPScore	60.7	67.1	76.9	51.1	51.5	63.1	34.4	17.8	42.4	50.6	51.3	67.9
RefCLIPScore	66.3	73.3	82.2	52.0	52.4	63.7	36.4	18.8	44.7	56.8	57.6	74.7
CLIP+CIDEr	66.8	73.8	82.6	53.1	53.4	65.3	33.9	17.5	41.8	54.3	55.1	72.2
MID (Kim et al., 2022)	-	-	-	-	54.9	-	37.3	-	-	-	-	-
InfoMetIC (Hu et al., 2023)	-	-	-	-	54.2	-	36.3	-	-	-	59.2	-
InfoMetIC ⁺ (Hu et al., 2023)	-	-	-	-	55.5	-	36.6	-	-	-	59.3	-
PR-MCS (Kim et al., 2023)	-	-	-	-	50.6	65.6	-	-	-	-	-	-
VICR (Narins et al., 2024)	-	75.8	-	-	53.1	-	-	-	-	-	-	-
PACScore (Sarto et al., 2023)	-	-	-	53.9	54.3	-	36.0	18.6	-	51.5	55.7	-
RefPACScore (Sarto et al., 2023)	-	-	-	55.5	55.9	-	37.6	19.5	-	53.0	57.3	-
CLIPScore variants												
English CLIPScore	67.4	74.4	83.1	56.3	56.6	68.4	38.5	19.9	47.1	55.1	55.9	72.8
English RefCLIPScore	68.3	75.4	83.8	56.8	57.1	68.9	38.6	19.9	47.3	56.4	57.2	74.0
English CLIP+CIDEr	67.6	74.8	83.3	56.0	56.4	68.5	35.6	18.4	43.8	53.9	54.7	71.7
Multilingual CLIPScore	67.6	73.0	82.4	57.4	54.3	67.3	38.2	19.4	46.2	55.5	56.2	73.2
Multilingual RefCLIPScore	69.1	74.6	83.6	58.1	55.0	67.9	38.8	19.7	46.9	57.3	58.1	75.0
Multilingual CLIP+CIDEr	67.6	74.7	83.4	55.5	55.9	67.8	36.3	18.8	44.7	55.6	56.4	73.5
Finetuned CLIPScore	68.7	75.4	84.1	59.8	57.1	70.2	37.8	19.3	45.9	47.4	48.0	64.7
Finetuned RefCLIPScore	69.7	76.5	84.8	60.2	57.5	70.6	37.6	19.1	45.6	53.3	54.0	70.8
Finetuned CLIP+CIDEr	68.5	75.8	84.2	57.3	57.7	70.0	35.8	18.5	44.0	52.5	53.2	70.0

Table 2: Comparison between published results and our best English, multilingual and finetuned (also multilingual) CLIPScore models, considering settings (a) without human references, (b) using human references, and (c) combining CIDEr with CLIPScore when using human references.

forming multilingual model in English-only data, and also our finetuned model version.

Table 3 displays the correlation between multilingual CLIPScore values and human ratings, across the different languages. We observe that our finetuned version achieves a better correlation with human judgments in both the reference and reference-free settings, across all evaluated languages and for all correlation metrics. The finetuned CLIPScore is strongly correlated with human preferences in high-resource languages (i.e., English, French, German,

Spanish, and Chinese), and it also exhibits excellent performance in medium- and low-resource languages. We achieve average correlations of 67.9, 74.6, and 83.5, respectively with the τ_b , τ_c , and Spearman ρ metrics in the reference-free setting, and an average correlation of 68.9, 75.7, and 84.3, respectively with the τ_b , τ_c , and Spearman ρ using references. An additional table is provided in Appendix G regarding different loss variants for model finetuning, and with different model sizes. The findings in the appendix support the idea that

	CIDEr			Without Finetuning						With Finetuning					
				Multilingual CLIPScore			Multilingual RefCLIPScore			Finetuned Multilingual CLIPScore			Finetuned Multilingual RefCLIPScore		
	τ_b	τ_c	ρ	τ_b	τ_c	ρ	τ_b	τ_c	ρ	τ_b	τ_c	ρ	τ_b	τ_c	ρ
English	63.1	69.8	79.3	67.6	73.0	82.4	69.1	74.6	83.6	68.7	75.4	84.1	69.7	76.5	84.8
German	58.3	64.4	74.4	66.1	72.3	81.5	68.0	74.4	83.2	68.0	74.8	83.6	69.1	75.9	84.4
French	60.0	66.4	76.3	65.9	71.8	81.3	67.6	73.6	82.7	68.0	74.7	83.6	69.1	75.9	84.4
Spanish	61.3	67.8	77.6	66.3	72.7	81.8	66.3	72.7	81.8	67.8	74.5	83.5	68.9	75.7	84.3
Chinese	58.5	64.6	74.6	65.0	71.4	80.8	66.7	73.2	82.3	67.4	74.0	83.1	68.5	75.2	83.9
Portuguese	61.3	67.8	77.6	66.1	72.2	81.6	67.8	74.1	83.0	67.9	74.6	83.5	69.0	75.8	84.3
Italian	60.7	67.1	77.0	65.8	72.1	81.4	67.5	74.0	82.9	67.9	74.6	83.5	68.9	75.7	84.3
Russian	54.0	59.6	69.4	65.0	71.4	80.7	66.7	73.3	82.2	67.6	74.2	83.2	68.8	75.6	84.2
Korean	58.1	64.3	74.1	64.3	70.5	80.0	66.3	72.7	81.8	67.4	74.1	83.1	68.3	75.0	83.8
Dutch	58.1	64.2	74.2	65.8	72.0	81.4	67.7	74.1	83.0	68.2	74.9	83.8	69.0	75.8	84.3
Average	59.3	65.6	75.5	65.8	71.9	81.3	67.4	73.7	82.7	67.9	74.6	83.5	68.9	75.7	84.3

Table 3: Correlation between M-CLIPScore values and human rankings, considering machine-translated versions of the validated image caption rating dataset into 9 different languages besides the original English. The last row presents macro-averaged correlation results across all the languages (including English).

smaller models can have higher human judgment correlation gains using our finetuning strategy, compared to the respective original models.

In an ideal scenario, i.e. assuming perfect machine translation results and that CLIP performs equally well across the languages, the correlations between CLIPScore values across the different languages would be equal one, signifying a perfect alignment. To explore deviations from this behaviour, we can use heatmaps to visually represent the interrelationships between CLIPScore values across languages. In Figure 1, we plot the Pearson correlations between the best multilingual CLIPScore and our finetuned version for different languages, considering either (a) the complete set of original/translated instances from the VICR dataset, (b) the subset of instances with COMETKiwi (Rei et al., 2022) translation quality scores below the 25th percentile value for each language, and (c) the subset of instances with COMETKiwi scores above the 75th percentile value for each language.

Looking at the left plot in Figure 1, which features the Pearson correlations when considering the entire VICR dataset, we observe consistently high values across all languages. As expected, the plots confirm that the CLIPScore correlations do depend on the quality of the translated captions, and that the most significant differences occur primarily in languages using non-Latin scripts. In turn, the upper diagonal of the second and third plots contains the correlations when considering only high-quality translations, corresponding to slightly higher values compared with the values of their respective lower diagonal. We also observe a significant improvement across the correlation between the languages when using our finetuned CLIPScore version. Although this improvement

is somewhat expected, since the finetuned model saw "in distribution" data, it may also be an indication of the quality of our multilingual data during training, thus supporting the fact that our translation strategy maintains the quality of the data in the different languages.

4.3.3 Multilingual Classification

This section explores the robustness of the multilingual CLIPScore assessments through different types of classification tasks. Inspired by previous work (Hessel et al., 2021; Sarto et al., 2023) which assessed accuracy in English-only benchmarks, our goal is to delve deeper into the nuanced realm of multilingual and multicultural understanding.

Robustness to linguistic phenomena: One of our experiments used machine translated versions of the VALSE dataset, explicitly designed to evaluate the robustness to particular phenomena, such as inconsistencies in numeric quantities or spatial relations (Parcalabescu et al., 2021). VALSE specifically comprises seven tests that encompass a range of linguistic structures. In each test, a model is presented with a visual input and is tasked with distinguishing genuine captions from altered versions (i.e., foils), where a word or phrase has been modified to exhibit a specific linguistic phenomenon. Additional information about the dataset is provided in Appendix D. Table 4 displays the average performance across different language variants, considering that the CLIPScore values regarding the true captions should be higher than the ones related to the corresponding foils. A more detailed breakdown of the performance across the different tasks within VALSE is available in Appendix G. The results show that our finetuned model delivers the highest average performance across nearly

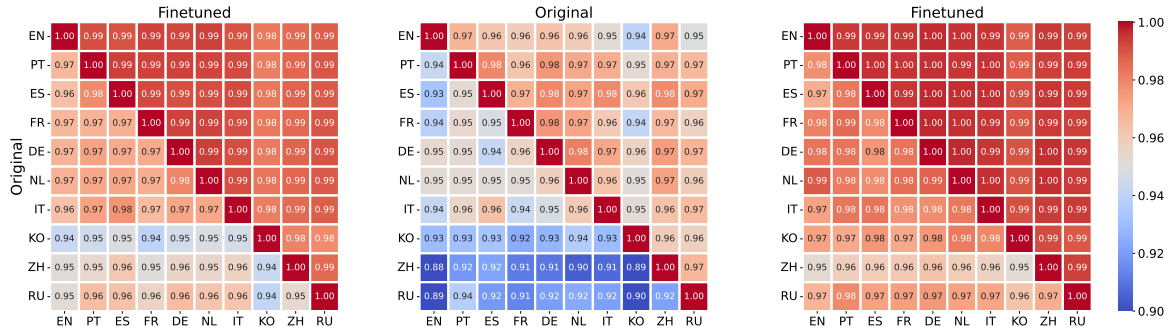


Figure 1: Pearson correlation scores between different languages, for the best multilingual CLIPScore model and our finetuned version. The first heatmap considers the complete set of instances from the VICR dataset, reporting results for both the original and finetuned model versions (lower/upper diagonal values). The second and third heatmaps consider the subset of instances with COMETKiwi scores below/above the 25th/75th percentile value for each language (lower/upper diagonal values) for the best performing multilingual CLIPScore model and our finetuned model version, respectively.

all languages for both model size variants. Compared to the models reported in the original VALSE dataset paper, our finetuned model was only outperformed by the multi-task ViLBERT 12-in-1 model proposed by Lu et al. (2020).

Classification of multicultural instances: We also used naturally multilingual datasets (i.e., XVNLI and MaRVL) to assess the multilingual and multicultural capabilities of CLIPScore models.

Each instance in the XVNLI dataset contains an image-caption pair and a categorical label associated with the relationship between the pair. This label can be either (a) contradiction, (b) neutral, or (c) entailment. Based on these labels, we defined three multilingual classification experiments under this scenario, leveraging concordant/discordant instances as illustrated in Figure 2:

Experiment 1: This setting only considers contradiction and entailment instances, under the assumption that the order of the CLIPScore values should match the order of the labels.

Experiment 2: In a more challenging scenario, we can consider a larger set of duplets and the ordering between the three possible labels.

Experiment 3: In this case, we also consider the three possible labels, but we now assess triples of instances A , B , and C from the dataset, sharing the same image. We only assume a correct classification when we achieve a perfect match between the order of the labels and the CLIPScore values.

In the case of the MaRVL dataset, each instance consists of a caption, a pair of images, and a Boolean label with the value *true* when the caption accurately matches the images and *false* when the caption is incorrect (e.g. because its contents only

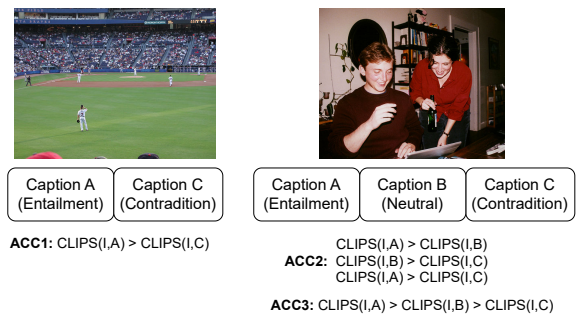


Figure 2: The three different XVNLI multilingual classification experiments, where accuracy is defined with basis on comparisons between CLIPScore values.

describe at maximum one of the images instead of both). The data can be analyzed considering four instances simultaneously, sharing the same caption but featuring distinct pairs of images.

Experiment 1: We consider the scenario where the CLIPScore for the image that best aligns with the caption should be higher than the CLIPScore for the image that least aligns with the caption.

Experiment 2: In a more challenging scenario, we consider sets of four instances sharing the same caption and decide on a correct classification only if all best-aligning captions have a CLIPScore higher than all the least aligning captions.

A more detailed explanation for both the XVNLI and MaRVL experiments is given in Appendix E.

Table 5 contains the results for the two multilingual classification tasks. In the XVNLI experiments, as expected, we notice a consistent trend where results for Experiment 3 are worse than those for Experiment 2, which in turn are worse than those for Experiment 1. The performance is generally high in Experiments 1 and 2, for all the languages except Swahili and Tamil. The lower perfor-

Finetuning method	Proposed Models										VALSE Models	
	English	German	French	Spanish	Chinese	Portuguese	Italian	Russian	Korean	Dutch	Model	English
Both	69.7	67.3	67.5	67.9	66.9	69.0	66.0	67.1	60.6	66.2	CLIP	64.0
Pearson	69.2	67.2	66.8	67.4	67.3	67.2	66.0	65.0	59.4	65.8	LXMERT	59.6
Contrastive	68.4	65.0	65.5	65.7	65.6	66.4	64.6	64.8	59.2	63.7	ViLBERT	46.4
None	67.6	64.8	64.2	65.6	64.0	65.4	64.4	63.4	58.4	62.5	12-in-1	75.1

Table 4: Average accuracy scores for the different classification tasks present in the VALSE dataset and its multilingual variants, considering different CLIP models.

	XVNLI						MaRVL			
	Accuracy 1		Accuracy 2		Accuracy 3		Accuracy 1		Accuracy 2	
	Original	Finetuned	Original	Finetuned	Original	Finetuned	Original	Finetuned	Original	Finetuned
English	92.3	91.5	80.7	81.6	47.7	51.0	91.1	91.7	81.4	80.7
Indonesian	89.9	92.0	80.5	80.8	50.0	49.0	92.4	93.4	82.3	83.3
Mandarin	88.1	90.7	78.8	79.5	46.7	46.7	89.4	91.5	80.1	81.3
Swahili	67.1	66.0	60.7	60.8	24.0	22.4	85.0	85.0	65.3	65.3
Tamil	71.6	68.7	62.8	62.7	25.3	25.7	85.7	88.0	74.7	77.5
Turkish	87.8	88.1	77.9	79.7	44.1	49.7	93.5	92.9	87.0	86.0
Arabic	84.6	85.4	76.3	76.3	44.7	41.8	-	-	-	-
French	86.7	90.5	77.9	79.9	45.4	47.7	-	-	-	-
Spanish	86.2	87.8	77.7	78.4	46.7	45.1	-	-	-	-
Russian	87.5	88.9	78.0	79.2	44.4	46.7	-	-	-	-
Overall	84.2	85.0	75.1	75.9	41.9	42.6	89.5	90.4	78.5	79.0

Table 5: Accuracy for different classification tasks defined over the datasets derived from XVNLI and MaRVL.

mance in these languages can perhaps be attributed to a lower quality in the machine translation results, and to the ability of the multilingual CLIP model in handling text in these lower resource languages. However, in the MaRVL experiments where the instances never involve machine translation, the task is performed with high accuracy across all languages. Once again, we observe that our finetuned CLIPScore version is capable of achieving moderate performance gains when compared to the original model version, across all evaluated multicultural classification tasks.

While the experiments with the XVNLI and MaRVL datasets provide interesting insights into the effectiveness of multilingual CLIPScore models, they also involve several important limitations. For instance, considering the XVNLI experiments, previous studies have reported good results in multimodal inference leveraging CLIP (Song et al., 2022). However, the authors of SNLI-VE (Xie et al., 2019), from which XVNLI is derived, noted that good performance (i.e., an accuracy up to 67%) can be achieved when looking only at the information in the textual hypothesis, without the visual premise. This points to significant biases in the XVNLI data. In the case of the MaRVL experiments, given that the captions refer to a pair instead of individual images, the CLIPScore values can be unreliable when attempting to match images to textual sentences. Previous studies have noted that CLIP models can treat inputs as a bag-of-words and suffer from a concept association bias (Yamada et al., 2022), e.g. ignoring the missing information

when two concepts are present in one of the inputs while the other only contains a single concept. Hence, while already useful, there is room for future work and analysis of the robustness of caption evaluation metrics, accounting for the discussed limitations.

5 Conclusions

We studied the potential of CLIPScore variants in multilingual image captioning evaluation, considering a variety of settings and demonstrating that multilingual CLIP maintains or even surpasses performance on English benchmarks, while enabling versatile multilingual assessments.

Performance improves with larger model sizes and increased training data, aligning with the scalability of CLIP. On English data, large multilingual CLIP models can even outperform more advanced methods that extend the original CLIPScore metric or involve specific training for captioning evaluation. Experiments with machine-translated data reveal that a finetuned version of multilingual CLIPScore strongly correlates with human judgments across languages of varying complexity, and enhances the robustness of CLIP to different linguistic phenomena across languages. Additionally, assessments with natively multilingual and multicultural datasets, specifically with data from the XVNLI and MaRVL benchmarks repurposed for evaluating captioning metrics, reaffirm the ability of multilingual CLIPScore models to consistently provide high-quality assessments in varied settings.

592 Limitations and Ethical Considerations

593 Although our work does not raise new ethical issues
594 within the domain of vision-language models (e.g.,
595 we conducted our experiments on public datasets,
596 carefully designed for academic research and ex-
597 tensively used in previous studies), there are some
598 general important concerns.

599 Models like CLIP are, for instance, notorious for
600 their internal biases, e.g. inherited from the train-
601 ing data itself. We therefore recommend caution
602 in the use of the approach proposed in this paper,
603 and anticipate further research into the specific is-
604 sue of model biases, before relying on our work
605 beyond research environments. Another important
606 limitation in the work reported on this paper con-
607 cerns the use of machine translated data in some
608 of the evaluation experiments, which despite our
609 best efforts to avoid translation errors can still lead
610 to different types of biases and to the reliance on
611 artificially impoverished language. The develop-
612 ment of manually curated benchmarks, specifically
613 designed for the assessment of multilingual met-
614 rics for image captioning evaluation, is left as an
615 important challenge for future work.

616 We also note that we used Github Copilot³ dur-
617 ing the development of our research work, and we
618 used ChatGPT⁴ for minor verifications during the
619 preparation of this manuscript.

620 References

621 Somak Aditya, Yezhou Yang, Chitta Baral, Cornelia
622 Fermuller, and Yiannis Aloimonos. 2015. From im-
623 ages to sentences through scene description graphs
624 using commonsense reasoning and knowledge. *arXiv*
625 *preprint arXiv:1511.03292*.

626 Saba Ahmadi and Aishwarya Agrawal. 2023. An ex-
627 amination of the robustness of reference-free im-
628 age captioning evaluation metrics. *arXiv preprint*
629 *arXiv:2305.14998*.

630 Duarte M Alves, José Pombal, Nuno M Guerreiro, Pe-
631 dro H Martins, João Alves, Amin Farajian, Ben Pe-
632 ters, Ricardo Rei, Patrick Fernandes, Sweta Agrawal,
633 et al. 2024. Tower: An open multilingual large
634 language model for translation-related tasks. *arXiv*
635 *preprint arXiv:2402.17733*.

636 Emanuele Bugliarello, Fangyu Liu, Jonas Pfeiffer, Siva
637 Reddy, Desmond Elliott, Edoardo Maria Ponti, and
638 Ivan Vulić. 2022. IGLUE: A benchmark for trans-
639 fer learning across modalities, tasks, and languages.

In Proceedings of the International Conference on
Machine Learning. 640
641

Angela Fan, Shruti Bhosale, Holger Schwenk, Zhiyi
Ma, Ahmed El-Kishky, Siddharth Goyal, Mandeep
Baines, Onur Celebi, Guillaume Wenzek, Vishrav
Chaudhary, Naman Goyal, Tom Birch, Vitaliy
Liptchinsky, Sergey Edunov, Michael Auli, and Ar-
mand Joulin. 2021. Beyond english-centric multilin-
gual machine translation. *Journal of Machine Learn-*
ing Research, 22(107). 642
643
644
645
646
647
648
649

Gregor Geigle, Abhay Jain, Radu Timofte, and
Goran Glavaš. 2023. mblip: Efficient bootstrap-
ping of multilingual vision-llms. *arXiv preprint*
arXiv:2307.06930. 650
651
652
653

Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le
Bras, and Yejin Choi. 2021. CLIPScore: A reference-
free evaluation metric for image captioning. *arXiv*
preprint arXiv:2104.08718. 654
655
656
657

Micah Hodosh, Peter Young, and Julia Hockenmaier.
2013. Framing image description as a ranking task:
Data, models and evaluation metrics. *Journal of*
Artificial Intelligence Research, 47:853–899. 658
659
660
661

Anwen Hu, Shizhe Chen, Liang Zhang, and Qin
Jin. 2023. InfoMetIC: An informative metric for
reference-free image caption evaluation. *arXiv*
preprint arXiv:2305.06002. 662
663
664
665

Maxime Kayser, Oana-Maria Camburu, Leonard
Salewski, Cornelius Emde, Virginie Do, Zeynep
Akata, and Thomas Lukasiewicz. 2021. e-vil: A
dataset and benchmark for natural language explana-
tions in vision-language tasks. In *Proceedings of the*
IEEE/CVF International Conference on Computer
Vision. 666
667
668
669
670
671
672

Jin-Hwa Kim, Yunji Kim, Jiyoung Lee, Kang Min Yoo,
and Sang-Woo Lee. 2022. Mutual information diver-
gence: A unified metric for multimodal generative
models. In *Proceedings of the Annual Meeting on*
Neural Information Processing Systems. 673
674
675
676
677

Yongil Kim, Yerin Hwang, Hyeonju Yun, Seunghyun
Yoon, Trung Bui, and Kyomin Jung. 2023. PR-MCS:
Perturbation robust metric for multilingual image
captioning. *arXiv preprint arXiv:2303.08389*. 678
679
680
681

Fangyu Liu, Emanuele Bugliarello, Edoardo Maria
Ponti, Siva Reddy, Nigel Collier, and Desmond
Elliott. 2021. Visually grounded reasoning
across languages and cultures. *arXiv preprint*
arXiv:2109.13238. 682
683
684
685
686

Yinhan Liu, Jiatao Gu, Naman Goyal, Xian Li, Sergey
Edunov, Marjan Ghazvininejad, Mike Lewis, and
Luke Zettlemoyer. 2020. Multilingual denoising pre-
training for neural machine translation. *Transactions*
of the Association for Computational Linguistics, 8. 687
688
689
690
691

Jiasen Lu, Vedanuj Goswami, Marcus Rohrbach, Devi
Parikh, and Stefan Lee. 2020. 12-in-1: Multi-task 692
693

³<https://github.com/features/copilot>

⁴<https://openai.com/chatgpt/>

694	vision and language representation learning. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pages 10437–10446.	750
695		751
696		752
697	Lothar D Narins, Andrew Scott, Aakash Gautam,	753
698	Anagha Kulkarni, Mar Castanon, Benjamin Kao,	754
699	Shasta Ihorn, Yue-Ting Siu, James M Mason, Alexander Blum, et al. 2024. Validated image caption rating dataset. In <i>Proceedings of the Annual Meeting on Neural Information Processing Systems</i> .	755
700		756
701		
702		
703	Letitia Parcalabescu, Michele Cafagna, Lilitta Muradjan, Anette Frank, Iacer Calixto, and Albert Gatt. 2021. Valse: A task-independent benchmark for vision and language models centered on linguistic phenomena. <i>arXiv preprint arXiv:2112.07566</i> .	
704		
705		
706		
707		
708	Haoyi Qiu, Zi-Yi Dou, Tianlu Wang, Asli Celikyilmaz, and Nanyun Peng. 2023. Gender biases in automatic evaluation metrics for image captioning. In <i>Proceedings of the Conference on Empirical Methods in Natural Language Processing</i> .	
709		
710		
711		
712		
713	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In <i>Proceedings of the International conference on machine learning</i> .	
714		
715		
716		
717		
718		
719	Rita Ramos, Emanuele Bugliarello, Bruno Martins, and Desmond Elliott. 2024. PAELLA: Parameter-efficient lightweight language-agnostic captioning model. In <i>Findings of the Annual Conference of the North American Chapter of the Association for Computational Linguistics</i> .	
720		
721		
722		
723		
724		
725	Rita Ramos, Bruno Martins, and Desmond Elliott. 2023a. LMCap: Few-shot multilingual image captioning by retrieval augmented language model prompting. In <i>Findings of the Association for Computational Linguistics</i> .	
726		
727		
728		
729		
730	Rita Ramos, Bruno Martins, Desmond Elliott, and Yova Kementchedjieva. 2023b. SmallCap: Lightweight image captioning prompted with retrieval augmentation. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> .	
731		
732		
733		
734		
735	Ricardo Rei, Nuno M Guerreiro, José Pombal, Daan van Stigt, Marcos Treviso, Luisa Coheur, José GC de Souza, and André FT Martins. 2023. Scaling up COMETKiwi: Unbabel-IST 2023 submission for the quality estimation shared task. <i>arXiv preprint arXiv:2309.11925</i> .	
736		
737		
738		
739		
740		
741	Ricardo Rei, Marcos Treviso, Nuno M Guerreiro, Chrysoula Zerva, Ana C Farinha, Christine Maroti, José GC De Souza, Taisiya Glushkova, Duarte M Alves, Alon Lavie, et al. 2022. Cometkiwi: Ist-unbabel 2022 submission for the quality estimation shared task. <i>arXiv preprint arXiv:2209.06243</i> .	
742		
743		
744		
745		
746		
747	Sara Sarto, Manuele Barraco, Marcella Cornia, Lorenzo Baraldi, and Rita Cucchiara. 2023. Positive-augmented contrastive learning for image and video captioning evaluation. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> .	753
748		754
749		755
		756
		757
		758
		759
		760
		761
		762
		763
		764
		765
		766
		767
		768
		769
		770
		771
		772
		773
		774
		775
		776
		777
		778
		779
		780
		781
		782
		783
		784

A The CLIPScore Metric

We now formally describe the CLIPScore and RefCLIPScore metrics (Hessel et al., 2021), which in our study are assessed in multilingual image captioning scenarios. In brief, we have that CLIPScore is based on a modified cosine similarity between representations for the input image and the caption under evaluation. The image and the caption are both passed through the respective feature extractors from a given CLIP model. Then, we compute the cosine similarity of the resultant embeddings, adjusting the resulting value through a re-scaling operation. For an image with visual CLIP embedding \mathbf{v} and a candidate caption with textual CLIP embedding \mathbf{c} , a re-scaling parameter is set as $w = 2.5$ and we compute the corresponding CLIPScore as follows:

$$\text{CLIPScore}(\mathbf{c}, \mathbf{v}) = w \times \max(\cos(\mathbf{c}, \mathbf{v}), 0). \quad (3)$$

To compute a corpus-level CLIPScore, e.g. for evaluating the overall quality of a captioning method over a given dataset of images, we can simply average over all the image-candidate pairs.

Note that CLIPScore does not depend on the availability of underlying references for each of the images in an evaluation dataset. However, an extension named RefCLIPScore was also proposed, which additionally extracts the vector representations \mathbf{R} of each available reference with the CLIP text encoder, and computes the harmonic mean of the CLIPScore value from Equation 3, and the maximal reference cosine similarity:

$$\begin{aligned} \text{RefCLIPScore}(\mathbf{c}, \mathbf{R}, \mathbf{v}) = \\ \text{H-Mean}(\text{CLIPScore}(\mathbf{c}, \mathbf{v}), \\ \max_{\mathbf{r} \in \mathbf{R}}(\max \cos(\mathbf{c}, \mathbf{r}), 0)). \end{aligned} \quad (4)$$

B The Correlation Metrics

This appendix presents a formal definition of the metrics used in the correlation experiments.

Seeing each of our evaluation datasets as sets of n observations with the form $(\hat{y}_1, y_1), \dots, (\hat{y}_n, y_n)$, for CLIPScore values \hat{y}_i and reference ratings y_i , the Spearman correlation coefficient ρ is defined as the Pearson correlation between the results of converting the scores \hat{y}_i and y_i to ranks.

Instead of using ranks, we can also define any pair of observations (\hat{y}_i, y_i) and (\hat{y}_j, y_j) , where $i < j$, as concordant (or otherwise discordant) if the sort order of the instances agrees (i.e. if

either both $\hat{y}_i > \hat{y}_j$ and $y_i > y_j$ holds, or both $\hat{y}_i < \hat{y}_j$ and $y_i < y_j$). Based on pairs, the Kendall τ correlation coefficient assesses the strength of association between the CLIPScore values and the reference ratings, with the τ_b variant accounting for ties and being defined as:

$$\tau_B = \frac{n_c - n_d}{\sqrt{(n_0 - n_1)(n_0 - n_2)}}, \quad (5)$$

where n_c is the number of concordant pairs, n_d the number of discordant pairs, $n_0 = n(n-1)/2$, $n_1 = \sum_i t_i(t_i-1)/2$, $n_2 = \sum_j u_j(u_j-1)/2$, t_i is the number of tied values in the i^{th} group of ties for the CLIPScore, and u_j is the number of tied values in the j^{th} group of ties for the reference ratings.

In turn, τ_c accounts with the fact that the underlying scales of the scores are different for CLIPScore and the reference ratings, being defined as:

$$\tau_c = \frac{n_c - n_d}{n_0} \times \frac{n-1}{n} \times \frac{m}{m-1}, \quad (6)$$

where m is the number of values in the ranking scale for the reference ratings.

C The Machine Translation Scheme

This appendix describes the translation scheme that was used to machine translate the datasets used in our experiments. This scheme is design to mitigate low-quality translations, or hallucinations generated by the machine translation model, thus providing a reliable dataset at the end. We specifically used a large (i.e., 1.2 billion parameters) open-access multilingual machine translation model named M2M100 (Fan et al., 2021), available on the HuggingFace⁵ model hub. M2M100 was trained on a range of high and low-resource languages from different families and using different scripts, achieving state-of-the-art performance across a diverse set of 100 languages.

While machine translated data allows us to assess multilingual captioning metrics, the results will depend not only on the performance of the metrics but also on the quality of the translations. Low-quality translations, or hallucinations generated by the translation model, will impact the caption and break our assumption that human ratings for the English data can be transferred across languages. To address this issue, we propose to use the COMETKiwi (Rei et al., 2023) machine translation quality estimation metric to control for translation

⁵https://huggingface.co/facebook/m2m100_1.2B

877 quality, assessing the impact of low quality transla- 924
878 tions on the observed results. 925

879 We specifically began by translating the VICR 926
880 dataset, featuring English captions with human rat- 927
881 ings and also reference captions originally from 928
882 the MSCOCO and Flickr8K datasets. For each 929
883 caption, whether a candidate or a reference, we 930
884 return 25 translations using a beam search tech- 931
885 nique with 100 beams. Subsequently, we filtered 932
886 the candidates with a language checker, to ensure 933
887 proper translation into the intended language. After 934
888 the language check, we selected for each instance 935
889 the translation that scored higher based on a large 936
890 COMETKiwi model⁶. 937

891 D Description of the Datasets

892 The following datasets were used in the tests that 938
893 assessed correlation with human judgment. 939

- 894 • **Flickr8K-Expert (Hodosh et al., 2013)**: This 942
895 dataset comprises 16,992 expert human judg- 943
896 ments for 5,664 image-caption pairs from 944
897 Flickr8K. Human assessors graded captions 945
898 on a scale of 1 to 4, where 4 indicates a caption 946
899 that accurately describes the image without er- 947
900 rors, and 1 signifies a caption unrelated to the 948
901 image. 949
- 902 • **Flickr8K-CF (Hodosh et al., 2013)**: This 950
903 dataset consists of 145,000 binary quality 951
904 judgments collected with CrowdFlower, in- 952
905 volving 47,830 image-caption pairs with 953
906 1,000 unique Flickr8K images. Each pair re- 954
907 ceived at least three binary judgments, and we 955
908 use the proportion of *yes* annotations as the 956
909 score for each pair. 957
- 910 • **Composite (Aditya et al., 2015)**: This dataset 958
911 contains 13,146 image-caption pairs taken 959
912 from MSCOCO (2007 images), Flickr8K (997 960
913 images), and Flickr30K (991 images). Each 961
914 image originally had five reference captions. 962
915 One of these references was chosen for human 963
916 rating and subsequently removed from the ref- 964
917 erence set that is to be used when assessing 965
918 evaluation metrics. 966
- 919 • **VICR (Narins et al., 2024)**: The Validated 967
920 Image Caption Rating (VICR) dataset features 968
921 68,217 ratings, collected through a gamified 969
922 approach, for 15,646 image-caption pairs in- 970
923 volving 9,990 distinct images. The authors of 971

⁶<https://huggingface.co/Unbabel/wmt23-cometkiwi-da-xxl>

the dataset demonstrated that it exhibits a su- 924
perior inter-rater agreement compared to other 925
alternatives (e.g., an improvement of 19% in 926
Fleiss’ κ when compared to the agreement for 927
the Flickr8K-Expert dataset), and it features a 928
more balanced distribution across various lev- 929
els of caption quality. In our tests, we used the 930
test split of the VICR dataset, which includes 931
3,161 image-caption pairs, with 2,000 images 932
from the MSCOCO 2014 validation dataset and 933
1,161 images from the Flickr8K dataset. 934
When using VICR to finetune CLIP models 935
with a contrastive loss, we used the original 936
image captions from MSCOCO or Flickr8K. 937

All the previous datasets are originally available 938
only for English, but we translated them to nine 939
other different languages using the approach de- 940
scribed in Appendix C. 941

For the experiments that assessed accuracy in 942
terms of distinguishing correct vs incorrect cap- 943
tions, we used the following datasets. 944

- 945 • **VALSE (Parcalabescu et al., 2021)**: VALSE 945
946 is designed to test visio-linguistic grounding 946
947 capabilities on specific linguistic phenomena. 947
948 It is composed by seven tasks, each with the 948
949 same structure: given a visual input, a model 949
950 is asked to distinguish real captions from foils, 950
951 where a foil is constructed from a caption 951
952 by altering a word or phrase that realizes a 952
953 specific linguistic phenomenon. The tests in- 953
954 clude: (a) existential quantifiers, where mod- 954
955 els need to differentiate between examples (i) 955
956 where there is no entity of a certain type or 956
957 (ii) where one or more of these entities are 957
958 visible in an image; (b) plurality, where mod- 958
959 els need to distinguish between noun phrases 959
960 denoting a single entity in an image (*exactly* 960
961 *one flower*), versus multiple entities (*some* 961
962 *flowers*); (c) counting, where models need 962
963 to differentiate between examples where the 963
964 specific number of entities in the associated 964
965 image is correct or incorrect, given the state- 965
966 ment; (d) spatial relations, where Models need 966
967 to distinguish between different spatial rela- 967
968 tions, with foils differing from the original 968
969 caption only by the replacement of a spatial 969
970 preposition; (e) actions, particularly (i) action 970
971 replacement and (ii) actant swapping, where 971
972 models need to (i) identify whether an action 972
973 mentioned in the text matches the action seen 973
974 in the image (e.g., *a man shouts* versus *smiles* 974

at a woman), and (ii) correctly identify the participants of an action and the roles they play (e.g., is it the man who is shouting or is it the woman); (f) coreference, where models need to perform pronominal coreference resolution, encompassing cases where (i) the pronoun has a noun (phrase) antecedent and pronoun and (noun) phrase are both grounded in the visual modality (e.g., in *a woman is driving a motorcycle*, is she wearing a helmet?), and cases where (ii) the pronoun refers to a region in the image or even to the entire image (e.g., *is this outside?*); (g) foil-it cases, in which the foil minimally differs from the original caption, only by swapping a important noun.

- **XVNLI (Bugliarello et al., 2022):** XVNLI is a multilingual dataset for evaluating vision-language inference, challenging models to predict entailment relationships between a textual hypothesis and an image premise. XVNLI includes high/mid-resource languages like Arabic, French, Spanish, Russian, and English. This dataset includes 1,164 instances per language, each featuring an image and two captions in different languages. There are 357 unique images in total.
- **MaRVL (Liu et al., 2021):** MaRVL follows a format similar to the English NLVR2 dataset (Suhr et al., 2018) and is designed as a multicultural vision-language reasoning dataset, where the goal is to determine the correctness of a sentence about a pair of images. MaRVL predominantly comprises very low-resource languages: Indonesian, Chinese, Swahili, Tamil, and Turkish. This dataset includes around one thousand instances per language, each featuring two image and one caption. There are 1,411 unique captions in total. The English content is composed by collecting the reverse translations from the low-resource languages into English, as provided by the authors on the original GitHub⁷.

For model training, besides instances in the training split from the aforementioned VICR dataset, we also used data from the natively multilingual CrossModal-3600 dataset (i.e., XM3600, in short).

- **XM3600 (Thapliyal et al., 2022):** This is a geographically-diverse set of 3600 images an-

notated with human-generated reference captions in 36 languages. The images were selected from all across the world, covering regions where the 36 languages are spoken, and consistently annotating captions in terms of style across all languages, while avoiding annotation artifacts due to direct translation.

E Multicultural Experiments

This appendix details the datasets and experimental settings that were considered for the tests including natively multilingual and multicultural data.

E.1 Settings for the XVNLI Experiments

Each instance in the XVNLI dataset contains an image-caption pair and a categorical label associated with the relationship between the pair. This label can be either (a) contradiction, (b) neutral, or (c) entailment. With basis on the labels, we defined three multilingual classification experiments under this scenario, leveraging concordant/discordant instances as illustrated in Figure 2:

Experiment 1: This setting only considers instances with the extreme label classes (i.e., contradiction and entailment), noting that some previous studies have pointed to the fact that SNLI-VE, from which XVNLI is derived, has some problems in the annotations for the neutral class (Kayser et al., 2021). We compare pairs of instances A and B with the same image, in which the label associated with A differs from the label associated with B . When computing CLIPScore values individually for the instances A and B , the order of the CLIPScore values should match the order of the labels (i.e., contradiction < entailment).

Experiment 2: In a more challenging scenario, we can consider a larger set of instances and the ordering between the three possible labels (i.e. entailment > neutral > contradiction), i.e. including also the neutral class. Similarly to the previous case, by fixing an image and comparing pairs of captions associated with that image with different labels, we assess the matching of the order between the labels against the CLIPScore values.

Experiment 3: In this case, we also consider the three possible labels, but we now assess triples of instances A , B , and C from the dataset, sharing the same image. We only assume a correct classification when we achieve a perfect match between the order of the labels and the CLIPScore values.

⁷<https://github.com/marvl-challenge/marvl-code/tree/master/data>



Figure 3: Multilingual CLIPScore values for image-caption pairs featuring concepts biased to particular languages.

E.2 Settings for the MaRVL Experiments

In the case of the MaRVL dataset, each instance consists of a caption, a pair of images, and a Boolean label with the value *true* when the caption accurately matches the images, and *false* when the caption is incorrect (e.g., because its contents only describe at maximum one of the images instead of both). The data can be analyzed considering four instances at a time, sharing the same caption but featuring distinct pairs of images. MaRVL was designed in such a way that, within these four instances, two of them are labeled as *true* while the remaining two are labeled as *false*. We consider two multilingual classification experiments under this scenario, defined as follows:

Experiment 1: We draw comparisons between pairs of instances with distinct labels. For the instance labeled as *true*, we compute the CLIPScore values for both images associated with the caption and select the maximum, obtaining the score for the image that best aligns with the caption. Conversely, we perform a similar computation for the instance labeled as *false*, this time choosing the minimum CLIPScore value, which results in the score for the image that least aligns with the caption, presumably the incorrect image. The maximum CLIPScore value in an instance labeled as *true* should be higher than the minimum CLIPScore value of an instance labeled as *false*.

Experiment 2: In a more challenging scenario, we consider sets of four instances sharing the same caption, and decide on a correct classification only if all maximum CLIPScore values of the *true* instances are higher than all the minimum CLIPScore values of the *false* instances.

F A Qualitative Study with Captions Featuring Culturally Related Concepts

We performed a small qualitative study on image-caption pairs that feature concepts where some languages should exhibit a particular bias (e.g., *codfish* in the case of Portugal, *paella* for Spain, *beer* for Germany, *croissant* for France, *ushanka* for Russia, and *cheongsam* for China). We attempted to see if the multilingual CLIPScore could distinguish between two plausible captions, where one mentions a specific concept that should better match the image. Figure 3 shows that the multilingual CLIPScore is indeed capable of distinguishing nuanced multicultural concepts and favouring culturally specific captions over generic ones.

G Additional Classification Results

Table 6 presents classification results on the different tasks from the VALSE dataset, separately for each of the considered languages and comparing CLIP models of different sizes under different finetuning strategies (i.e., without model finetuning, considering only the contrastive loss, only the Pearson correlation loss, or the combined loss function). The results show that the smaller CLIP model can achieve significantly higher gains from finetuning, approaching the performance of the larger CLIP model. Better results are, in general, obtained when considering the combined loss function.

In turn, Table 7 presents correlation results on the VICR dataset, separately for each language and comparing the same two CLIP models under the different finetuning strategies. Results again show that the smaller CLIP model approach the performance of the larger model, with better results consistently obtained when considering the combined loss function.

	Multilingual CLIP ViT-B									Multilingual CLIP ViT-H								
	Contrastive			Pearson			Combined			Contrastive			Pearson			Combined		
	τ_b	τ_c	ρ	τ_b	τ_c	ρ	τ_b	τ_c	ρ	τ_b	τ_c	ρ	τ_b	τ_c	ρ	τ_b	τ_c	ρ
ENG	65.8	72.2	81.5	66.7	73.2	82.3	67.2	73.5	82.6	68.6	75.0	83.8	64.5	70.6	80.6	68.7	75.4	84.1
GER	65.1	71.4	80.8	65.6	72.1	81.4	66.5	72.6	81.9	68.0	74.3	83.2	62.9	69.1	79.3	68.0	74.8	83.6
FRE	65.4	71.7	81.1	66.0	72.4	81.7	66.7	72.9	82.1	68.0	74.4	83.3	62.8	69.0	79.3	68.0	74.7	83.6
SPA	65.1	71.3	80.8	65.9	72.6	81.7	66.4	72.7	81.9	67.8	74.3	83.2	62.4	68.6	78.9	67.8	74.5	83.5
CHI	64.4	70.6	80.2	64.7	71.3	80.7	65.7	71.9	81.3	67.4	73.8	82.8	62.1	68.3	78.6	67.4	74.0	83.1
POR	65.2	71.5	80.9	65.7	72.1	81.4	66.5	72.7	81.9	67.9	74.3	83.2	62.7	68.9	79.2	67.9	74.6	83.5
ITA	65.0	71.3	80.7	65.6	72.2	81.4	66.3	72.6	81.8	67.9	74.3	83.2	62.5	68.7	79.0	67.9	74.6	83.5
RUS	64.6	70.8	80.3	65.1	71.6	80.9	65.9	71.9	81.4	67.4	73.8	82.8	62.7	68.9	79.2	67.6	74.2	83.2
KOR	63.8	70.0	79.6	64.3	70.8	80.2	65.1	71.2	80.7	67.3	73.7	82.7	62.1	68.3	78.6	67.4	74.1	83.1
DUT	65.2	71.4	80.8	65.8	72.4	81.6	66.6	72.8	82.0	68.1	74.5	83.4	62.9	69.1	79.3	68.2	74.9	83.8
AVG	65.0	71.2	80.7	65.5	72.1	81.3	66.2	72.5	81.8	67.8	74.2	83.2	62.8	68.9	79.2	67.9	74.6	83.5

Table 7: Correlation between CLIPScore values and human rankings, using multimodal models finetuned with different loss functions. The last row presents macro-averaged correlation results across all the languages (including English). Non-English results above the English baseline are shown in bold.