

# Towards Cross-Lingual Explanation of Artwork in Large-scale Vision Language Models

Anonymous ACL submission

## Abstract

As the performance of Large-scale Vision Language Models (LVLMs) improves, they are increasingly capable of responding in multiple languages, and there is an expectation that the demand for explanations generated by LVLMs will grow. However, pre-training of Vision Encoder and the integrated training of LLMs with Vision Encoder are mainly conducted using English training data, leaving it uncertain whether LVLMs can completely handle their potential when generating explanations in languages other than English. In addition, multilingual QA benchmarks that create datasets using machine translation have cultural differences and biases, remaining issues for use as evaluation tasks. To address these challenges, this study created an extended dataset in multiple languages without relying on machine translation. This dataset that takes into account nuances and country-specific phrases was then used to evaluate the generation explanation abilities of LVLMs. Furthermore, this study examined whether Instruction-Tuning in resource-rich English improves performance in other languages. Our findings indicate that LVLMs perform worse in languages other than English compared to English. In addition, it was observed that LVLMs struggle to effectively manage the knowledge learned from English data<sup>1</sup>.

## 1 Introduction

Each artwork, e.g., image, has a unique title, making it suitable for evaluating Large-scale Vision Language Models (LVLMs) that handle both the image and the text. Hayashi et al. (2024) focused on artwork explanation generation to investigate the relationship between language-based and vision-based knowledge of LVLMs using English data. When using LVLMs for creative support, explanation generation abilities are required based on the

<sup>1</sup>Our data is publicly available at <https://huggingface.co/anonymous/anonymous-data>



Figure 1: An example of situations that require multilingual and explanation skills.

composition and ingenuity of the image, e.g., comparisons with other works, historical background, and deep artistic knowledge. LVLMs enable image and text aware tasks exactly, e.g., determining the color of traffic lights in the image and judging if it is possible to proceed, by integrating Vision Encoder (Li et al., 2023b), e.g., Vision Transformer (ViT) (Dosovitskiy et al., 2020), which processes image data into high-dimensional features, and Large Language Models (LLMs), which can handle natural language, through additional training. This allows LVLMs to understand instructions with image inputs by humans and generate responses based on those instructions and they have archived remarkable performance on Vision & Language (V&L) benchmarks (Liu et al., 2023c; Li et al., 2023a).

However, there are remaining issues with training current LVLMs when dealing with multilin-

059 gual data. Training and evaluation of LVLMs  
060 often use English data, leaving questions on the  
061 performance on other languages, and there exists  
062 no standard protocol especially when evaluating  
063 the performance of multilingual image understand-  
064 ing tasks. Several multilingual image QA tasks  
065 do exist (Changpinyo et al., 2023; Nguyen et al.,  
066 2023), but they primarily rely on machine transla-  
067 tion, making it uncertain whether country-specific  
068 cultural nuances or biases are completely consid-  
069 ered. Specifically, when creating multilingual QA  
070 tasks, Sakai et al. (2024) pointed out that multiple  
071 concepts e.g., 'roast', 'grill', 'broil', 'toast' and  
072 'bake' in English could be potentially translated  
073 into only one expression e.g., '焼く' in Japanese.  
074 Thus, it is necessary to construct a completely fair  
075 multilingual evaluation dataset for explanation gen-  
076 eration abilities. The issue is, in particular, com-  
077 pounded in the field of art, since an explanation  
078 of an image may vary across countries, leading to  
079 different explanations due to the impression of the  
080 image in other countries. Simply translating from  
081 resource-rich languages like English into other lan-  
082 guages using machine translation to create datasets  
083 fails to account for cultural nuances. For example,  
084 "Mona Lisa" is translated directly into Chinese and  
085 Japanese correctly, but in Spanish, it is translated as  
086 "Mona Lisa" even though it is called "La Gioconda"  
087 in Spanish. Moreover, since these QA datasets do  
088 not evaluate the ability to generate explanations,  
089 there are no appropriate metrics to evaluate the  
090 explanation generation abilities of LVLMs across  
091 different languages.

092 To solve the lack of datasets that can evaluate the  
093 ability to generate explanations in other languages  
094 and the inability to account for country-specific  
095 cultural nuances or biases by simply using machine  
096 translation to create datasets, we created datasets  
097 that allow you to evaluate the ability to generate  
098 explanations in other languages without machine  
099 translation using Wikipedia. Hayashi et al. (2024)  
100 focused only on English, but our study expanded  
101 this work to ten languages (Chinese, Dutch, En-  
102 glish, French, German, Italian, Japanese, Russian,  
103 Spanish, and Swedish).

104 We utilized these datasets to analyze the multilin-  
105 gual performance of current LVLMs in generation  
106 explanation abilities related to artworks with three  
107 settings which are Alignment-10, Alignment-5, and  
108 Full tasks, and investigated whether LVLMs can  
109 maintain equal generation explanation abilities in

110 artworks when extended to ten languages. To inves-  
111 tigate LVLMs' multilingual generation explanation  
112 abilities, we hypothesized that "the integrated train-  
113 ing of LVLMs and the pre-training of Vision En-  
114 coder are mainly trained in English data, limiting  
115 their ability to achieve optimal performance when  
116 handling other languages." Moreover, we also con-  
117 ducted Instruction-Tuning in English-only training  
118 data for two of the models so that validate the extent  
119 to which these two models can acquire explanation  
120 generation capabilities in other languages solely  
121 from English training data.

122 We found that LVLMs perform best when given  
123 instructions in English and generating output in  
124 English, while their performance declines when  
125 instructions or output are in languages other than  
126 English. Moreover, we observed that outputting in  
127 the same language as the instructions like Japanese  
128 instruction with its Japanese response leads to bet-  
129 ter performance than the response in English for  
130 Japanese instruction, indicating that LVLMs strug-  
131 gle to effectively utilize the knowledge learned in  
132 English when applied to other languages. The  
133 result also showed that performance was further  
134 worse with Instruction-Tuning conducted in En-  
135 glish. These findings support our hypothesis and  
136 suggests that it is necessary to let Vision Encoder  
137 train not only English training data but also other  
138 language data.

## 139 2 Related Work

140 **LVLMs** In general, an LVLM comprises a Vision  
141 Encoder that processes visual information and an  
142 LLM pre-trained on a large amount of textual data.  
143 They are trained using contrastive learning (Chen  
144 et al., 2020), aiming to integrate visual and linguis-  
145 tic information. Vision Encoder is a model trained  
146 to encode images and visual data, typically using  
147 architectures such as ResNet (He et al., 2015) or Vi-  
148 sion Transformer (ViT) (Dosovitskiy et al., 2020).  
149 On the other hand, LLMs are models pre-trained  
150 on a large text dataset, with prominent examples  
151 including Qwen (Bai et al., 2023a; Yang et al.,  
152 2024), LLaMA (Touvron et al., 2023a,b; Dubey  
153 et al., 2024), Gemini (Team et al., 2023; Reid et al.,  
154 2024) and GPT (Brown et al., 2020; Ouyang et al.,  
155 2022; Achiam et al., 2023). LVLMs such as Qwen-  
156 VL (Bai et al., 2023b) and LLaVA-NeXT (Liu et al.,  
157 2024) are examples of integrated models. These  
158 models achieve visual and natural language integra-  
159 tion by acquiring features from images through the

Language	Type	Template	Instruction	Output
English	Section	Explain the <b>{Section}</b> of this artwork, <b>{Title}</b> .	Explain the <b>History</b> of this artwork, <b>Mona Lisa</b> .	Of Leonardo da Vinci's works, the Mona Lisa is the only portrait whose authenticity...
	Subsection	Explain the <b>{Subsection}</b> regarding the <b>{Section}</b> of this artwork, <b>{Title}</b> .	Explain the <b>Creation and date</b> regarding the <b>History</b> of this artwork, <b>Mona Lisa</b> .	The record of an October 1517 visit by Luis d'Aragon states that the Mona Lisa...
Japanese	Section	<b>{Title}</b> の作品に関して、この作品の <b>{Section}</b> を説明してください。	<b>モナリザ</b> の作品に関して、この作品の <b>歴史</b> について説明してください。	レオナルド・ダ・ヴィンチの作品の中で、「モナ・リザ」は唯一、その真偽が不確かな肖像画であり...
	Subsection	<b>{Title}</b> の作品に関して、この作品の <b>{Section}</b> に関する <b>{Subsection}</b> を説明してください。	<b>モナリザ</b> の作品に関して、この作品の <b>歴史</b> に関する <b>制作と日付</b> を説明してください。	1517年10月のルイ・ド・アラゴンの訪問の記録には、「モナ・リザ」について...
Chinese	Section	解释这件艺术品的 <b>{Section}</b> ， <b>{Title}</b> 。	解释这件艺术品的 <b>历史</b> ， <b>蒙娜丽莎</b> 。	在达芬奇的作品中，蒙娜丽莎是唯一一幅真伪有争议的肖像画...
	Subsection	解释关于这件艺术品的 <b>{Section}</b> 的 <b>{Subsection}</b> 的 <b>{Section}</b> ， <b>{Title}</b> 。	解释关于这件艺术品的 <b>历史的创作和日期</b> ， <b>蒙娜丽莎</b> 。	路易·德·阿拉贡在1517年10月访问的记录中提到，“蒙娜丽莎”...
Spanish	Section	Explica la <b>{Section}</b> de esta obra de arte, <b>{Title}</b> .	Explica la <b>Historia</b> de esta obra de arte, <b>Mona Lisa</b> .	De las obras de Leonardo da Vinci, la Mona Lisa es el único retrato cuya autenticidad...
	Subsection	Explica la <b>{Subsection}</b> sobre la <b>{Section}</b> de esta obra de arte, <b>{Title}</b> .	Explica la <b>Creación y fecha</b> sobre la <b>Historia</b> de esta obra de arte, <b>Mona Lisa</b> .	El registro de una visita en octubre de 1517 de Luis de Aragón menciona que la Mona Lisa...

Table 1: Examples of templates and instructions for the proposed task. The blue part indicates the artwork’s title and the red part indicates the names of sections and subsections in the original Wikipedia articles that correspond to their explanations. We prepared such templates for ten languages and asked native speakers to make sure they are on the same level as English.

Vision Encoder and textual features through LLMs and then performing additional training with the goal of integrating vision and language.

**LVLMS & Knowledge** Whether the visual knowledge learned by the Vision Encoder and the linguistic knowledge learned by LLMs are properly aligned remains mostly unclear (Li et al., 2022, 2023b). Especially for generating explanations involving knowledge about artwork, which this study focuses on, it is essential to systematically align and utilize both types of knowledge (Hayashi et al., 2024). This requires the integration of visual knowledge (e.g., visual features of specific artworks) and linguistic knowledge (e.g., historical background and technical details about those artworks). In LVLMS, the integration of Vision Encoder and LLMs are achieved by adding partial networks, but this alone makes it challenging to properly align visual and linguistic knowledge. In domains requiring sophisticated knowledge, such as artwork, improper alignment can degrade the quality of generated explanations. Thus, while this study aims to integrate visual and linguistic information and build efficient models using contrastive learning, it also indicates that further research is necessary to achieve proper alignment of visual and linguistic knowledge.

**LVLMS & Multilingual** As we mentioned earlier, LVLMS follow human instructions through integrated learning of Vision Encoder (Li et al., 2023b) and LLMs trained by a large amount of English training data. However, it is unclear whether

LVLMS are able to really understand and output properly when input from languages other than English. On the other hand, as far as evaluation tasks such as XGQA (Pfeiffer et al., 2022) they expanded the English GQA dataset into seven languages through translation. However, because this expansion relies on translations from English, it likely includes QA pairs that do not consider the cultural contexts of the target languages. For instance, MaXM (Changpinoy et al., 2023) collects large data sets by translating non-English language data into English, which is then back-translated into seven languages. Similarly, EVJVQA (Nguyen et al., 2023) creates around 33,000 QA pairs from approximately 5,000 images taken in Vietnam, but the translations still retain biases unique to Vietnamese culture and norms. In our research, we mitigated these biases by focusing on artworks, preventing the introduction of a specific culture to any country within the images. (i.e., There are countries where cars drive on the right lane and others where they drive on the left.) Since artworks have unique and definitive relationships between the title and its image, we also create datasets from relatively resource-rich Wikipedia in various languages without relying on machine translation. Our study is not a Question Answering task, such as VQA (Antol et al., 2015), but an explanation task, which requires LVLMS to explain images correctly. We evaluated an explanation-generation task in ten languages expanding Hayashi et al. (2024) work.

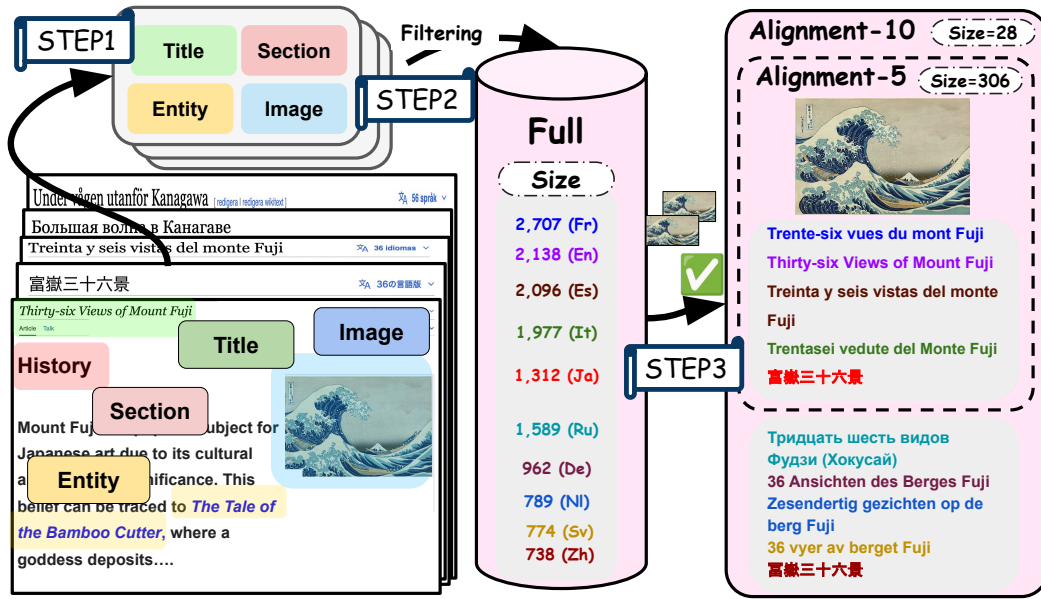


Figure 2: How to make datasets from Wikipedia. As shown in Section 4, we extracted and filtered Wikipedia pages about artworks. We then manually identified pages with titles and images common across ten languages.

### 3 Task

Our task is “Analyzing the multilingual performance of LVLMs in explaining artworks”. To tackle this, we measured explanatory capabilities using three settings (Alignment-10, Alignment-5, Full) which was described below. In addition, we prepared templates for input when evaluating explanation abilities with datasets we created.

**Alignment-10** We created datasets composed only of data with the same images and titles across ten languages from the Full task, which total number is 28. In other words, this dataset contains the same images with titles represented in the language of each country, allowing for an equal evaluation of description generation capabilities across the ten languages.

**Alignment-5** To mitigate the data scarcity issue in Alignment-10, Alignment-5 restricts the target languages to five specific languages. The total number of data is 306, and this dataset is used to compare explanation generation abilities across the five languages. To cover a diverse range of language families, we selected English, Spanish, French, Italian, and Japanese as Alignment-5 task.

**Full** To further mitigate the data scarcity issues in the above settings, Full ignores the correspondence of artworks between languages and treats each language independently. For details on the number of data, refer to Table 7 or Figure 4. By using the

Full task, we aimed to evaluate the differences in performance.

**Templates** We prepared templates for evaluating explanation generation abilities using the datasets created from three tasks mentioned above. The process is as follows: 1) We prepared four patterns of templates for each of ten languages. In templates, we referred to the study by Hayashi et al. (2024), selecting four patterns with clearly different grammatical structures to avoid a lack of diversity. Sakai et al. (2024) noted that not choosing distinctly different patterns may result in differences originally present in English being lost in translation; 2) We let ChatGPT<sup>2</sup> translate the obtained templates into ten languages. We chose to use LLMs rather than translation tools because LLMs are thought to better understand and translate including nuances; 3) Even with translations taking into nuances by ChatGPT, there may be variations in quality between languages. To solve this, we asked nine native speakers of ten languages, to check whether the templates translated back into English maintained the same nuance and level of difficulty. This process ensured that all 10 language templates created in this study have the same level of difficulty; Of course, it might be possible to crowdsource this task using platforms like MTurk<sup>3</sup>, but asking annotators simply “Is this translation correct including

<sup>2</sup><https://openai.com/chatgpt/>

<sup>3</sup><https://www.mturk.com/>

Input	Output	LVLM	Entity Cov.		Entity F1	Entity Cooccurrence				Avg. Length
			exact	partial		n=0	n=1	n=2	n=∞	
En	En	LLaVA-NeXT	<b>26.49</b>	<b>31.54</b>	<b>26.07</b>	1.35	1.65	1.66	1.70	252
		LLaVA-NeXT (FT)	16.98	22.70	19.95	3.02	3.33	3.23	3.11	83
		Qwen-VL	21.11	27.03	21.78	1.60	1.59	1.56	1.52	155
		Qwen-VL (FT)	21.12	24.87	21.95	<b>3.57</b>	<b>3.83</b>	<b>3.78</b>	<b>3.68</b>	177
		mPLUG-Owl2	12.79	17.08	13.48	2.07	1.68	1.59	1.56	151

Table 2: Results of LVLMs in Alignment-10 Task (the instruction and the output in English, {En}-{En}). Bold fonts indicate the best scores. The red and blue figures shown in the following figures are the different figures compared to this Table. "(FT)" indicates the model conducted LoRA-Tuning.

nuances?" may not lead to serious engagement with the translation checking task. For examples of the each language template, refer to Table 1.

## 4 Dataset Creation

For each of ten languages, the following steps were taken to create the dataset. Ten languages were determined based on having a higher number of Wikipedia articles than the total number of articles.

**STEP1 Extracting Data from Wikipedia** We collected Artwork articles from the English Wikipedia Infobox. Articles with the same title in nine other languages are identified to create corresponding articles in those languages. Hyperlinked strings within the articles are extracted as entities related to artworks. The description includes four types of information: the image, the title, hierarchical information from the article (Section, Subsection, Sub subsection), and the extracted entities.

**STEP2 Filtering and Formatting** From the collected articles, those without images were excluded. Any articles that had domains but no actual pages on Wikipedia were also removed. This process completes the dataset used for the Full task.

**STEP3 Adjusting** For the Alignment-10 task and Alignment-5 task, to achieve alignment across languages, only articles with the same English-translated titles are selected for both datasets. To eliminate differences between languages, a manual verification is conducted to ensure that all articles contain images of the same artwork. Variations in image size are permitted, but all images must represent the same artwork across languages. The datasets for Alignment-10 and Alignment-5 are prepared accordingly, using images from the English articles for alignment.

**STEP4 Data Splitting** To measure the explanation generation abilities of LVLMs, the following

approach is used: (1) For the Alignment task, all data was treated as test set. (2) For the Full task, nine non-English languages are used for test set, while English data is divided into train, dev, and test sets. To avoid biases arising from the popularity of artworks in the LVLM’s training data, we shuffled the English data based on six indicators: page views, number of links, number of edits, number of references, number of language versions, and article length (Hayashi et al., 2024). The data was ranked according to these indicators, and the test, valid, and train data were split in a 2:2:6 ratio to maintain average rankings. The data used in the Alignment task was included in the test set.

## 5 Experiments

### 5.1 Evaluation Metrics

This study adopted three evaluation metrics proposed by Hayashi et al. (2024) and also described these metrics more details in Appendix E. We also utilize popular metrics in NLG for evaluation, i.e., BLEU (Papineni et al., 2002), ROUGE (Lin, 2004), and BERTScore (Zhang et al., 2019).

**Entity Coverage, Entity F1, and Entity Cooccurrence** These metrics evaluate how well the generated text incorporates entities related to the artwork and how accurately it reflects the relationships between these entities proposed by Hayashi et al. (2024). Entity Coverage measures the inclusion of relevant entities in both exact and partial matches. Entity F1 assesses the frequency and appropriateness of entity usage by comparing the generated text with reference explanations, inspired by the ROUGE metric. Entity Cooccurrence goes a step further by examining how entities are contextually combined across sentences, considering their co-occurrence within the entire text, and applying brevity penalties to avoid inflated coverage in longer explanations.

Input	Output	LVLML	Entity Cov.		Entity F1	Entity Cooccurrence				Avg. Length
			exact	partial		n=0	n=1	n=2	n=∞	
En	Es	LLaVA-NeXT	0.00 (-26.49)	2.24 (-29.30)	0.00 (-26.07)	0.00 (-1.35)	0.00 (-1.65)	0.00 (-1.66)	0.00 (-1.70)	137
		LLaVA-NeXT (FT)	6.23 (-10.75)	9.05 (-13.65)	6.87 (-13.08)	1.27 (-1.75)	1.21 (-2.12)	1.09 (-2.13)	1.06 (-2.05)	83
		Qwen-VL	<b>10.81 (-10.29)</b>	<b>15.18 (-11.85)</b>	<b>11.42 (-10.36)</b>	<b>1.48 (-0.12)</b>	<b>1.41 (-0.18)</b>	<b>1.34 (-0.22)</b>	<b>1.27 (-0.25)</b>	109
		Qwen-VL (FT)	4.25 (-16.87)	7.86 (-17.01)	5.40 (-16.55)	0.36 (-3.21)	0.36 (-3.47)	0.31 (-3.46)	0.29 (-3.39)	190
		mPLUG-Owl2	8.95 (-3.83)	11.95 (-5.13)	9.62 (-3.86)	0.93 (-1.14)	1.13 (-0.55)	1.07 (-0.52)	1.02 (-0.54)	108
En	Fr	LLaVA-NeXT	1.00 (-25.49)	7.42 (-24.12)	1.33 (-24.74)	0.00 (-1.35)	0.00 (-1.65)	0.00 (-1.66)	0.00 (-1.70)	179
		LLaVA-NeXT (FT)	8.39 (-8.59)	11.41 (-11.29)	8.71 (-11.24)	1.43 (-1.59)	<b>1.99 (-1.34)</b>	<b>1.96 (-1.27)</b>	<b>1.95 (-1.16)</b>	68
		Qwen-VL	<b>12.11 (-9.00)</b>	<b>17.23 (-9.80)</b>	<b>13.05 (-8.73)</b>	1.44 (-0.16)	1.45 (-0.14)	1.40 (-0.16)	1.34 (-0.18)	96
		Qwen-VL (FT)	7.19 (-13.92)	11.24 (-13.63)	8.34 (-13.61)	0.45 (-3.12)	0.90 (-2.93)	0.88 (-2.89)	0.89 (-2.79)	175
		mPLUG-Owl2	10.26 (-2.53)	15.51 (-1.57)	10.99 (-2.49)	<b>1.72 (-0.35)</b>	1.33 (-0.35)	1.20 (-0.39)	1.16 (-0.40)	109
En	De	LLaVA-NeXT	<b>14.03 (-12.46)</b>	<b>17.90 (-13.64)</b>	<b>16.51 (-9.56)</b>	<b>1.73 (+0.38)</b>	<b>1.70 (+0.05)</b>	<b>1.67 (+0.01)</b>	<b>1.82 (+0.13)</b>	169
		LLaVA-NeXT (FT)	6.83 (-10.15)	9.54 (-13.16)	8.23 (-11.72)	0.86 (-2.15)	0.74 (-2.59)	0.77 (-2.46)	0.78 (-2.33)	82
		Qwen-VL	10.64 (-10.46)	13.95 (-13.08)	13.21 (-8.56)	1.16 (-0.44)	1.24 (-0.35)	1.21 (-0.35)	1.40 (-0.12)	111
		Qwen-VL (FT)	7.98 (-13.14)	11.08 (-13.79)	9.86 (-12.09)	0.80 (-2.77)	0.65 (-3.18)	0.74 (-3.03)	0.74 (-2.94)	203
		mPLUG-Owl2	8.81 (-3.98)	12.12 (-4.97)	10.54 (-2.94)	0.72 (-1.35)	0.76 (-0.92)	0.74 (-0.85)	0.70 (-0.86)	98
En	It	LLaVA-NeXT	<b>8.53 (-17.95)</b>	<b>13.33 (-18.21)</b>	<b>9.37 (-16.70)</b>	0.86 (-0.48)	0.87 (-0.79)	1.06 (-0.60)	1.05 (-0.65)	171
		LLaVA-NeXT (FT)	5.89 (-11.09)	8.90 (-13.80)	6.61 (-13.34)	0.96 (-2.06)	<b>1.32 (-2.01)</b>	<b>1.32 (-1.91)</b>	<b>1.31 (-1.80)</b>	66
		Qwen-VL	7.23 (-13.87)	11.43 (-15.59)	8.71 (-13.06)	0.51 (-1.08)	0.62 (-0.97)	0.65 (-0.91)	0.63 (-0.89)	107
		Qwen-VL (FT)	5.51 (-15.61)	8.17 (-16.70)	6.53 (-15.42)	<b>1.14 (-2.44)</b>	0.82 (-3.01)	0.85 (-2.93)	0.84 (-2.84)	170
		mPLUG-Owl2	3.97 (-8.82)	8.50 (-8.58)	4.50 (-8.98)	0.15 (-1.92)	0.14 (-1.53)	0.16 (-1.43)	0.15 (-1.41)	107
En	Nl	LLaVA-NeXT	12.21 (-14.28)	<b>17.83 (-13.71)</b>	14.60 (-11.46)	0.36 (-0.99)	1.81 (+0.15)	1.70 (+0.04)	<b>1.83 (+0.13)</b>	178
		LLaVA-NeXT (FT)	9.41 (-7.56)	15.01 (-7.69)	12.14 (-7.81)	<b>1.21 (-1.81)</b>	1.07 (-2.27)	0.91 (-2.32)	1.02 (-2.09)	119
		Qwen-VL	11.07 (-10.04)	16.44 (-10.59)	12.73 (-9.05)	0.89 (-0.71)	<b>1.90 (+0.32)</b>	<b>1.78 (+0.22)</b>	1.80 (+0.28)	132
		Qwen-VL (FT)	<b>12.67 (-8.45)</b>	17.03 (-7.84)	<b>16.91 (-5.04)</b>	1.02 (-2.55)	0.96 (-2.88)	0.95 (-2.83)	1.01 (-2.67)	181
		mPLUG-Owl2	8.27 (-4.51)	13.46 (-3.62)	9.06 (-4.42)	0.46 (-1.61)	0.43 (-1.25)	0.41 (-1.18)	0.41 (-1.14)	100
En	Sv	LLaVA-NeXT	<b>15.01 (-11.48)</b>	<b>18.65 (-12.89)</b>	<b>13.56 (-12.51)</b>	<b>1.29 (-0.05)</b>	0.97 (-0.69)	<b>1.15 (-0.51)</b>	<b>1.09 (-0.61)</b>	174
		LLaVA-NeXT (FT)	10.00 (-6.97)	12.43 (-10.27)	10.54 (-9.41)	0.84 (-2.17)	<b>1.08 (-2.26)</b>	0.97 (-2.26)	0.87 (-2.24)	115
		Qwen-VL	10.37 (-10.74)	14.08 (-12.94)	10.15 (-11.62)	0.84 (-0.76)	0.86 (-0.72)	0.83 (-0.73)	0.80 (-0.72)	123
		Qwen-VL (FT)	8.97 (-12.14)	12.25 (-12.61)	9.66 (-12.29)	0.87 (-2.70)	0.94 (-2.89)	0.92 (-2.86)	0.90 (-2.78)	164
		mPLUG-Owl2	10.21 (-2.57)	13.03 (-4.05)	9.07 (-4.41)	0.35 (-1.72)	0.35 (-1.33)	0.34 (-1.25)	0.34 (-1.22)	88
En	Ru	LLaVA-NeXT	<b>10.32 (-16.17)</b>	<b>15.15 (-16.39)</b>	<b>8.53 (-17.54)</b>	<b>0.32 (-1.02)</b>	<b>0.36 (-1.30)</b>	<b>0.31 (-1.35)</b>	<b>0.32 (-1.38)</b>	203
		LLaVA-NeXT (FT)	0.55 (-16.42)	1.87 (-20.83)	0.49 (-19.46)	0.00 (-3.02)	0.02 (-3.32)	0.02 (-3.21)	0.01 (-3.10)	85
		Qwen-VL	4.59 (-16.52)	8.05 (-18.97)	3.51 (-18.26)	0.02 (-1.58)	0.07 (-1.52)	0.07 (-1.49)	0.07 (-1.45)	113
		Qwen-VL (FT)	0.00 (-21.12)	0.95 (-23.91)	0.00 (-21.95)	0.00 (-3.57)	0.00 (-3.83)	0.00 (-3.78)	0.00 (-3.68)	169
		mPLUG-Owl2	5.99 (-6.80)	8.68 (-8.40)	4.88 (-8.60)	0.00 (-2.07)	0.02 (-1.66)	0.01 (-1.57)	0.01 (-1.54)	99
En	Ja	LLaVA-NeXT	<b>8.68 (-17.81)</b>	<b>8.68 (-22.86)</b>	<b>11.47 (-14.60)</b>	<b>0.80 (-0.54)</b>	<b>0.80 (-0.85)</b>	<b>0.80 (-0.86)</b>	<b>0.80 (-0.90)</b>	211
		LLaVA-NeXT (FT)	0.29 (-16.68)	0.30 (-22.40)	0.38 (-19.57)	0.04 (-2.98)	0.04 (-3.29)	0.04 (-3.19)	0.04 (-3.07)	85
		Qwen-VL	3.52 (-17.59)	3.53 (-23.49)	4.78 (-17.00)	0.32 (-1.28)	0.32 (-1.27)	0.32 (-1.24)	0.32 (-1.20)	132
		Qwen-VL (FT)	0.00 (-21.12)	0.03 (-24.84)	0.00 (-21.95)	0.00 (-3.57)	0.00 (-3.83)	0.00 (-3.78)	0.00 (-3.68)	188
		mPLUG-Owl2	3.75 (-9.04)	3.75 (-13.33)	4.98 (-8.49)	0.39 (-1.68)	0.39 (-1.28)	0.39 (-1.20)	0.39 (-1.17)	112
En	Zh	LLaVA-NeXT	<b>14.00 (-1.86)</b>	<b>14.09 (-6.86)</b>	<b>16.69 (+0.19)</b>	0.66 (-0.42)	0.66 (-0.58)	0.66 (-0.56)	0.66 (-0.59)	228
		LLaVA-NeXT (FT)	0.14 (-11.49)	0.39 (-15.08)	0.15 (-13.97)	0.00 (-2.42)	0.00 (-2.60)	0.00 (-2.51)	0.00 (-2.43)	92
		Qwen-VL	10.69 (-1.45)	10.70 (-5.71)	12.71 (+0.52)	<b>0.74 (-0.59)</b>	<b>0.73 (-0.44)</b>	<b>0.73 (-0.39)</b>	<b>0.73 (-0.35)</b>	138
		Qwen-VL (FT)	0.37 (-13.40)	0.75 (-16.88)	0.51 (-12.91)	0.01 (-2.96)	0.01 (-3.09)	0.01 (-3.04)	0.01 (-2.99)	154
		mPLUG-Owl2	6.38 (-6.45)	6.40 (-10.74)	7.75 (-5.77)	0.32 (-1.75)	0.32 (-1.36)	0.32 (-1.27)	0.32 (-1.24)	108

Table 3: Results of LVLMLs in Alignment-10 Task ({En}-{Lang}). Bold fonts indicate the best score for that language combination. The values are noted next to the output of the difference by the same model in the method with instruction and output in English ({En}-{En}). Red indicates a higher value than that method; blue indicates a lower value.

## 5.2 Models and Others

We chose five models with relatively high performance: mPLUG-Owl2 (Ye et al., 2024), LLaVA-NeXT (Liu et al., 2023b, 2024, 2023a), XComposer2 (Dong et al., 2024), Phi-3 (Abdin et al., 2024), and Qwen-VL (Bai et al., 2023b). In addition, LLaVA-NeXT and Qwen-VL were conducted LoRA Tuning (Hu et al., 2022) with English train data and included in the evaluation. Detailed experimental settings are described in Appendix A.1. This approach is based on the observation that current LLMs perform better when instructions are given in English (Putri et al., 2024). As far as Alignment tasks, we validated four patterns of input: {En, Lang}-{En, Lang}. This indicates that when the input is English, the output can be di-

rected to English or another language. The same thing can also be done when the input is another language, and these four patterns were tested in this study. By testing these patterns, we verify whether or not LVLMLs perform better when supported in English, and whether or not having the output in English is a meaningful instruction. As far as tokenizing words, we used SpaCy<sup>4</sup> as a multilingual tokenizer, tokenizing each language to perform segmentation. Thus, each language is expected to be divided into optimal token units.

## 5.3 Results

From the experiments conducted with Alignment-10, the method let LVLMLs generate in English with

<sup>4</sup><https://spacy.io/>

Input	Output	LVLML	Entity Cov.		Entity F1	Entity Cooccurrence				Avg. Length
			exact	partial		n=0	n=1	n=2	n=∞	
Es	Es	LLaVA-NeXT	<b>17.26 (-9.23)</b>	<b>21.30 (-10.25)</b>	<b>17.05 (-9.01)</b>	2.13 (+0.78)	2.32 (+0.67)	2.17 (+0.51)	2.10 (+0.40)	186
		LLaVA-NeXT (FT)	12.82 (-4.15)	16.84 (-5.86)	12.77 (-7.18)	0.88 (-2.14)	1.03 (-2.31)	1.26 (-1.97)	1.11 (-1.99)	147
		Qwen-VL	14.68 (-6.43)	18.43 (-8.59)	14.35 (-7.43)	2.40 (+0.81)	<b>2.42 (+0.83)</b>	<b>2.57 (+1.01)</b>	<b>2.56 (+1.04)</b>	150
		Qwen-VL (FT)	4.09 (-17.02)	7.10 (-17.77)	4.52 (-17.42)	0.15 (-3.43)	0.16 (-3.68)	0.15 (-3.63)	0.14 (-3.54)	301
		mPLUG-Owl2	10.91 (-1.87)	15.06 (-2.02)	11.91 (-1.57)	<b>2.47 (+0.40)</b>	2.07 (+0.40)	2.02 (+0.44)	1.99 (+0.44)	135
Fr	Fr	LLaVA-NeXT	<b>24.35 (-2.14)</b>	<b>29.27 (-2.27)</b>	24.38 (-1.69)	0.95 (-0.40)	0.90 (-0.75)	0.88 (-0.78)	0.90 (-0.80)	211
		LLaVA-NeXT (FT)	16.63 (-0.35)	20.13 (-2.57)	16.09 (-3.86)	1.18 (-1.83)	0.93 (-2.41)	1.00 (-2.23)	0.98 (-2.13)	98
		Qwen-VL	19.38 (-1.73)	24.71 (-2.32)	18.30 (-3.47)	1.07 (-0.53)	1.03 (-0.55)	0.96 (-0.60)	0.96 (-0.56)	165
		Qwen-VL (FT)	24.15 (+3.04)	28.59 (+3.73)	<b>24.79 (+2.85)</b>	<b>3.83 (+0.26)</b>	<b>4.41 (+0.58)</b>	<b>4.51 (+0.73)</b>	<b>4.51 (+0.83)</b>	219
		mPLUG-Owl2	17.43 (+4.64)	22.48 (+5.40)	17.78 (+4.30)	0.85 (-1.22)	0.65 (-1.02)	0.75 (-0.84)	0.73 (-0.83)	158
De	De	LLaVA-NeXT	<b>17.45 (-9.04)</b>	<b>20.66 (-10.89)</b>	<b>21.05 (-5.02)</b>	2.11 (+0.77)	<b>2.20 (+0.55)</b>	<b>2.22 (+0.56)</b>	<b>2.11 (+0.41)</b>	204
		LLaVA-NeXT (FT)	10.53 (-6.44)	13.10 (-9.60)	13.32 (-6.63)	1.53 (-1.49)	1.09 (-2.25)	1.15 (-2.07)	1.16 (-1.95)	123
		Qwen-VL	15.10 (-6.00)	18.20 (-8.82)	17.97 (-3.81)	<b>2.12 (+0.52)</b>	1.99 (+0.41)	2.08 (+0.52)	1.99 (+0.47)	160
		Qwen-VL (FT)	7.74 (-13.38)	9.58 (-15.28)	9.23 (-12.72)	0.37 (-3.20)	0.40 (-3.43)	0.43 (-3.34)	0.40 (-3.28)	287
		mPLUG-Owl2	14.33 (+1.55)	17.63 (+0.55)	16.73 (+3.25)	1.99 (-0.08)	1.92 (+0.25)	1.94 (+0.35)	1.81 (+0.25)	143
It	It	LLaVA-NeXT	<b>10.34 (-16.14)</b>	<b>15.43 (-16.11)</b>	<b>11.33 (-14.74)</b>	1.16 (-0.19)	<b>0.93 (-0.72)</b>	<b>0.96 (-0.70)</b>	<b>0.96 (-0.74)</b>	185
		LLaVA-NeXT (FT)	5.73 (-11.25)	9.84 (-12.86)	6.45 (-13.50)	0.31 (-2.71)	0.25 (-3.08)	0.25 (-2.98)	0.23 (-2.88)	91
		Qwen-VL	9.97 (-11.13)	14.20 (-12.82)	11.09 (-10.68)	<b>1.16 (-0.44)</b>	0.93 (-0.65)	0.94 (-0.62)	0.90 (-0.62)	126
		Qwen-VL (FT)	3.15 (-17.96)	6.95 (-17.92)	3.42 (-18.53)	0.15 (-3.42)	0.18 (-3.65)	0.23 (-3.54)	0.21 (-3.47)	253
		mPLUG-Owl2	8.69 (-4.10)	12.66 (-4.42)	9.54 (-3.94)	0.51 (-1.56)	0.32 (-1.36)	0.35 (-1.24)	0.33 (-1.23)	111
Nl	Nl	LLaVA-NeXT	17.66 (-8.83)	23.56 (-7.99)	<b>19.78 (-6.28)</b>	0.79 (-0.56)	<b>3.55 (+1.89)</b>	3.61 (+1.95)	3.88 (+2.18)	199
		LLaVA-NeXT (FT)	15.57 (-1.40)	20.79 (-1.91)	16.87 (-3.08)	1.66 (-1.35)	3.38 (+0.05)	3.32 (+0.09)	3.47 (+0.37)	183
		Qwen-VL	<b>19.41 (-1.69)</b>	<b>24.45 (-2.58)</b>	19.65 (-2.13)	<b>2.13 (+0.53)</b>	3.27 (+1.69)	<b>3.89 (+2.33)</b>	<b>4.04 (+2.52)</b>	172
		Qwen-VL (FT)	12.68 (-8.43)	18.46 (-6.41)	16.72 (-5.22)	1.09 (-2.48)	1.66 (-2.18)	1.81 (-1.96)	1.80 (-1.88)	300
		mPLUG-Owl2	10.78 (-2.01)	15.43 (-1.66)	12.81 (-0.67)	0.15 (-1.92)	1.08 (-0.60)	1.05 (-0.54)	1.12 (-0.43)	114
Sv	Sv	LLaVA-NeXT	<b>27.51 (+1.02)</b>	<b>29.61 (-1.93)</b>	16.71 (-9.36)	2.10 (+0.75)	0.87 (-0.78)	0.89 (-0.77)	0.90 (-0.79)	206
		LLaVA-NeXT (FT)	22.83 (+5.86)	25.10 (+2.40)	12.17 (-7.78)	2.82 (-0.20)	1.11 (-2.22)	1.17 (-2.06)	1.16 (-1.94)	169
		Qwen-VL	24.02 (+2.92)	26.69 (-0.34)	<b>19.18 (-2.60)</b>	<b>3.60 (+2.00)</b>	<b>1.53 (-0.06)</b>	<b>1.54 (-0.02)</b>	<b>1.50 (-0.02)</b>	147
		Qwen-VL (FT)	16.04 (-5.07)	18.10 (-6.77)	6.15 (-15.80)	0.23 (-3.35)	0.18 (-3.65)	0.20 (-3.57)	0.21 (-3.47)	242
		mPLUG-Owl2	21.40 (+8.61)	23.51 (+6.43)	13.84 (+0.36)	2.01 (-0.06)	1.07 (-0.61)	1.06 (-0.52)	1.05 (-0.51)	111
Ru	Ru	LLaVA-NeXT	<b>14.38 (-12.11)</b>	<b>17.43 (-14.11)</b>	<b>9.81 (-16.26)</b>	0.26 (-1.08)	<b>0.45 (-1.20)</b>	<b>0.42 (-1.24)</b>	<b>0.41 (-1.29)</b>	219
		LLaVA-NeXT (FT)	10.74 (-6.24)	13.67 (-9.03)	6.55 (-13.40)	0.32 (-2.70)	0.37 (-2.96)	0.36 (-2.87)	0.36 (-2.75)	184
		Qwen-VL	6.80 (-14.31)	9.68 (-17.34)	4.63 (-17.15)	0.31 (-1.29)	0.32 (-1.27)	0.30 (-1.26)	0.31 (-1.21)	170
		Qwen-VL (FT)	1.76 (-19.35)	3.60 (-21.27)	1.52 (-20.42)	0.14 (-3.43)	0.14 (-3.69)	0.14 (-3.64)	0.14 (-3.54)	324
		mPLUG-Owl2	7.07 (-5.72)	8.92 (-8.16)	5.57 (-7.91)	<b>0.51 (-1.56)</b>	0.34 (-1.33)	0.31 (-1.28)	0.35 (-1.21)	129
Ja	Ja	LLaVA-NeXT	<b>13.38 (-13.11)</b>	<b>13.38 (-18.17)</b>	<b>17.68 (-8.39)</b>	0.73 (-0.61)	0.83 (-0.83)	0.83 (-0.83)	0.83 (-0.87)	249
		LLaVA-NeXT (FT)	7.51 (-9.46)	7.51 (-15.19)	7.80 (-12.15)	1.14 (-1.88)	1.14 (-2.19)	1.14 (-2.09)	1.14 (-1.97)	167
		Qwen-VL	10.89 (-10.22)	10.90 (-16.13)	14.56 (-7.22)	0.92 (-0.68)	0.92 (-0.67)	0.92 (-0.64)	0.92 (-0.60)	154
		Qwen-VL (FT)	0.86 (-20.26)	0.88 (-23.99)	1.12 (-20.83)	0.03 (-3.55)	0.03 (-3.81)	0.03 (-3.75)	0.03 (-3.65)	278
		mPLUG-Owl2	6.91 (-5.88)	6.93 (-10.15)	9.34 (-4.14)	<b>1.20 (-0.87)</b>	<b>1.21 (-0.46)</b>	<b>1.21 (-0.38)</b>	<b>1.21 (-0.35)</b>	144
Zh	Zh	LLaVA-NeXT	13.78 (-2.08)	13.78 (-7.17)	17.00 (+0.50)	0.54 (-0.54)	0.53 (-0.70)	0.53 (-0.69)	0.53 (-0.72)	246
		LLaVA-NeXT (FT)	6.93 (-4.71)	6.97 (-8.50)	7.31 (-6.81)	0.78 (-1.64)	0.78 (-1.83)	0.78 (-1.73)	0.78 (-1.65)	170
		Qwen-VL	<b>17.90 (+5.76)</b>	<b>17.90 (+1.48)</b>	<b>22.12 (+9.93)</b>	<b>3.31 (+1.97)</b>	<b>3.30 (+2.13)</b>	<b>3.30 (+2.13)</b>	<b>3.30 (+2.22)</b>	155
		Qwen-VL (FT)	0.22 (-13.55)	0.33 (-17.29)	0.27 (-13.16)	0.00 (-2.97)	0.00 (-3.10)	0.00 (-3.06)	0.00 (-3.00)	249
		mPLUG-Owl2	9.03 (-3.80)	9.05 (-8.08)	12.98 (-0.55)	0.77 (-1.31)	0.77 (-0.91)	0.77 (-0.82)	0.77 (-0.80)	150

Table 4: Results of LVLMLs in Alignment-10 Task (the format with instruction and output in each of the ten languages, {Lang}-{Lang}). Bold fonts indicate the best score for that language combination. The values are noted next to the differences output by the same model in the format with instruction and output in English ({En}-{En}). Red indicates a higher value than {En}-{En}; blue indicates a lower value.

English ({En}-{En}) results are listed in Table 2, the method which is instruction in English and output in other languages ({En}-{Lang}) results in Table 3, and the instruction and output in other same languages ({Lang}-{Lang}) results in Table 4. The results for Phi-3 and XComposer2 are described in the Table 10 in Appendix. Overall, the results confirm that giving instructions in English and letting them generate output in English (i.e., {En}-{En}) maximizes the performance of LVLMLs. On the other hand, LoRA Tuning increased the value of Entity Cooccurrence, while other values decreased. This suggests that LoRA Tuning enabled LVLMLs to understand and explain the context, but prevented entities from appearing in the generated sentences. Furthermore, looking at the results of Alignment-5 in Table 9 in Appendix, where the number of

data was expanded, the outputs that used English instructions and outputs were generally higher, followed by those using instructions and outputs in other languages. This is consistent with the results of Alignment-10. In addition, Figure 3 includes results where instructions were given in other languages and outputs were produced in English.

## 6 Analysis and Discussion

**Which Instruction and Output Language is Best?** We confirmed that the pattern which instruction and output are English ({En}-{En}) performed the best ability, whereas the performance is lower for the pattern in which instruction in English and output in other languages ({En}-{Lang}, i.e., Please generate the output in Chinese). This

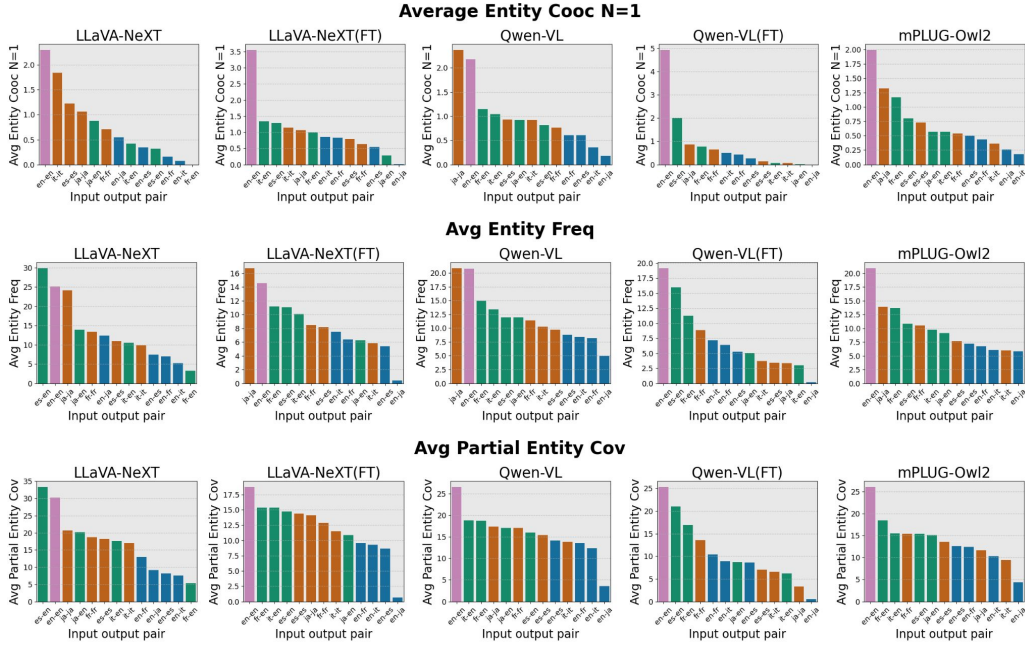


Figure 3: Some of the results in the Alignment-5 task. Purple bin indicates the method which is the instruction and the output in English ( $\{En\}$ - $\{En\}$ ), Green bin indicates the instruction in languages other than English and the output in English ( $\{Lang\}$ - $\{En\}$ ), Brown bin indicates the instruction and output in languages other than English ( $\{Lang\}$ - $\{Lang\}$ ) and Blue bin indicates the instruction in English and the output in languages other than English ( $\{En\}$ - $\{Lang\}$ ). From this figure, it can be seen that the English instructions are optimal, even if the number of data is expanded. We described further detailed results in Table 9 including Phi-3 and XComposer2. You can see the rest of the results in Figure 6 in Appendix.

suggests that “LVLMs have a poor ability to successfully transfer knowledge learned in English to other languages”. We also confirmed that this effect was more pronounced in the LoRA-Tuning model (LLaVA-NeXT(FT) and Qwen-VL(FT)).

**LVLMs’ Ability to Explain Artworks in Other Languages** Considering the multilingual explanation generation capabilities of LVLMs, a comparison between Table 3 and Table 4 reveals that performing the method is instruction and output in other same languages ( $\{Lang\}$ - $\{Lang\}$ ) generally yields better results than in the instruction in English and output in other language ( $\{En\}$ - $\{Lang\}$ ). When explaining in the native language using data trained in that language, the model effectively manages the knowledge. However, when explaining in other languages using knowledge trained in English, the model struggles to handle the information adequately. This result shows particularly clear in the cases of Qwen-VL’s results between the method is instruction and output in Chinese ( $\{Zh\}$ - $\{Zh\}$ ) and the instruction in English and output in English ( $\{Zh\}$ - $\{En\}$ ) pairs. In addition, using English

training data for LoRA Tuning likely leads to the forgetting of original performance, resulting in a decline in effectiveness. From these observations, it is clear that LVLMs currently exhibit their maximum capabilities only when instructed and output in English ( $\{En\}$ - $\{En\}$ ). Thus, future research should focus on training LVLMs in multiple languages.

## 7 Conclusion

This study focused on artworks, which have a unique image and name regardless of the language, to evaluate the explanation generation abilities of LVLMs in multilingual contexts. We created datasets compiled from Wikipedia pages in ten languages without using machine translations to evaluate their abilities across multilingual languages. The results indicate that LVLMs perform optimally when input and output are both in English, while their performance declines when using languages other than English. Thus, our hypothesis, that “Vision Encoder needs to be learned in other languages as part of its pre-training,” is correct, and might need to train Vision Transformer using multilingual data.



## 463 **Limitations**

### 464 **Data Collection and Crawling Consistency**

465 Our initial data collection was conducted through  
466 web crawling on June 30th, 2024. It is important  
467 to note that subsequent crawls may yield different  
468 results due to page updates, such as an increase  
469 in the number of pages or the addition of images.  
470 As a result, the data retrieved through repeated  
471 crawling may not consistently match the original  
472 dataset. This introduces a level of variability in the  
473 data, which must be considered when replicating  
474 or extending this research.

### 475 **Necessity of Human Evaluation Across 476 Multiple Languages**

477 To validate the effectiveness and accuracy of LLMs,  
478 especially when dealing with complex and diverse  
479 linguistic features across multiple languages, hu-  
480 man evaluation is indispensable. In this study, we  
481 conducted manual evaluations across ten languages.  
482 This step is crucial for assessing the model’s real-  
483 world applicability and ensuring that automated  
484 evaluations do not overlook nuanced errors that  
485 only human evaluators can identify.

## 486 **Ethical Considerations**

### 487 **Linguistic Considerations and Ethical 488 Implications**

489 In several languages, nouns are gendered, meaning  
490 they are classified as either masculine or feminine  
491 such as Spanish and Italian. For this study, we  
492 assumed that LLMs are capable of accurately dis-  
493 tinguishing between these gendered forms. This  
494 assumption is crucial, as it reflects the model’s  
495 ability to handle linguistic nuances, particularly in  
496 gendered languages. This raises ethical consider-  
497 ations, as any failure of the model to accurately  
498 represent gendered language could result in biased  
499 or incorrect outputs.

### 500 **Wikipedia Resources among Ten Languages**

501 Regarding Wikipedia pages, non-English versions  
502 are often less well-maintained, and whether entities  
503 are as well-organized as in English is debatable.  
504 In addition, Chinese Wikipedia contains a mix of  
505 traditional and simplified characters, which seems  
506 less standardized. In this study, since we crawled  
507 pages from Wikipedia and evaluated using their  
508 entities, it’s possible that the correct answers are  
509 included in the outputs of LVLs.

510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539  
540  
541  
542  
543  
544  
545  
546  
547  
548  
549  
550  
551  
552  
553  
554  
555  
556  
557  
558  
559  
560  
561  
562  
563  
564  
565  
566

## References

Marah Abdin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti Aneja, Ahmed Awadallah, Hany Awadalla, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Harkirat Behl, et al. 2024. Phi-3 technical report: A highly capable language model locally on your phone. *arXiv preprint arXiv:2404.14219*.

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.

Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C. Lawrence Zitnick, and Devi Parikh. 2015. VQA: Visual Question Answering. In *International Conference on Computer Vision (ICCV)*.

Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu, Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi Tan, Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang Wu, Benfeng Xu, Jin Xu, An Yang, Hao Yang, Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu, Hongyi Yuan, Zheng Yuan, Jianwei Zhang, Xingxuan Zhang, Yichang Zhang, Zhenru Zhang, Chang Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang Zhu. 2023a. Qwen technical report. *arXiv preprint arXiv:2309.16609*.

Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. 2023b. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond. *arXiv preprint arXiv:2308.12966*.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.

Soravit Changpinyo, Linting Xue, Michal Yarom, Ashish Thapliyal, Idan Szpektor, Julien Amelot, Xi Chen, and Radu Soricut. 2023. **MaXM: Towards multilingual visual question answering**. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 2667–2682, Singapore. Association for Computational Linguistics.

Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. 2020. A simple framework for contrastive learning of visual representations. *arXiv preprint arXiv:2002.05709*.

Tim Dettmers, Mike Lewis, Younes Belkada, and Luke Zettlemoyer. 2022. Gpt3. int8 (): 8-bit matrix multiplication for transformers at scale. *Advances in*

*Neural Information Processing Systems*, 35:30318–30332. 567  
568

Xiaoyi Dong, Pan Zhang, Yuhang Zang, Yuhang Cao, Bin Wang, Linke Ouyang, Xilin Wei, Songyang Zhang, Haodong Duan, Maosong Cao, et al. 2024. Internlm-xcomposer2: Mastering free-form text-image composition and comprehension in vision-language large model. *arXiv preprint arXiv:2401.16420*. 569  
570  
571  
572  
573  
574  
575

Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. 2020. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*. 576  
577  
578  
579  
580  
581  
582

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*. 583  
584  
585  
586  
587

Kazuki Hayashi, Yusuke Sakai, Hidetaka Kamigaito, Katsuhiko Hayashi, and Taro Watanabe. 2024. **Towards artwork explanation in large-scale vision language models**. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 705–729, Bangkok, Thailand. Association for Computational Linguistics. 588  
589  
590  
591  
592  
593  
594  
595

Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2015. Deep residual learning for image recognition. arxiv e-prints. *arXiv preprint arXiv:1512.03385*, 10. 596  
597  
598  
599

Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. **LoRA: Low-rank adaptation of large language models**. In *International Conference on Learning Representations*. 600  
601  
602  
603  
604

Bohao Li, Rui Wang, Guangzhi Wang, Yuying Ge, Yixiao Ge, and Ying Shan. 2023a. Seed-bench: Benchmarking multimodal llms with generative comprehension. *arXiv preprint arXiv:2307.16125*. 605  
606  
607  
608

Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023b. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *International conference on machine learning*, pages 19730–19742. PMLR. 609  
610  
611  
612  
613

Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. 2022. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International conference on machine learning*, pages 12888–12900. PMLR. 614  
615  
616  
617  
618

Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81. 619  
620  
621

622	Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2023a. Improved baselines with visual instruction tuning.	highly capable multimodal models. <i>arXiv preprint arXiv:2312.11805</i> .	677 678
625	Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. 2024. <a href="#">Llava-next: Improved reasoning, ocr, and world knowledge</a> .	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023a. Llama: Open and efficient foundation language models. <i>arXiv preprint arXiv:2302.13971</i> .	679 680 681 682 683 684
628	Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023b. Visual instruction tuning.	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023b. Llama 2: Open foundation and fine-tuned chat models. <i>arXiv preprint arXiv:2307.09288</i> .	685 686 687 688 689 690
630	Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, et al. 2023c. Mm-bench: Is your multi-modal model an all-around player? <i>arXiv preprint arXiv:2307.06281</i> .	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. <a href="#">Transformers: State-of-the-art natural language processing</a> . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations</i> , pages 38–45, Online. Association for Computational Linguistics.	691 692 693 694 695 696 697 698 699 700 701 702
635	Ngan Luu-Thuy Nguyen, Nghia Hieu Nguyen, Duong TD Vo, Khanh Quoc Tran, and Kiet Van Nguyen. 2023. Vlsq2022-evjvqa challenge: Multilingual visual question answering. <i>arXiv preprint arXiv:2302.11752</i> .	An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jin Xu, Jingren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zhihao Fan. 2024. Qwen2 technical report. <i>arXiv preprint arXiv:2407.10671</i> .	703 704 705 706 707 708 709 710 711 712 713 714 715 716 717 718
640	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. <i>Advances in neural information processing systems</i> , 35:27730–27744.	Jonas Pfeiffer, Gregor Geigle, Aishwarya Kamath, Jan-Martin O. Steitz, Stefan Roth, Ivan Vulić, and Iryna Gurevych. 2022. <a href="#">xGQA: Cross-lingual visual question answering</a> . In <i>Findings of the Association for Computational Linguistics: ACL 2022</i> , pages 2497–2511, Dublin, Ireland. Association for Computational Linguistics.	719 720 721 722 723 724 725
646	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In <i>Proceedings of the 40th annual meeting of the Association for Computational Linguistics</i> , pages 311–318.	Rifki Afina Putri, Faiz Ghifari Haznitrama, Dea Adhista, and Alice Oh. 2024. Can llm generate culturally relevant commonsense qa data? case study in indonesian and sundanese. <i>arXiv preprint arXiv:2402.17302</i> .	726 727 728 729
651	Jonas Pfeiffer, Gregor Geigle, Aishwarya Kamath, Jan-Martin O. Steitz, Stefan Roth, Ivan Vulić, and Iryna Gurevych. 2022. <a href="#">xGQA: Cross-lingual visual question answering</a> . In <i>Findings of the Association for Computational Linguistics: ACL 2022</i> , pages 2497–2511, Dublin, Ireland. Association for Computational Linguistics.	Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy Lillicrap, Jean-baptiste Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, et al. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. <i>arXiv preprint arXiv:2403.05530</i> .	730 731 732 733 734 735 736 737 738 739
658	Rifki Afina Putri, Faiz Ghifari Haznitrama, Dea Adhista, and Alice Oh. 2024. Can llm generate culturally relevant commonsense qa data? case study in indonesian and sundanese. <i>arXiv preprint arXiv:2402.17302</i> .	Yusuke Sakai, Hidetaka Kamigaito, and Taro Watanabe. 2024. mcsqa: Multilingual commonsense reasoning dataset with unified creation strategy by language models and humans. <i>arXiv preprint arXiv:2406.04215</i> .	740 741 742 743 744 745 746 747 748 749
662	Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy Lillicrap, Jean-baptiste Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, et al. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. <i>arXiv preprint arXiv:2403.05530</i> .	Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. 2023. Gemini: a family of	750 751 752 753 754 755 756 757 758 759 760 761 762 763 764 765 766 767 768 769 770 771 772 773 774 775 776 777 778 779 780 781 782 783 784 785 786 787 788 789 790 791 792 793 794 795 796 797 798 799 800

## A Appendix

### A.1 Inference and LoRA Tuning Settings

#### Inference setting

In this study, as far as inference which needs to use GPUs, all experiments were conducted on a single NVIDIA RTX A6000 GPU and NVIDIA A100-SXM4-40GB, with 8-bit quantization utilized for model generation. However, there is no InternLM-XComposer-2 with 8-bit, this model was loaded and inferred in 4-bit mode. To standardize the length of tokens generated across all models, the maximum token length was set to 1024. The same settings were applied to each model for performance comparison purposes.

#### LoRA Tuning setting

We conducted LoRA (Hu et al., 2022) Tuning with two models: LLaVA-NeXT and Qwen-VL. Both were trained using two NVIDIA A100-SXM4-40GB GPUs. Detailed parameters are provided in Table 5 and Table 6.

### B Explanation Generation Abilities from Other Languages to English ({Lang}-{En})

When considering output in English from other languages, we found this method also performs less abilities. This suggests that LVLMs have relatively less training data in languages other than English, and they may not properly understand instructions given in other languages. Thus, it is difficult to say that the integrated learning of LLMs and Vision Encoder work properly.

Hyper Parameter	Value
torch_dtype	bfloat16
seed	42
max length	2048
warmup ratio	0.01
learning rate	1e-5
batch size	4
epoch	1
lora r	64
lora alpha	16
lora dropout	0.05
lora target modules	c_attn, attn.c_proj, w1, w2

Table 5: The hyper-parameters of Qwen-VL used in the experiment, and others, were set to default settings. The implementation used Transformers (Wolf et al., 2020) and bitsandbytes (Dettmers et al., 2022).

Hyper Parameter	Value
seed	42
max length	2048
lora enable	True
learning rate	2e-5
warmup ratio	0.05
lora r	16
lora alpha	32
torch_dtype	float16

Table 6: The hyper-parameters of LLaVA-NeXT used in the experiment, and others were also set to default settings.

### C Details of experimental settings

Model	Base Model	HuggingFace Name
mPLUG-Owl2	LLaMA2-7B	MAGAAer13/mplug-owl2-llama2-7b
Qwen-VL-Chat	Qwen	Qwen/Qwen-VL-Chat
LLaVA-NeXT	LLaMA3-8B	lmms-lab/llama3-llava-next-8b
Phi-3	Phi-3-Vision-128K-Instruct	microsoft/Phi-3-vision-128k-instruct
XComposer2	internlm-xcomposer2-7B	internlm/internlm-xcomposer2-7B

### D Details of Creating Datasets or Training Data

#### D.1 How to Choose Ten Languages?

We selected ten languages with the highest number of articles from the statistics of all language versions of Wikipedia<sup>5</sup>. Of the top 10 prefectures, Cebuano, Egyptian dialects of Arabic, and Polish were deemed difficult to identify by sampling during the evaluation, so we added the runners-up, Chinese and Japanese.

#### D.2 How to Split Train, Valid, and Test Data in English?

For English, a language rich resource, we split the data into train, valid, and test data using six metrics proposed by Hayashi et al. (2024) (six metrics: page views, number of links, number of edits, number of references, number of language versions, and article length.) were used in this study as well, and the data were divided equally considering famous artworks. All data included in the alignment were used as test data so that data used in the alignment task were not included in the train. We described the number of all data in Table 7.

#### D.3 License

In our study, we created a dataset from Wikipedia articles regarding artworks. Each image is available under the Creative Commons License (CC)

<sup>5</sup>[https://en.wikipedia.org/wiki/Wikipedia:Multilingual\\_statistics](https://en.wikipedia.org/wiki/Wikipedia:Multilingual_statistics)

Language	Alignment-10 Test	Alignment-5 Test	Train	Full Valid	Test
English	28	306	6,413	2,138	2,138
French	28	306	-	-	2,707
Spanish	28	306	-	-	2,096
Italian	28	306	-	-	1,977
Russian	28	306	-	-	1,589
Japanese	28	306	-	-	1,312
German	28	306	-	-	962
Dutch	28	306	-	-	789
Swedish	28	306	-	-	774
Chinese	28	306	-	-	738

Table 7: The number of each language data in Alignment-10, Alignment-5, and Full task, split by train, valid, and test sets. We split train, valid and test sets only English due to the number of data in English.

or other licenses. Specific license information for each image can be found on the Wikipedia page or the image description page for that image. The images in this study are used under the terms of these licenses, and links to the images are provided in the datasets we publish so that users can download the images directly. The images themselves are not directly published. Thus, our data does not infringe upon the licenses.

## E Evaluation Metrics Formulation

This section describes on the evaluation metrics used in Section 5 using mathematical expressions (Hayashi et al., 2024). An explanation consisting of  $n$  sentences generated by the model is denoted as  $G = \{g_1, \dots, g_n\}$ , and a reference explanation consisting of  $m$  sentences is denoted as  $R = \{r_1, \dots, r_m\}$ . The function  $\text{Entity}(\cdot)$  is defined to extract entities contained in the input text. The notation  $|G|$  represents the total number of tokens in the generated explanation, and  $|R|$  represents the total number of tokens in the reference explanation.

**Entity Coverage (EC)** is calculated as follows:

$$EC(G, R) = Cov(G, R) \quad (1)$$

Here,  $Cov(G, R)$  is a function returning the proportion of entities in  $R$  that are covered by  $G$ . For partial matches, the Lowest Common Subsequence (LCS) is employed to calculate the longest matching length ratio in the generated explanation relative to the length of the reference entity.

**Entity F1 ( $EF_1$ )** is computed as follows:

$$EF_1 = \frac{2 \times P \times R}{P + R} \quad (2)$$

$$P = \frac{\sum_{e_i \in \text{Entity}(G)} \text{Count}_{\text{clip}}(e_i, G, R)}{\sum_{e_j \in \text{Entity}(G)} \#(e_j, G)} \quad (3)$$

$$R = \frac{\sum_{e_i \in \text{Entity}(R)} \text{Count}_{\text{clip}}(e_i, G, R)}{\sum_{e_j \in \text{Entity}(R)} \#(e_j, R)}, \quad (4)$$

where  $\#(e_j, G)$ ,  $\#(e_j, R)$  are functions that count the occurrences of entity  $e_j$  in  $G$  and  $R$  respectively, and  $\text{Count}_{\text{clip}}(e_i, G, R)$  returns the lesser frequency of occurrence of  $e_i$  in either  $G$  or  $R$ .

**Entity Cooccurrence (ECooC)** is calculated using  $BP$  from equation (6) as follows:

$$ECooC(G, R) = BP(G, R) \times Cov(Co(G), Co(R)), \quad (5)$$

where  $BP(G, R)$  is given by:

$$BP(G, R) = \exp(\max(0.0, \frac{|G|}{|R|} - 1)) \quad (6)$$

and the function  $Co(\cdot)$  returns pairs of co-occurring entities within a context window comprising a sentence and its adjacent  $n$  sentences. Sentence segmentation was performed using the nltk sentence splitter for this purpose.<sup>6</sup>

## F Filtered Sections

The following section was filtered in this study. Approximately 30 instances from the Alignment-10 task were reviewed, and sections without informative content.

<sup>6</sup>Sentence segmentation was performed using the NLTK sentence splitter.

### English

References, See also, External links, Sources, Further reading, Bibliography, Gallery, Footnotes, Notes References, References Sources, Bibliography (In Spanish), Bibliography (In Italian), Bibliography (In German), Bibliography (In French), Images, Note, Links, Notes, List, Notes and references, List by location

### Japanese

外部リンク, 参考文献, 関連項目, 脚注, 出典, ギャラリー, バージョン, 注釈, 関連する作品

### Italian

Collegamenti esterni, Altri progetti, Bibliografia, Note, Omaggi, Voci correlate, Bibliografia, Musica, Fumetti, Letteratura, Filmografia, Nella cultura di massa, Altri progetti, Galleria d'immagini, Curiosità, Calendario

### French

Liens externes, Articles connexes, Bibliographie et ressources en ligne, Annexes, Notes et références, Divers, Littérature, Peinture et sculpture, Déclinaisons et détournements, Bases de données et dictionnaires, Italien, Français, Ouvrages, Articles, Bibliographie, Théâtre, Cinéma, Article connexe, Annexe, Notes et référence, Voir aussi, Divers, Pour approfondir, Versions, Références, Sources secondaires, Sources originales, Références de l'expression dans l'art, Ouvrages, Ailleurs, Notes, Films, Dans la culture, Postérité, Données techniques, Galerie, Historique

### Spanish

Enlaces externos, Bibliografía, Referencias, Fuentes, Enlaces externos, Bibliografía, Véase también, Notas, Información, Galería, Galería de imágenes, Filmografía

### Chinese (Traditional)

外部連結, 延伸, 参考文献, 參考文獻, 參見, 參見, 書目, 注与参考文献, 來源, 擴展讀, 參考来源, 外部接, 延伸, 引用, 注, 參考資料, 參考料, 相關條目, 參考來源, 參見條目, 其他事項, 參考, 註解, 媒體, 紀錄片, 書籍, 近似作品, 相關作品, 德文, 注, 擴展讀, 吉米·威士的声明

### Chinese (Simplified)

外部链接, 延伸阅读, 参考文献, 参见, 注释与参考文献, 来源, 扩展阅读, 参考来源, 引用, 注释, 参考资料, 相关条目, 参见条目, 其他事项, 参考, 近似作品, 媒体, 纪录片, 书籍, 注释, 吉米·威尔士的声明

### Swedish

Noter, Referenser, Se även, Externa länkar, Allmänna källor, Galleri, Källor, Bilder, Kalenderfunktionen, Relaterade målningar

### Dutch

Zie ook, Literatuur, Externe links, Bewerkingen, Andere, Latere edities, Trivia, Zie ook, Galerij, Originele gietingen, Stanza dell'incendio del Borgo, Stanza della Segnatura, Noten, Literatuur en bronnen

### Russian

Ссылки, Примечания, См. также, Документалистика, Литература, Источники, Отражение в искусстве

### German

Anmerkungen, Weblinks, Literatur, Anmerkungen und Einzelnachweise, Einzelbelege, Einzelnachweise, Chronologie, Quellen, Übersicht, Literatur (Auswahl), Siehe auch, Rezeption, Dokumentarfilme, Ausstellungen, Siehe auch

850

851

852

853

854

855

845

846

847

848

849

## G Instruction to Native Speakers

856

We asked native speaker to prepare the instruction to check if the above template is equal in difficulty compared to the English text.

857

858

### Instruction

#### # What we research

We are conducting a study to measure LLMs' ability to understand the arts. Previous studies have been done only for English, and we are now trying to extend and validate it for multiple languages.

The text presented has been translated from English into your language using DeepL.

I want you to make sure that the sentence you translate has the same meaning as the English sentence.

The time I assume will not take more than 5 minutes and that's about OK for a check. I also use back translation to check it, so I believe it is not that broken.

#### # Keep in mind

- My final goal is to have the sentences corrected to be as natural as English sentences.
- Depending on {title} and {section}, and in some countries, you may need to be concerned about masculine and feminine nouns. If that is the case, choose whichever you type into the LLM in your native language (i.e., the more natural one).
- Please do not change the entire text.
- Changing, deleting or adding words is acceptable.

#### # Examples of {title} and {section}, {subsection} and {subsubsection}

We use Wikipedia for our research.

Here is one of the example: [https://en.wikipedia.org/wiki/Mona\\_Lisa](https://en.wikipedia.org/wiki/Mona_Lisa)

In this case, {title} will contain "Mona Lisa".

In addition, {section} contains "Description", "History", and so on.

{subsection} refers to a smaller frame within {section}, such as "Creation and date".

↓ Below is the text I would like you to review.

({lang}\_temp1\_sec is translated from en\_temp1\_sec using DeepL) ({lang}\_temp2\_subsec is translated from en\_temp2\_subsec using DeepL)

#### #English (source)

This sentence is a sample.

#### # Your native language (target I translated from DeepL.)

This sentence is a sample.

859

## H Other Results and Visualizations

Input	Output	LVL	Entity Cov.		Entity F1	Entity Cooccurrence				Avg. Length
			exact	partial		n=0	n=1	n=2	n=∞	
En	En	LLaVA-NeXT	17.66	26.05	18.55	1.31	1.23	1.20	1.20	242
		LLaVA-NeXT (FT)	17.92	23.65	19.20	<b>5.67</b>	5.66	5.63	5.60	81
		Qwen-VL	14.60	21.51	15.39	1.36	1.33	1.28	1.27	110
		Qwen-VL (FT)	<b>20.09</b>	<b>26.27</b>	<b>20.84</b>	5.67	<b>5.78</b>	<b>5.77</b>	<b>5.72</b>	171
		mPLUG-Owl2	14.41	21.96	15.71	1.27	1.17	1.14	1.10	121
En	Es	LLaVA-NeXT	<b>10.40 (-7.26)</b>	<b>16.05 (-10.00)</b>	<b>10.86 (-7.69)</b>	<b>0.79 (-0.52)</b>	<b>0.78 (-0.45)</b>	<b>0.83 (-0.37)</b>	<b>0.83 (-0.37)</b>	181
		LLaVA-NeXT (FT)	4.96 (-12.96)	8.42 (-15.23)	5.40 (-13.80)	0.56 (-5.10)	0.57 (-5.10)	0.58 (-5.04)	0.57 (-5.02)	90
		Qwen-VL	8.11 (-6.49)	13.18 (-8.33)	8.66 (-6.73)	0.53 (-0.83)	0.50 (-0.83)	0.52 (-0.76)	0.51 (-0.76)	103
		Qwen-VL (FT)	4.23 (-15.86)	8.47 (-17.80)	4.66 (-16.17)	0.23 (-5.43)	0.23 (-5.56)	0.24 (-5.53)	0.24 (-5.48)	195
		mPLUG-Owl2	7.26 (-7.14)	12.13 (-9.83)	7.55 (-8.16)	0.45 (-0.82)	0.49 (-0.68)	0.52 (-0.63)	0.51 (-0.59)	100
En	Fr	LLaVA-NeXT	<b>9.71 (-7.95)</b>	<b>16.17 (-9.88)</b>	<b>9.49 (-9.06)</b>	<b>0.57 (-0.74)</b>	<b>0.57 (-0.66)</b>	<b>0.57 (-0.63)</b>	<b>0.55 (-0.64)</b>	168
		LLaVA-NeXT (FT)	7.02 (-10.90)	10.37 (-13.29)	7.60 (-11.60)	<b>0.84 (-4.83)</b>	<b>0.84 (-4.82)</b>	<b>0.82 (-4.81)</b>	<b>0.81 (-4.79)</b>	60
		Qwen-VL	7.64 (-6.96)	12.82 (-8.68)	7.71 (-7.68)	0.51 (-0.85)	0.46 (-0.87)	0.45 (-0.83)	0.43 (-0.83)	86
		Qwen-VL (FT)	6.42 (-13.68)	11.17 (-15.10)	6.88 (-13.95)	0.43 (-5.24)	0.56 (-5.22)	0.55 (-5.22)	0.54 (-5.18)	155
		mPLUG-Owl2	6.99 (-7.42)	12.55 (-9.41)	6.91 (-8.79)	0.41 (-0.86)	0.38 (-0.79)	0.37 (-0.77)	0.35 (-0.75)	95
En	De	LLaVA-NeXT	<b>10.32 (-7.35)</b>	<b>13.84 (-12.21)</b>	<b>12.28 (-6.27)</b>	<b>0.90 (-0.41)</b>	<b>0.88 (-0.34)</b>	<b>0.88 (-0.32)</b>	<b>0.86 (-0.34)</b>	161
		LLaVA-NeXT (FT)	5.52 (-12.40)	7.80 (-15.86)	5.93 (-13.26)	0.52 (-5.15)	0.48 (-5.19)	0.46 (-5.17)	0.45 (-5.15)	75
		Qwen-VL	7.75 (-6.85)	10.60 (-10.91)	8.69 (-6.69)	0.63 (-0.73)	0.59 (-0.74)	0.58 (-0.70)	0.56 (-0.71)	99
		Qwen-VL (FT)	4.79 (-15.30)	7.40 (-18.87)	5.17 (-15.67)	0.23 (-5.44)	0.25 (-5.53)	0.24 (-5.53)	0.24 (-5.48)	177
		mPLUG-Owl2	6.87 (-7.53)	9.66 (-12.30)	7.69 (-8.01)	0.60 (-0.67)	0.54 (-0.63)	0.53 (-0.61)	0.50 (-0.60)	91
En	It	LLaVA-NeXT	<b>9.57 (-8.10)</b>	<b>16.52 (-9.53)</b>	<b>10.72 (-7.83)</b>	0.72 (-0.60)	0.72 (-0.50)	0.74 (-0.46)	0.72 (-0.47)	168
		LLaVA-NeXT (FT)	6.21 (-11.71)	9.51 (-14.15)	7.59 (-11.61)	<b>0.79 (-4.88)</b>	<b>0.86 (-4.80)</b>	<b>0.85 (-4.77)</b>	<b>0.85 (-4.74)</b>	87
		Qwen-VL	7.08 (-7.52)	12.73 (-8.77)	8.26 (-7.13)	0.34 (-1.02)	0.38 (-0.95)	0.38 (-0.90)	0.38 (-0.88)	112
		Qwen-VL (FT)	6.08 (-14.01)	10.10 (-16.17)	7.39 (-13.44)	0.49 (-5.17)	0.58 (-5.20)	0.59 (-5.19)	0.59 (-5.13)	187
		mPLUG-Owl2	6.54 (-7.86)	12.20 (-9.76)	7.44 (-8.27)	0.42 (-0.85)	0.40 (-0.77)	0.39 (-0.75)	0.39 (-0.71)	102
En	Nl	LLaVA-NeXT	<b>7.91 (-9.76)</b>	<b>13.25 (-12.80)</b>	8.63 (-9.92)	0.31 (-1.01)	0.44 (-0.79)	0.42 (-0.78)	0.43 (-0.77)	175
		LLaVA-NeXT (FT)	7.89 (-10.03)	11.66 (-12.00)	<b>8.81 (-10.39)</b>	<b>1.22 (-4.44)</b>	<b>1.13 (-4.53)</b>	<b>1.12 (-4.51)</b>	<b>1.11 (-4.49)</b>	102
		Qwen-VL	7.41 (-7.19)	12.33 (-9.18)	7.93 (-7.46)	0.35 (-1.01)	0.49 (-0.84)	0.50 (-0.78)	0.53 (-0.74)	137
		Qwen-VL (FT)	6.67 (-13.42)	10.07 (-16.21)	7.67 (-13.16)	0.68 (-4.98)	0.73 (-5.05)	0.70 (-5.08)	0.71 (-5.01)	166
		mPLUG-Owl2	4.61 (-9.80)	8.96 (-13.00)	4.84 (-10.87)	0.20 (-1.06)	0.26 (-0.91)	0.25 (-0.89)	0.25 (-0.85)	106
En	Sv	LLaVA-NeXT	<b>13.08 (-4.59)</b>	<b>17.19 (-8.85)</b>	<b>12.38 (-6.18)</b>	<b>0.89 (-0.42)</b>	<b>0.82 (-0.41)</b>	<b>0.82 (-0.39)</b>	<b>0.75 (-0.44)</b>	172
		LLaVA-NeXT (FT)	9.44 (-8.47)	12.79 (-10.87)	9.62 (-9.58)	0.73 (-4.94)	0.64 (-5.03)	0.60 (-5.03)	0.58 (-5.01)	94
		Qwen-VL	10.59 (-4.01)	14.72 (-6.79)	10.75 (-4.64)	0.58 (-0.78)	0.61 (-0.72)	0.66 (-0.62)	0.61 (-0.66)	124
		Qwen-VL (FT)	9.47 (-10.62)	13.20 (-13.07)	9.93 (-10.91)	0.72 (-4.95)	0.65 (-5.13)	0.63 (-5.14)	0.58 (-5.14)	155
		mPLUG-Owl2	9.37 (-5.03)	12.82 (-9.14)	8.53 (-7.17)	0.40 (-0.86)	0.36 (-0.81)	0.36 (-0.78)	0.33 (-0.77)	79
En	Ru	LLaVA-NeXT	<b>7.86 (-9.81)</b>	<b>10.75 (-15.29)</b>	<b>6.39 (-12.16)</b>	<b>0.22 (-1.09)</b>	<b>0.26 (-0.97)</b>	<b>0.28 (-0.92)</b>	<b>0.28 (-0.92)</b>	203
		LLaVA-NeXT (FT)	0.42 (-17.50)	1.51 (-22.14)	0.31 (-18.89)	0.01 (-5.66)	0.01 (-5.65)	0.01 (-5.62)	0.01 (-5.58)	72
		Qwen-VL	3.05 (-11.55)	4.81 (-16.69)	2.35 (-13.04)	0.05 (-1.31)	0.07 (-1.26)	0.08 (-1.20)	0.08 (-1.18)	112
		Qwen-VL (FT)	0.15 (-19.94)	1.09 (-25.19)	0.09 (-20.74)	0.00 (-5.67)	0.00 (-5.78)	0.00 (-5.77)	0.00 (-5.72)	203
		mPLUG-Owl2	3.69 (-10.71)	5.33 (-16.64)	2.83 (-12.88)	0.11 (-1.16)	0.10 (-1.07)	0.09 (-1.05)	0.10 (-1.00)	107
En	Ja	LLaVA-NeXT	<b>8.65 (-9.01)</b>	<b>8.70 (-17.35)</b>	<b>12.34 (-6.21)</b>	<b>0.44 (-0.87)</b>	<b>0.44 (-0.79)</b>	<b>0.44 (-0.76)</b>	<b>0.44 (-0.76)</b>	213
		LLaVA-NeXT (FT)	0.46 (-17.45)	0.61 (-23.04)	0.51 (-18.69)	0.02 (-5.65)	0.02 (-5.65)	0.02 (-5.61)	0.02 (-5.58)	67
		Qwen-VL	3.10 (-11.50)	3.16 (-18.35)	4.37 (-11.02)	0.12 (-1.24)	0.12 (-1.21)	0.12 (-1.16)	0.12 (-1.14)	127
		Qwen-VL (FT)	0.21 (-19.88)	0.46 (-25.82)	0.12 (-20.72)	0.00 (-5.67)	0.00 (-5.78)	0.00 (-5.77)	0.00 (-5.72)	152
		mPLUG-Owl2	4.00 (-10.40)	4.06 (-17.90)	5.39 (-10.32)	0.25 (-1.01)	0.25 (-0.92)	0.25 (-0.89)	0.25 (-0.85)	104
En	Zh	LLaVA-NeXT	<b>10.81 (-6.86)</b>	<b>10.90 (-15.15)</b>	<b>13.00 (-5.56)</b>	0.60 (-0.71)	0.60 (-0.62)	0.60 (-0.60)	0.60 (-0.59)	220
		LLaVA-NeXT (FT)	0.64 (-17.27)	0.89 (-22.76)	0.75 (-18.45)	0.08 (-5.59)	0.08 (-5.59)	0.08 (-5.55)	0.08 (-5.52)	71
		Qwen-VL	8.60 (-6.00)	8.65 (-12.85)	10.34 (-5.05)	<b>0.80 (-0.56)</b>	<b>0.79 (-0.54)</b>	<b>0.79 (-0.49)</b>	<b>0.79 (-0.47)</b>	133
		Qwen-VL (FT)	0.35 (-19.74)	0.64 (-25.63)	0.27 (-20.57)	0.01 (-5.66)	0.01 (-5.77)	0.01 (-5.77)	0.01 (-5.71)	155
		mPLUG-Owl2	4.99 (-9.42)	5.04 (-16.92)	6.08 (-9.62)	0.52 (-0.75)	0.52 (-0.65)	0.52 (-0.63)	0.52 (-0.58)	107

Table 8: Results of LVLs in Full Task. Bold fonts indicate the best score for that language combination. This result shows that no matter how much the amount of data is increased, the best performance is achieved by having instructions given and output in English. The values are noted next to the output of the difference by the same model in the method with instruction and output in English ({En}-{En}). Red indicates a higher value than that method; blue indicates a lower value.





Input	Output	LVLM	Entity Cov.		Entity F1	Entity Cooccurrence				Avg. Length
			exact	partial		n=0	n=1	n=2	n=∞	
En	En	Phi-3	9.74	13.75	8.06	0.95	0.86	0.92	0.90	108
		XComposer2	<b>16.57</b>	<b>21.56</b>	<b>16.53</b>	<b>1.51</b>	<b>1.47</b>	<b>1.43</b>	<b>1.36</b>	223
En	Es	Phi-3	<b>10.35 (+0.61)</b>	<b>14.46 (+0.72)</b>	<b>11.67 (+3.61)</b>	<b>1.20 (+0.25)</b>	<b>1.37 (+0.51)</b>	<b>1.43 (+0.51)</b>	<b>1.38 (+0.48)</b>	141
		XComposer2	10.03 (-6.55)	13.75 (-7.81)	10.04 (-6.49)	0.79 (-0.73)	0.96 (-0.51)	0.92 (-0.50)	0.89 (-0.47)	116
Es	En	Phi-3	7.99 (-1.75)	14.53 (+0.78)	7.35 (-0.71)	<b>0.90 (-0.05)</b>	<b>0.79 (-0.07)</b>	<b>0.80 (-0.12)</b>	<b>0.75 (-0.16)</b>	91
		XComposer2	<b>9.58 (-6.99)</b>	<b>15.05 (-6.51)</b>	<b>9.03 (-7.50)</b>	0.67 (-0.85)	0.55 (-0.92)	0.61 (-0.82)	0.51 (-0.85)	69
Es	Es	Phi-3	<b>12.81 (+3.07)</b>	<b>16.25 (+2.50)</b>	<b>13.04 (+4.98)</b>	1.10 (+0.16)	1.23 (+0.37)	<b>1.24 (+0.32)</b>	<b>1.21 (+0.31)</b>	190
		XComposer2	9.19 (-7.38)	11.93 (-9.63)	8.68 (-7.85)	<b>1.11 (-0.40)</b>	<b>1.24 (-0.24)</b>	1.17 (-0.26)	1.14 (-0.21)	77
En	Fr	Phi-3	9.00 (-0.74)	14.17 (+0.42)	<b>9.84 (+1.77)</b>	<b>1.60 (+0.65)</b>	<b>1.12 (+0.26)</b>	<b>1.11 (+0.19)</b>	<b>1.03 (+0.13)</b>	151
		XComposer2	<b>9.32 (-7.25)</b>	<b>14.32 (-7.24)</b>	9.28 (-7.25)	1.09 (-0.42)	0.89 (-0.59)	0.81 (-0.62)	0.74 (-0.62)	79
Fr	En	Phi-3	10.37 (+0.63)	16.77 (+3.03)	9.83 (+1.77)	<b>1.22 (+0.27)</b>	<b>1.13 (+0.26)</b>	<b>1.10 (+0.17)</b>	<b>1.11 (+0.21)</b>	154
		XComposer2	<b>10.79 (-5.78)</b>	<b>17.10 (-4.46)</b>	<b>10.12 (-6.41)</b>	0.72 (-0.79)	0.64 (-0.83)	0.62 (-0.81)	0.63 (-0.73)	121
Fr	Fr	Phi-3	<b>11.06 (+1.32)</b>	<b>16.04 (+2.30)</b>	<b>10.21 (+2.15)</b>	<b>0.95 (+0.00)</b>	<b>0.87 (+0.01)</b>	<b>0.84 (-0.08)</b>	<b>0.86 (-0.04)</b>	220
		XComposer2	9.90 (-6.67)	14.07 (-7.49)	8.93 (-7.60)	0.73 (-0.79)	0.61 (-0.86)	0.56 (-0.87)	0.50 (-0.85)	116
En	De	Phi-3	<b>9.55 (-0.19)</b>	<b>13.14 (-0.61)</b>	<b>11.88 (+3.82)</b>	<b>0.80 (-0.15)</b>	<b>0.80 (-0.07)</b>	<b>0.80 (-0.13)</b>	<b>0.95 (+0.04)</b>	216
		XComposer2	8.02 (-8.55)	11.26 (-10.30)	9.48 (-7.05)	0.78 (-0.73)	0.77 (-0.70)	0.75 (-0.68)	0.78 (-0.57)	107
De	En	Phi-3	9.68 (-0.06)	16.53 (+2.78)	<b>10.01 (+1.95)</b>	0.67 (-0.27)	0.62 (-0.24)	0.70 (-0.22)	0.72 (-0.19)	183
		XComposer2	<b>11.05 (-5.52)</b>	<b>17.07 (-4.49)</b>	9.39 (-7.13)	<b>0.95 (-0.56)</b>	<b>0.89 (-0.58)</b>	<b>0.90 (-0.52)</b>	<b>0.91 (-0.45)</b>	86
De	De	Phi-3	<b>13.37 (+3.63)</b>	<b>16.55 (+2.80)</b>	<b>15.57 (+7.51)</b>	<b>1.72 (+0.77)</b>	<b>1.50 (+0.64)</b>	<b>1.53 (+0.61)</b>	<b>1.55 (+0.65)</b>	240
		XComposer2	9.89 (-6.68)	11.93 (-9.63)	11.73 (-4.80)	1.17 (-0.34)	0.91 (-0.56)	0.81 (-0.62)	0.80 (-0.56)	107
En	It	Phi-3	<b>4.98 (-4.76)</b>	<b>8.61 (-5.14)</b>	<b>5.63 (-2.43)</b>	0.08 (-0.87)	0.12 (-0.74)	0.14 (-0.78)	0.12 (-0.78)	150
		XComposer2	4.58 (-11.99)	8.40 (-13.16)	5.44 (-11.09)	<b>0.35 (-1.16)</b>	<b>0.19 (-1.28)</b>	<b>0.19 (-1.23)</b>	<b>0.19 (-1.16)</b>	87
It	En	Phi-3	7.11 (-2.63)	13.70 (-0.05)	7.04 (-1.03)	<b>0.89 (-0.06)</b>	<b>0.95 (+0.08)</b>	<b>0.95 (+0.03)</b>	<b>0.93 (+0.03)</b>	143
		XComposer2	<b>9.82 (-6.75)</b>	<b>16.56 (-5.00)</b>	<b>8.72 (-7.81)</b>	0.58 (-0.93)	0.64 (-0.83)	0.68 (-0.75)	0.64 (-0.71)	94
It	It	Phi-3	<b>7.93 (-1.81)</b>	<b>11.61 (-2.14)</b>	<b>9.17 (+1.11)</b>	0.14 (-0.81)	0.11 (-0.75)	0.11 (-0.81)	0.11 (-0.79)	183
		XComposer2	4.29 (-12.29)	6.92 (-14.64)	5.27 (-11.26)	<b>0.22 (-1.29)</b>	<b>0.18 (-1.29)</b>	<b>0.18 (-1.25)</b>	<b>0.18 (-1.18)</b>	65
En	Nl	Phi-3	4.07 (-5.67)	8.41 (-5.34)	5.03 (-3.03)	0.04 (-0.91)	0.04 (-0.82)	0.03 (-0.89)	0.05 (-0.85)	240
		XComposer2	<b>5.56 (-11.01)</b>	<b>10.07 (-11.49)</b>	<b>5.84 (-10.69)</b>	<b>0.29 (-1.22)</b>	<b>0.52 (-0.95)</b>	<b>0.50 (-0.92)</b>	<b>0.52 (-0.84)</b>	78
Nl	En	Phi-3	4.09 (-5.65)	9.41 (-4.34)	4.17 (-3.89)	0.50 (-0.45)	0.35 (-0.51)	0.32 (-0.61)	0.31 (-0.59)	213
		XComposer2	<b>9.74 (-6.83)</b>	<b>17.17 (-4.39)</b>	<b>9.44 (-7.09)</b>	<b>1.21 (-0.30)</b>	<b>0.99 (-0.48)</b>	<b>0.98 (-0.45)</b>	<b>0.93 (-0.43)</b>	92
Nl	Nl	Phi-3	8.37 (-1.37)	13.12 (-0.63)	6.26 (-1.81)	0.01 (-0.94)	<b>1.38 (+0.52)</b>	<b>1.38 (+0.46)</b>	<b>1.38 (+0.47)</b>	273
		XComposer2	<b>10.73 (-5.85)</b>	<b>14.59 (-6.97)</b>	<b>9.95 (-6.58)</b>	<b>0.08 (-1.44)</b>	0.92 (-0.55)	0.94 (-0.49)	0.94 (-0.42)	73
En	Sv	Phi-3	6.55 (-3.19)	9.20 (-4.55)	6.18 (-1.88)	0.05 (-0.89)	0.04 (-0.82)	0.04 (-0.88)	0.03 (-0.87)	235
		XComposer2	<b>8.03 (-8.55)</b>	<b>10.90 (-10.66)</b>	<b>6.91 (-9.62)</b>	<b>0.31 (-1.20)</b>	<b>0.32 (-1.16)</b>	<b>0.30 (-1.12)</b>	<b>0.29 (-1.07)</b>	76
Sv	En	Phi-3	4.69 (-5.05)	10.04 (-3.71)	4.58 (-3.48)	0.51 (-0.44)	0.46 (-0.41)	0.51 (-0.41)	0.45 (-0.45)	176
		XComposer2	<b>11.35 (-5.23)</b>	<b>13.91 (-7.65)</b>	<b>8.60 (-7.93)</b>	<b>1.23 (-0.28)</b>	<b>0.48 (-0.99)</b>	<b>0.56 (-0.87)</b>	<b>0.52 (-0.83)</b>	78
Sv	Sv	Phi-3	<b>14.03 (+4.29)</b>	<b>15.53 (+1.78)</b>	<b>9.92 (+1.85)</b>	0.74 (-0.21)	0.27 (-0.60)	0.26 (-0.66)	0.26 (-0.65)	194
		XComposer2	11.58 (-4.99)	13.07 (-8.50)	8.41 (-8.12)	<b>1.26 (-0.26)</b>	<b>0.44 (-1.03)</b>	<b>0.44 (-0.99)</b>	<b>0.43 (-0.93)</b>	63
En	Ru	Phi-3	0.61 (-9.13)	2.17 (-11.58)	0.31 (-7.75)	<b>0.00 (-0.95)</b>	0.00 (-0.86)	0.00 (-0.92)	0.00 (-0.90)	194
		XComposer2	<b>3.70 (-12.88)</b>	<b>6.50 (-15.06)</b>	<b>3.07 (-13.46)</b>	<b>0.00 (-1.51)</b>	<b>0.00 (-1.47)</b>	<b>0.01 (-1.42)</b>	<b>0.01 (-1.34)</b>	73
Ru	En	Phi-3	<b>6.62 (-3.12)</b>	<b>12.82 (-0.92)</b>	<b>6.57 (-1.50)</b>	<b>0.31 (-0.63)</b>	<b>0.47 (-0.39)</b>	<b>0.47 (-0.45)</b>	<b>0.46 (-0.44)</b>	147
		XComposer2	4.69 (-11.89)	8.56 (-13.00)	4.07 (-12.46)	0.19 (-1.32)	0.11 (-1.37)	0.14 (-1.28)	0.10 (-1.25)	62
Ru	Ru	Phi-3	2.42 (-7.32)	4.45 (-9.30)	1.58 (-6.48)	<b>0.21 (-0.74)</b>	<b>0.21 (-0.66)</b>	<b>0.20 (-0.72)</b>	<b>0.20 (-0.70)</b>	269
		XComposer2	<b>2.73 (-13.84)</b>	<b>4.85 (-16.71)</b>	<b>2.22 (-14.31)</b>	0.17 (-1.35)	0.16 (-1.31)	0.15 (-1.28)	0.14 (-1.21)	45
En	Ja	Phi-3	2.53 (-7.21)	2.53 (-11.22)	3.26 (-4.80)	0.06 (-0.89)	0.06 (-0.80)	0.06 (-0.86)	0.06 (-0.84)	202
		XComposer2	<b>3.27 (-13.31)</b>	<b>3.27 (-18.29)</b>	<b>3.78 (-12.75)</b>	<b>0.21 (-1.30)</b>	<b>0.21 (-1.26)</b>	<b>0.21 (-1.22)</b>	<b>0.21 (-1.15)</b>	109
Ja	En	Phi-3	8.17 (-1.57)	15.38 (+1.64)	7.46 (-0.60)	0.37 (-0.58)	0.40 (-0.46)	0.44 (-0.48)	0.38 (-0.52)	168
		XComposer2	<b>10.59 (-5.99)</b>	<b>18.39 (-3.17)</b>	<b>9.71 (-6.82)</b>	<b>0.61 (-0.90)</b>	<b>0.57 (-0.90)</b>	<b>0.54 (-0.89)</b>	<b>0.50 (-0.86)</b>	159
Ja	Ja	Phi-3	<b>8.73 (-1.01)</b>	<b>8.74 (-5.01)</b>	<b>12.19 (+4.13)</b>	0.88 (-0.06)	0.88 (+0.02)	0.88 (-0.04)	0.88 (-0.02)	214
		XComposer2	6.04 (-10.53)	6.04 (-15.52)	7.18 (-9.34)	<b>1.22 (-0.29)</b>	<b>1.22 (-0.25)</b>	<b>1.22 (-0.20)</b>	<b>1.22 (-0.13)</b>	133
En	Zh	Phi-3	4.48 (-5.26)	4.52 (-9.23)	5.14 (-2.92)	0.13 (-0.81)	0.14 (-0.73)	0.14 (-0.79)	0.14 (-0.77)	145
		XComposer2	<b>13.35 (-3.23)</b>	<b>13.38 (-8.18)</b>	<b>16.00 (-0.53)</b>	<b>0.69 (-0.83)</b>	<b>0.68 (-0.80)</b>	<b>0.67 (-0.75)</b>	<b>0.67 (-0.69)</b>	124
Zh	En	Phi-3	3.74 (-6.00)	3.74 (-10.01)	6.11 (-1.95)	0.14 (-0.81)	0.14 (-0.72)	0.14 (-0.78)	0.14 (-0.76)	186
		XComposer2	<b>9.27 (-7.30)</b>	<b>9.27 (-12.29)</b>	<b>11.95 (-4.58)</b>	<b>0.25 (-1.26)</b>	<b>0.24 (-1.23)</b>	<b>0.24 (-1.18)</b>	<b>0.24 (-1.11)</b>	215
Zh	Zh	Phi-3	2.53 (-7.21)	2.53 (-11.22)	<b>4.94 (-3.12)</b>	0.00 (-0.95)	0.00 (-0.86)	0.00 (-0.92)	0.00 (-0.90)	55
		XComposer2	<b>2.87 (-13.70)</b>	<b>2.87 (-18.69)</b>	4.44 (-12.09)	<b>0.09 (-1.42)</b>	<b>0.09 (-1.38)</b>	<b>0.09 (-1.34)</b>	<b>0.09 (-1.26)</b>	55

Table 10: Results for Phi-3 and XComposer2 in the Alignment-10 task. Bold fonts indicate the best score for that language combination. The values are noted next to the output of the difference by the same model in the method with instruction and output in English ({En}-{En}). Red indicates a higher value than that method; blue indicates a lower value.

Input	Output	LVL	Entity Cov.		Entity F1	Entity Cooccurrence				Avg. Length
			exact	partial		n=0	n=1	n=2	n=∞	
En	En	LLaVA-NeXT	<b>26.49</b>	<b>31.54</b>	<b>26.07</b>	1.35	1.65	1.66	1.70	252
		LLaVA-NeXT (FT)	16.98	22.70	19.95	3.02	3.33	3.23	3.11	83
		Qwen-VL	21.11	27.03	21.78	1.60	1.59	1.56	1.52	155
		Qwen-VL (FT)	21.12	24.87	21.95	<b>3.57</b>	<b>3.83</b>	<b>3.78</b>	<b>3.68</b>	177
		mPLUG-Owl2	12.79	17.08	13.48	2.07	1.68	1.59	1.56	151
Es	En	LLaVA-NeXT	4.22 (-22.26)	6.22 (-25.33)	4.17 (-21.90)	0.00 (-1.34)	0.00 (-1.65)	0.00 (-1.66)	0.00 (-1.70)	180
		LLaVA-NeXT (FT)	13.64 (-3.34)	18.60 (-4.10)	13.54 (-6.41)	1.35 (-1.67)	1.44 (-1.89)	1.35 (-1.88)	1.35 (-1.76)	123
		Qwen-VL	10.38 (-10.73)	16.82 (-10.21)	10.40 (-11.38)	0.72 (-0.88)	0.62 (-0.96)	0.64 (-0.92)	0.55 (-0.97)	60
		Qwen-VL (FT)	<b>15.50 (-5.62)</b>	<b>22.52 (-2.34)</b>	<b>16.19 (-5.76)</b>	1.31 (-2.27)	<b>1.70 (-2.13)</b>	<b>1.76 (-2.01)</b>	<b>1.65 (-2.02)</b>	199
		mPLUG-Owl2	11.97 (-0.82)	18.26 (+1.18)	11.68 (-1.80)	<b>1.72 (-0.35)</b>	1.28 (-0.40)	1.17 (-0.42)	1.11 (-0.45)	54
Fr	En	LLaVA-NeXT	<b>21.11 (-5.38)</b>	<b>28.40 (-3.14)</b>	<b>21.22 (-4.85)</b>	<b>2.01 (+0.67)</b>	<b>2.04 (+0.38)</b>	<b>2.09 (+0.43)</b>	<b>1.97 (+0.27)</b>	232
		LLaVA-NeXT (FT)	15.34 (-1.64)	21.24 (-1.46)	15.24 (-4.71)	1.21 (-1.81)	1.11 (-2.22)	1.08 (-2.15)	1.00 (-2.11)	101
		Qwen-VL	14.24 (-6.86)	20.10 (-6.93)	15.24 (-6.54)	1.43 (-0.17)	1.38 (-0.21)	1.32 (-0.24)	1.27 (-0.25)	144
		Qwen-VL (FT)	12.92 (-8.19)	21.49 (-3.38)	11.85 (-10.10)	0.35 (-3.22)	0.66 (-3.17)	0.64 (-3.14)	0.62 (-3.05)	275
		mPLUG-Owl2	15.56 (+2.77)	22.04 (+4.96)	14.34 (+0.86)	1.69 (-0.38)	1.61 (-0.06)	1.65 (+0.06)	1.47 (-0.09)	106
De	En	LLaVA-NeXT	<b>21.88 (-4.60)</b>	<b>30.48 (-1.07)</b>	<b>21.76 (-4.31)</b>	0.89 (-0.46)	<b>1.33 (-0.33)</b>	<b>1.53 (-0.13)</b>	1.50 (-0.20)	239
		LLaVA-NeXT (FT)	14.15 (-2.83)	19.99 (-2.71)	12.43 (-7.52)	0.68 (-2.33)	0.90 (-2.43)	0.86 (-2.37)	0.85 (-2.26)	128
		Qwen-VL	17.36 (-3.74)	25.21 (-1.82)	16.82 (-4.95)	<b>1.11 (-0.49)</b>	1.31 (-0.28)	1.49 (-0.07)	<b>1.51 (-0.01)</b>	109
		Qwen-VL (FT)	13.60 (-7.51)	21.30 (-3.57)	13.17 (-8.77)	0.80 (-2.77)	1.12 (-2.72)	1.14 (-2.64)	1.07 (-2.61)	265
		mPLUG-Owl2	13.29 (+0.50)	19.74 (+2.66)	12.42 (-1.06)	0.78 (-1.28)	0.77 (-0.91)	0.80 (-0.79)	0.77 (-0.79)	75
It	En	LLaVA-NeXT	7.98 (-18.51)	11.14 (-20.41)	5.40 (-20.67)	0.22 (-1.12)	0.28 (-1.38)	0.28 (-1.38)	0.28 (-1.42)	137
		LLaVA-NeXT (FT)	<b>12.42 (-4.56)</b>	18.65 (-4.05)	<b>11.97 (-7.98)</b>	1.15 (-1.86)	1.11 (-2.23)	1.07 (-2.16)	1.00 (-2.10)	105
		Qwen-VL	11.02 (-10.09)	<b>18.70 (-8.32)</b>	10.89 (-10.89)	<b>1.18 (-0.42)</b>	<b>1.17 (-0.42)</b>	<b>1.14 (-0.41)</b>	<b>1.10 (-0.42)</b>	100
		Qwen-VL (FT)	0.00 (-21.12)	0.00 (-24.87)	0.00 (-21.95)	0.00 (-3.57)	0.00 (-3.83)	0.00 (-3.78)	0.00 (-3.68)	83
		mPLUG-Owl2	10.03 (-2.75)	17.28 (+0.20)	9.77 (-3.71)	1.09 (-0.98)	0.90 (-0.78)	0.84 (-0.74)	0.82 (-0.74)	55
Nl	En	LLaVA-NeXT	<b>15.81 (-10.68)</b>	<b>24.80 (-6.75)</b>	<b>21.13 (-4.94)</b>	0.11 (-1.24)	<b>2.04 (+0.38)</b>	<b>2.29 (+0.63)</b>	<b>1.83 (+0.13)</b>	223
		LLaVA-NeXT (FT)	9.92 (-7.06)	15.45 (-7.25)	9.73 (-10.22)	0.91 (-2.11)	0.98 (-2.36)	0.91 (-2.32)	0.88 (-2.23)	153
		Qwen-VL	11.65 (-9.46)	18.87 (-8.16)	13.19 (-8.59)	1.54 (-0.06)	1.47 (-0.12)	1.48 (-0.08)	1.44 (-0.08)	136
		Qwen-VL (FT)	10.35 (-10.76)	16.35 (-8.52)	10.79 (-11.16)	0.70 (-2.87)	1.13 (-2.70)	1.07 (-2.71)	1.04 (-2.64)	331
		mPLUG-Owl2	12.19 (-0.59)	20.44 (+3.36)	12.97 (-0.51)	<b>1.81 (-0.25)</b>	1.61 (-0.06)	1.56 (-0.03)	1.50 (-0.06)	82
Sv	En	LLaVA-NeXT	<b>18.70 (-7.79)</b>	<b>25.48 (-6.07)</b>	<b>18.98 (-7.09)</b>	1.79 (+0.44)	<b>1.86 (+0.20)</b>	<b>1.86 (+0.20)</b>	<b>1.80 (+0.10)</b>	246
		LLaVA-NeXT (FT)	9.30 (-7.68)	14.68 (-8.02)	9.12 (-10.83)	0.80 (-2.22)	0.76 (-2.57)	0.73 (-2.50)	0.71 (-2.40)	141
		Qwen-VL	11.77 (-9.33)	17.73 (-9.30)	13.03 (-8.75)	1.57 (-0.02)	1.36 (-0.23)	1.30 (-0.26)	1.25 (-0.27)	107
		Qwen-VL (FT)	11.00 (-10.11)	17.97 (-6.89)	9.83 (-12.12)	0.69 (-2.88)	0.74 (-3.10)	0.73 (-3.05)	0.63 (-3.05)	233
		mPLUG-Owl2	12.49 (-0.29)	18.51 (+1.43)	11.40 (-2.08)	<b>1.90 (-0.17)</b>	1.29 (-0.39)	1.29 (-0.30)	1.21 (-0.35)	81
Ru	En	LLaVA-NeXT	<b>18.31 (-8.18)</b>	<b>26.30 (-5.25)</b>	<b>18.43 (-7.64)</b>	<b>1.68 (+0.34)</b>	<b>1.64 (-0.01)</b>	<b>1.65 (-0.01)</b>	<b>1.59 (-0.11)</b>	241
		LLaVA-NeXT (FT)	9.61 (-7.37)	13.42 (-9.28)	8.54 (-11.41)	1.01 (-2.01)	0.97 (-2.36)	0.94 (-2.29)	0.91 (-2.20)	125
		Qwen-VL	13.36 (-7.75)	20.75 (-6.28)	13.97 (-7.81)	1.02 (-0.57)	1.13 (-0.45)	1.25 (-0.31)	1.22 (-0.30)	128
		Qwen-VL (FT)	9.66 (-11.45)	15.91 (-8.96)	9.12 (-12.83)	0.90 (-2.67)	0.87 (-2.96)	0.98 (-2.80)	0.87 (-2.81)	258
		mPLUG-Owl2	12.56 (-0.22)	19.45 (+2.37)	12.60 (-0.88)	1.60 (-0.47)	1.60 (-0.07)	1.51 (-0.07)	1.41 (-0.15)	96
Ja	En	LLaVA-NeXT	<b>15.36 (-11.13)</b>	<b>24.41 (-7.13)</b>	<b>16.18 (-9.89)</b>	<b>1.12 (-0.23)</b>	<b>1.15 (-0.51)</b>	<b>1.28 (-0.38)</b>	<b>1.11 (-0.59)</b>	208
		LLaVA-NeXT (FT)	7.69 (-9.28)	12.61 (-10.09)	8.29 (-11.66)	0.85 (-2.17)	0.54 (-2.80)	0.47 (-2.76)	0.45 (-2.65)	68
		Qwen-VL	10.32 (-10.78)	17.64 (-9.38)	9.75 (-12.03)	0.97 (-0.63)	0.75 (-0.83)	0.78 (-0.78)	0.76 (-0.76)	108
		Qwen-VL (FT)	0.73 (-20.38)	3.14 (-21.72)	0.00 (-21.95)	0.00 (-3.57)	0.00 (-3.83)	0.00 (-3.78)	0.00 (-3.68)	153
		mPLUG-Owl2	10.02 (-2.76)	17.10 (+0.02)	8.27 (-5.21)	1.11 (-0.96)	0.70 (-0.98)	0.71 (-0.88)	0.67 (-0.89)	76
Zh	En	LLaVA-NeXT	<b>13.44 (-2.42)</b>	<b>21.98 (+1.04)</b>	<b>12.83 (-3.67)</b>	<b>0.96 (-0.12)</b>	<b>1.43 (+0.19)</b>	<b>1.58 (+0.37)</b>	<b>1.42 (+0.17)</b>	168
		LLaVA-NeXT (FT)	6.71 (-4.93)	13.24 (-2.23)	6.44 (-7.69)	0.88 (-1.54)	0.59 (-2.01)	0.57 (-1.94)	0.54 (-1.89)	94
		Qwen-VL	8.98 (-3.15)	16.16 (-0.26)	9.65 (-2.54)	0.48 (-0.85)	0.45 (-0.72)	0.46 (-0.66)	0.40 (-0.67)	138
		Qwen-VL (FT)	8.90 (-4.87)	16.99 (-0.64)	8.79 (-4.64)	0.14 (-2.84)	0.10 (-3.00)	0.10 (-2.96)	0.09 (-2.91)	242
		mPLUG-Owl2	5.25 (-7.57)	11.81 (-5.33)	4.32 (-9.20)	0.15 (-1.93)	0.14 (-1.54)	0.15 (-1.45)	0.15 (-1.42)	34

Table 11: Results of LVLs in Alignment-10 Task, which the method is an instruction in languages other than English and output in English ({Lang}-{En}). Bold fonts indicate the best score for that language combination. The values are noted next to the output of the difference by the same model in the method with instruction and output in English ({En}-{En}). Red indicates a higher value than that method; blue indicates a lower value.

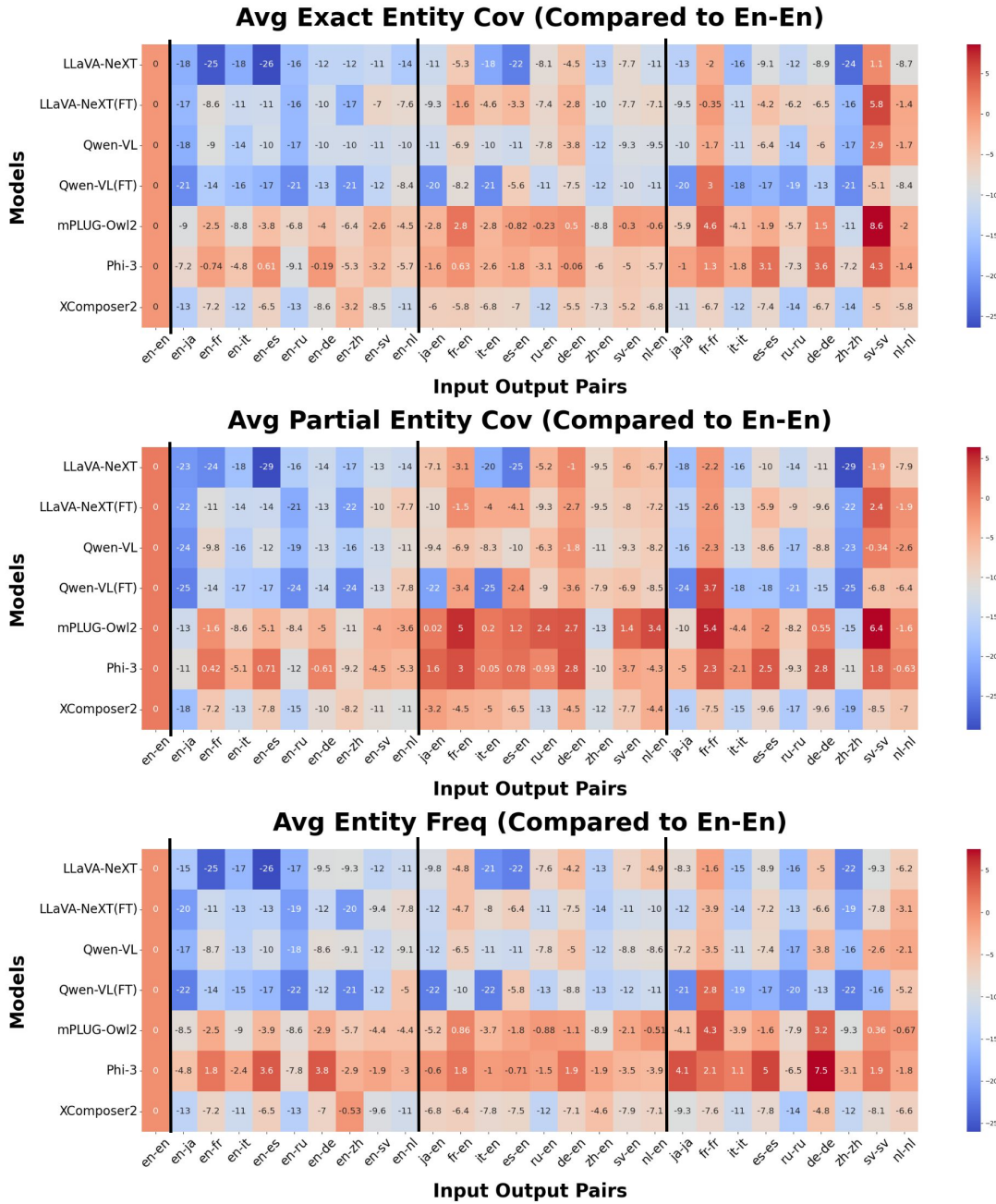


Figure 4: Visualization of Alignment-10 results in a heat map. We made the visualization based on when we had LVLMS give instructions and output in English.

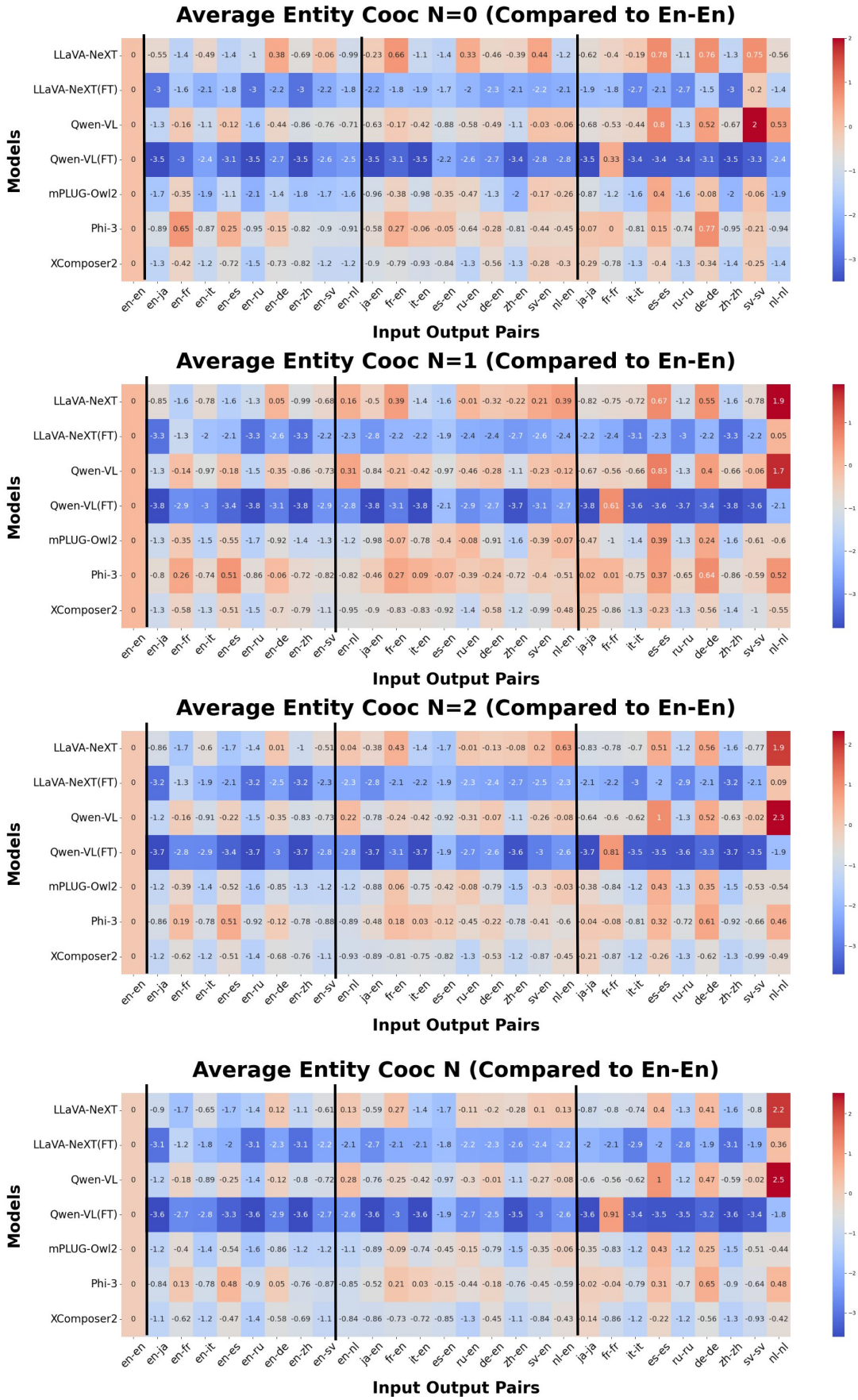


Figure 5: Visualization of Alignment-10 results in a heat map. We made the visualization based on when we had LVLMS give instructions and output in English.

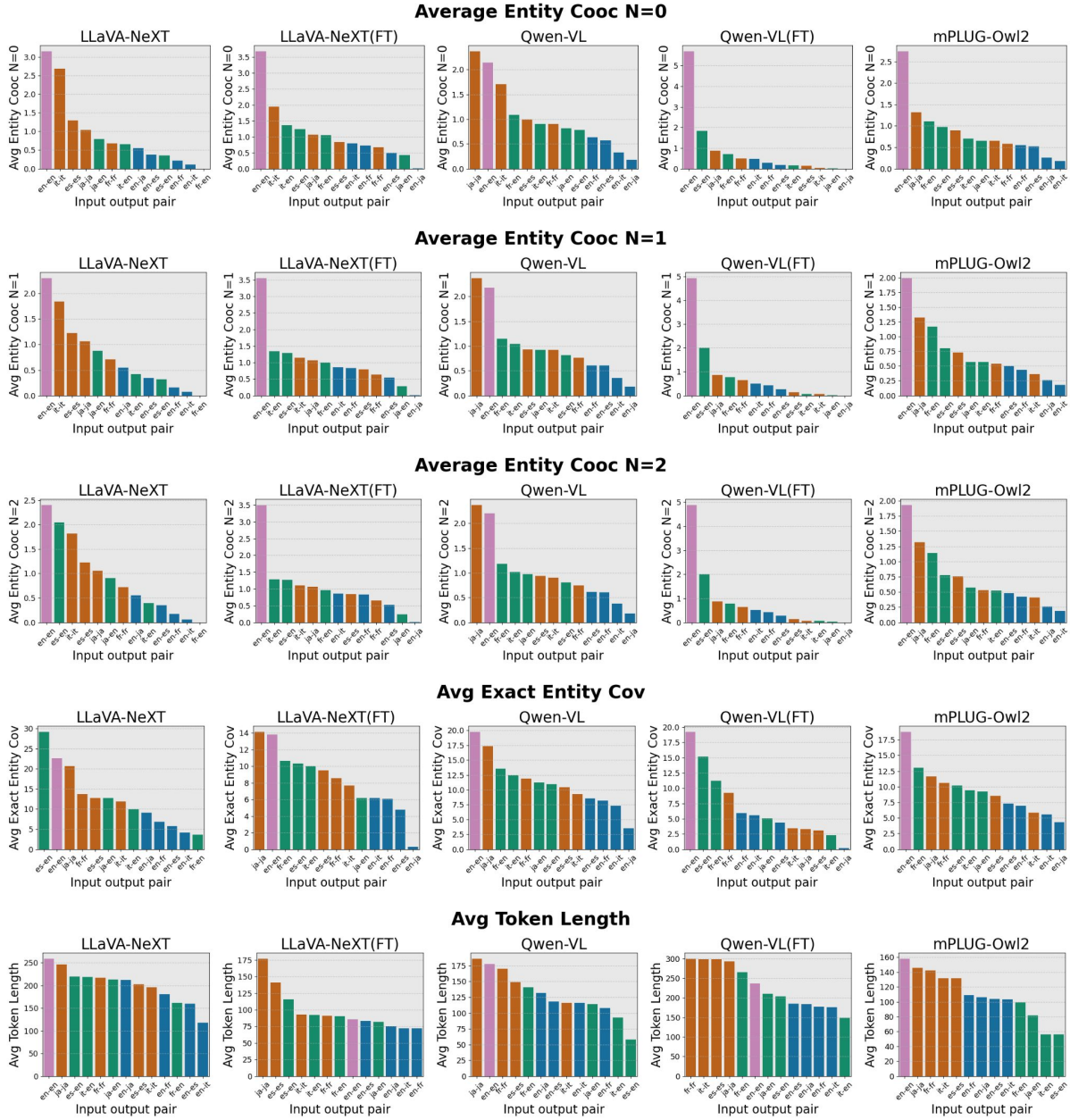


Figure 6: The rest of the results in the Alignment-5 task. From this figure, it can also be seen that the English instructions are optimal, even if the number of data is expanded. Purple bin indicates the method which is the instruction and the output in English ( $\{En\}-\{En\}$ ), Green bin indicates the instruction in languages other than English and the output in English ( $\{Lang\}-\{En\}$ ), Brown bin indicates the instruction and output in languages other than English ( $\{Lang\}-\{Lang\}$ ) and Blue bin indicates the instruction in English and the output in languages other than English ( $\{En\}-\{Lang\}$ ).

Input	Output	LVL	BLEU	ROUGE			BertScore
				1	2	L	
En	En	LLaVA-NeXT	0.01	0.24	0.05	0.15	0.82
		LLaVA-NeXT (FT)	<b>0.07</b>	0.28	<b>0.13</b>	<b>0.22</b>	<b>0.85</b>
		Qwen-VL	0.01	0.22	0.05	0.14	0.82
		Qwen-VL (FT)	0.06	<b>0.28</b>	0.12	0.22	0.84
		mPLUG-Owl2	0.01	0.24	0.05	0.15	0.82
		Phi-3	0.01	0.20	0.04	0.12	0.82
		XComposer2	0.01	0.24	0.05	0.14	0.82
En	Es	LLaVA-NeXT	<b>0.01 (-0.00)</b>	<b>0.28 (+0.04)</b>	<b>0.06 (+0.01)</b>	<b>0.16 (+0.01)</b>	<b>0.81 (-0.01)</b>
		LLaVA-NeXT (FT)	0.00 (-0.07)	0.05 (-0.23)	0.01 (-0.12)	0.04 (-0.18)	0.78 (-0.07)
		Qwen-VL	0.00 (-0.01)	0.20 (-0.03)	0.04 (-0.01)	0.12 (-0.02)	0.80 (-0.02)
		Qwen-VL (FT)	0.00 (-0.06)	0.03 (-0.25)	0.00 (-0.11)	0.03 (-0.19)	0.77 (-0.07)
		mPLUG-Owl2	0.00 (-0.01)	0.22 (-0.03)	0.04 (-0.01)	0.13 (-0.02)	0.80 (-0.02)
		Phi-3	0.00 (-0.00)	0.21 (+0.01)	0.04 (+0.00)	0.13 (+0.00)	0.79 (-0.02)
		XComposer2	0.00 (-0.01)	0.18 (-0.06)	0.04 (-0.02)	0.11 (-0.03)	0.80 (-0.02)
En	Fr	LLaVA-NeXT	<b>0.00 (-0.01)</b>	<b>0.20 (-0.04)</b>	<b>0.04 (-0.02)</b>	<b>0.12 (-0.03)</b>	<b>0.79 (-0.02)</b>
		LLaVA-NeXT (FT)	0.00 (-0.07)	0.06 (-0.21)	0.02 (-0.11)	0.05 (-0.17)	0.78 (-0.06)
		Qwen-VL	0.00 (-0.01)	0.15 (-0.08)	0.03 (-0.02)	0.09 (-0.05)	0.79 (-0.03)
		Qwen-VL (FT)	0.00 (-0.06)	0.03 (-0.25)	0.00 (-0.11)	0.03 (-0.19)	0.77 (-0.07)
		mPLUG-Owl2	0.00 (-0.01)	0.16 (-0.08)	0.03 (-0.02)	0.10 (-0.05)	0.79 (-0.03)
		Phi-3	0.00 (-0.00)	0.15 (-0.04)	0.02 (-0.01)	0.09 (-0.03)	0.78 (-0.03)
		XComposer2	0.00 (-0.01)	0.03 (-0.21)	0.01 (-0.05)	0.03 (-0.12)	0.78 (-0.04)
En	De	LLaVA-NeXT	<b>0.00 (-0.01)</b>	<b>0.20 (-0.05)</b>	<b>0.03 (-0.02)</b>	<b>0.11 (-0.03)</b>	<b>0.80 (-0.02)</b>
		LLaVA-NeXT (FT)	0.00 (-0.07)	0.04 (-0.23)	0.01 (-0.12)	0.03 (-0.19)	0.76 (-0.08)
		Qwen-VL	0.00 (-0.01)	0.14 (-0.08)	0.02 (-0.03)	0.09 (-0.06)	0.79 (-0.03)
		Qwen-VL (FT)	0.00 (-0.06)	0.03 (-0.25)	0.00 (-0.11)	0.03 (-0.19)	0.76 (-0.08)
		mPLUG-Owl2	0.00 (-0.01)	0.14 (-0.10)	0.02 (-0.03)	0.09 (-0.07)	0.79 (-0.03)
		Phi-3	0.00 (-0.00)	0.14 (-0.05)	0.02 (-0.02)	0.09 (-0.03)	0.78 (-0.03)
		XComposer2	0.00 (-0.01)	0.14 (-0.10)	0.02 (-0.03)	0.09 (-0.06)	0.79 (-0.03)
En	It	LLaVA-NeXT	<b>0.00 (-0.01)</b>	<b>0.19 (-0.05)</b>	<b>0.02 (-0.03)</b>	<b>0.11 (-0.04)</b>	<b>0.80 (-0.01)</b>
		LLaVA-NeXT (FT)	0.00 (-0.07)	0.04 (-0.23)	0.01 (-0.12)	0.04 (-0.19)	0.77 (-0.08)
		Qwen-VL	0.00 (-0.01)	0.14 (-0.08)	0.02 (-0.03)	0.09 (-0.06)	0.80 (-0.02)
		Qwen-VL (FT)	0.00 (-0.06)	0.04 (-0.24)	0.01 (-0.11)	0.04 (-0.18)	0.76 (-0.07)
		mPLUG-Owl2	0.00 (-0.01)	0.14 (-0.10)	0.02 (-0.04)	0.09 (-0.07)	0.80 (-0.02)
		Phi-3	0.00 (-0.00)	0.10 (-0.09)	0.01 (-0.03)	0.07 (-0.05)	0.78 (-0.03)
		XComposer2	0.00 (-0.01)	0.10 (-0.14)	0.01 (-0.04)	0.07 (-0.07)	0.80 (-0.02)
En	NI	LLaVA-NeXT	0.00 (-0.01)	<b>0.23 (-0.01)</b>	<b>0.04 (-0.01)</b>	<b>0.15 (-0.00)</b>	<b>0.81 (-0.01)</b>
		LLaVA-NeXT (FT)	<b>0.01 (-0.06)</b>	0.12 (-0.15)	0.03 (-0.10)	0.09 (-0.13)	0.78 (-0.07)
		Qwen-VL	0.00 (-0.01)	0.20 (-0.03)	0.04 (-0.01)	0.13 (-0.01)	0.80 (-0.02)
		Qwen-VL (FT)	0.00 (-0.06)	0.06 (-0.23)	0.01 (-0.11)	0.05 (-0.17)	0.76 (-0.08)
		mPLUG-Owl2	0.00 (-0.01)	0.17 (-0.07)	0.03 (-0.02)	0.11 (-0.04)	0.80 (-0.03)
		Phi-3	0.00 (-0.00)	0.10 (-0.10)	0.01 (-0.02)	0.08 (-0.05)	0.77 (-0.05)
		XComposer2	0.00 (-0.01)	0.15 (-0.09)	0.03 (-0.03)	0.11 (-0.04)	0.80 (-0.02)
En	Sv	LLaVA-NeXT	<b>0.00 (-0.01)</b>	<b>0.21 (-0.04)</b>	<b>0.04 (-0.02)</b>	<b>0.12 (-0.02)</b>	<b>0.81 (-0.01)</b>
		LLaVA-NeXT (FT)	0.00 (-0.07)	0.08 (-0.20)	0.02 (-0.11)	0.06 (-0.16)	0.78 (-0.07)
		Qwen-VL	0.00 (-0.01)	0.15 (-0.07)	0.02 (-0.03)	0.09 (-0.05)	0.79 (-0.03)
		Qwen-VL (FT)	0.00 (-0.06)	0.03 (-0.26)	0.01 (-0.11)	0.02 (-0.20)	0.76 (-0.08)
		mPLUG-Owl2	0.00 (-0.01)	0.14 (-0.11)	0.02 (-0.03)	0.09 (-0.07)	0.80 (-0.03)
		Phi-3	0.00 (-0.01)	0.05 (-0.14)	0.01 (-0.03)	0.04 (-0.08)	0.76 (-0.05)
		XComposer2	0.00 (-0.01)	0.11 (-0.13)	0.02 (-0.04)	0.08 (-0.07)	0.79 (-0.03)
En	Ru	LLaVA-NeXT	<b>0.00 (-0.01)</b>	<b>0.03 (-0.22)</b>	<b>0.00 (-0.05)</b>	<b>0.02 (-0.12)</b>	<b>0.89 (+0.07)</b>
		LLaVA-NeXT (FT)	0.00 (-0.07)	0.01 (-0.27)	0.00 (-0.13)	0.01 (-0.21)	0.72 (-0.13)
		Qwen-VL	0.00 (-0.01)	0.02 (-0.21)	0.00 (-0.05)	0.02 (-0.13)	0.85 (+0.03)
		Qwen-VL (FT)	0.00 (-0.06)	0.01 (-0.27)	0.00 (-0.12)	0.01 (-0.21)	0.70 (-0.14)
		mPLUG-Owl2	0.00 (-0.01)	0.01 (-0.23)	0.00 (-0.05)	0.01 (-0.14)	0.86 (+0.04)
		Phi-3	0.00 (-0.01)	0.01 (-0.19)	0.00 (-0.04)	0.01 (-0.12)	0.71 (-0.10)
		XComposer2	0.00 (-0.01)	0.02 (-0.22)	0.00 (-0.05)	0.02 (-0.13)	0.87 (+0.05)
En	Ja	LLaVA-NeXT	<b>0.01 (-0.00)</b>	<b>0.03 (-0.21)</b>	<b>0.01 (-0.05)</b>	<b>0.03 (-0.11)</b>	<b>0.84 (+0.03)</b>
		LLaVA-NeXT (FT)	0.00 (-0.07)	0.01 (-0.26)	0.00 (-0.13)	0.01 (-0.21)	0.73 (-0.12)
		Qwen-VL	0.00 (-0.01)	0.02 (-0.20)	0.00 (-0.05)	0.02 (-0.13)	0.83 (+0.00)
		Qwen-VL (FT)	0.00 (-0.06)	0.01 (-0.27)	0.00 (-0.12)	0.01 (-0.21)	0.72 (-0.12)
		mPLUG-Owl2	0.00 (-0.01)	0.02 (-0.23)	0.00 (-0.05)	0.02 (-0.14)	0.83 (+0.01)
		Phi-3	0.00 (-0.00)	0.02 (-0.18)	0.00 (-0.03)	0.02 (-0.11)	0.82 (+0.01)
		XComposer2	0.00 (-0.01)	0.02 (-0.22)	0.00 (-0.05)	0.02 (-0.12)	0.83 (+0.01)
En	Zh	LLaVA-NeXT	<b>0.00 (-0.01)</b>	<b>0.03 (-0.21)</b>	<b>0.01 (-0.05)</b>	<b>0.03 (-0.12)</b>	<b>0.83 (+0.01)</b>
		LLaVA-NeXT (FT)	0.00 (-0.07)	0.02 (-0.25)	0.01 (-0.12)	0.02 (-0.20)	0.73 (-0.12)
		Qwen-VL	0.00 (-0.01)	0.03 (-0.19)	0.01 (-0.04)	0.03 (-0.11)	0.83 (+0.01)
		Qwen-VL (FT)	0.00 (-0.06)	0.02 (-0.26)	0.00 (-0.11)	0.02 (-0.20)	0.72 (-0.12)
		mPLUG-Owl2	0.00 (-0.01)	0.02 (-0.22)	0.01 (-0.05)	0.02 (-0.14)	0.83 (+0.00)
		Phi-3	0.00 (-0.00)	0.02 (-0.18)	0.01 (-0.03)	0.02 (-0.11)	0.81 (-0.00)
		XComposer2	0.00 (-0.01)	0.03 (-0.21)	0.01 (-0.04)	0.03 (-0.12)	0.83 (+0.01)

Table 12: Other metrics results of LVLs in Full Task. Bold fonts indicate the best score for that language combination. We also measured outputs with existing NLG evaluation methods, BLEU (Papineni et al., 2002), ROUGE (Lin, 2004), and BertScore (Zhang et al., 2019).

861  
862  
863  
864  
865

## I Details of Each Language Templates

As indicated in Table 1, we created Templates for ten languages. Ten language templates are shown below. These templates were modified by nine native speakers of the country’s first language, who were asked to modify the sentences to have the same nuance and level of detail as in English. We described these templates from Table 13 to Table 22.

Language	Type	Template
English	<b>Template 1</b>	
	Section	Focus on {title} and explain the {section}.
	Subsection	In the context of {title}, explain the {subsection} and the {section}.
	Sub subsection	Focusing on the {section} of {title}, explain the {subsubsection} about the {subsection}.
	<b>Template 2</b>	
	Section	Explain the {section} of this artwork, {title}.
	Subsection	Explain the {subsection} about the {section} of this artwork, {title}.
	Sub subsection	Explain the {subsubsection} about the {section} of the {section} in this work, {title}.
	<b>Template 3</b>	
	Section	How does {title} explain its {section}?
	Subsection	In {title}, how is the {subsection} of the {section} explained?
	Sub subsection	Regarding {title}, how does the {section}'s {subsection} incorporate the {subsubsection}.
<b>Template 4</b>		
Section	In {title}, how is the {section} discussed?	
Subsection	Describe the characteristics of the {subsection} in {title}'s {section}.	
Sub subsection	When looking at the {section} of {title}, how do you discuss its {subsection}'s {subsubsection}?	

Table 13: Prompt Templates in English

Language	Type	Template
Japanese	<b>Template 1</b>	
	Section	{title}に焦点を当てて、その{section}を説明してください。
	Subsection	{title}の文脈で、{subsection}と{section}を説明してください。
	Sub subsection	{title}の{section}に焦点を当てて、{subsection}についての{subsubsection}を説明してください。
	<b>Template 2</b>	
	Section	{title}の{section}を説明してください。
	Subsection	{title}の{section}に関する{subsection}を説明してください。
	Sub subsection	{title}の{section}の{subsection}に関する{subsubsection}を説明してください。
	<b>Template 3</b>	
	Section	{title}はどのように{section}を説明していますか？
	Subsection	{title}では、どのように{section}の{subsection}が説明されていますか？
	Sub subsection	{title}に関して、{section}の{subsection}は{subsubsection}をどのように取り入れていますか？
<b>Template 4</b>		
Section	{title}に関して、どのように{section}が議論されていますか？	
Subsection	{title}の{section}における{subsection}の特徴を説明してください。	
Sub subsection	{title}の{section}について見たとき、その{subsection}の{subsubsection}をどのように議論しますか？	

Table 14: Prompt Templates in Japanese



Language	Type	Template
Spanish	<b>Template 1</b>	
	Section	Concéntrate en {title} y explora la {section}.
	Subsection	En el contexto de {title}, explora la {subsection} y {section}.
	Sub subsection	Concentrándote en la {section} de {title}, explora la {subsubsection} sobre la {subsection}.
	<b>Template 2</b>	
	Section	Explora la {section} de esta obra de arte, {title}.
	Subsection	Explora la {subsection} sobre la {section} de esta obra de arte, {title}.
	Sub subsection	Explora la {subsubsection} sobre {subsection} de la {section} en esta obra de arte, {title}.
	<b>Template 3</b>	
	Section	¿Cómo aclara {title} su {section}?
	Subsection	En {title}, ¿cómo se aclara la {subsection} de la {section}?
	Sub subsection	Con respecto a {title}, ¿cómo la {subsection} de la {section} incorpora a la {subsubsection}?
<b>Template 4</b>		
Section	En {title}, ¿cómo se discute la {section}?	
Subsection	Describe las características de la {subsection} en la {section} de {title}.	
Sub subsection	Al observar la {section} de {title}, ¿cómo discutes la {subsubsection} de su {subsection}?	

Table 15: Prompt Templates in Spanish

Language	Type	Template
Italian	<b>Template 1</b>	
	Section	Concentrati su {title} ed esplora la {section}.
	Subsection	Nel contesto di {title}, esplora la {subsection} e la {section}.
	Sub subsection	Concentrandosi sulla {section} di {title}, esplora la {subsubsection} sulla {subsection}.
	<b>Template 2</b>	
	Section	Esplora la {section} di questa opera d'arte, {title}.
	Subsection	Esplora la {subsection} sulla {section} di questa opera d'arte, {title}.
	Sub subsection	Esplora la {subsubsection} sulla {subsection} della {section} in questa opera, {title}.
	<b>Template 3</b>	
	Section	Come chiarisce {title} la sua {section}?
	Subsection	In {title}, come viene chiarita la {subsection} della {section}?
	Sub subsection	Per quanto riguarda {title}, come la {section} incorpora la {subsection} con la {subsubsection}?
<b>Template 4</b>		
Section	Come viene discussa la {section} in {title}?	
Subsection	Descrivi le caratteristiche della {subsection} nella {section} di {title}.	
Sub subsection	Osservando la {section} di {title}, come discuti la {subsection} della {subsubsection}?	

Table 16: Prompt Templates in Italian

Language	Type	Template
French	<b>Template 1</b>	
	Section	Concentrez-vous sur {title} et expliquez la {section}.
	Subsection	Dans le contexte de {title}, expliquez la {subsection} et la {section}.
	Sub subsection	En vous concentrant sur la {section} de {title}, expliquez la {subsubsection} concernant la {subsection}.
	<b>Template 2</b>	
	Section	Expliquer la {section} de cette œuvre d'art, {title}.
	Subsection	Expliquer la {subsection} concernant la {section} de cette œuvre d'art, {title}.
	Sub subsection	Expliquer la {subsubsection} concernant la {subsection} de la {section} dans cette œuvre, {title}.
	<b>Template 3</b>	
	Section	Comment {title} explique-t-il sa {section}?
	Subsection	Dans {title}, comment la {subsection} de la {section} est-elle expliquée?
	Sub subsection	Concernant {title}, comment la {subsection} de la {section} intègre-t-elle la {subsubsection}?
<b>Template 4</b>		
Section	Dans {title}, comment est discutée la {section}?	
Subsection	Décrivez les caractéristiques de la {subsection} dans la {section} de {title}.	
Sub subsection	En examinant la {section} de {title}, comment discutez-vous la {subsubsection} de la {subsection}?	

Table 17: Prompt Templates in French

Language	Type	Template
Chinese (Simplified)	<b>Template 1</b>	
	Section	专注于{title}并探索{section}。
	Subsection	在{title}的背景下，探索{subsection}和{section}。
	Sub subsection	专注于{title}的{section}，探索关于{subsection}的{subsubsection}。
	<b>Template 2</b>	
	Section	探索艺术作品{title}的{section}。
	Subsection	探索艺术作品{title}中关于{section}的{subsection}。
	Sub subsection	探索作品{title}中{section}的{subsection}的{subsubsection}。
	<b>Template 3</b>	
	Section	{title}是如何阐明其{section}的？
	Subsection	在{title}中，{section}的{subsection}是如何被阐明的？
	Sub subsection	关于{title}，{section}的{subsection}是如何结合{subsubsection}的？
<b>Template 4</b>		
Section	在{title}中，{section}是如何被讨论的？	
Subsection	描述{title}的{section}中{subsection}的特点。	
Sub subsection	在查看{title}的{section}时，你如何讨论其{subsection}的{subsubsection}？	

Table 18: Prompt Templates in Chinese (Simplified)

Language	Type	Template
Dutch	<b>Template 1</b>	
	Section	Focus op {title} en leg de {section} uit.
	Subsection	In de context van {title}, leg de {subsection} en de {section} uit.
	Sub subsection	Gefocust op de {section} van {title}, leg de {subsubsection} over de {subsection} uit.
	<b>Template 2</b>	
	Section	Leg de {section} van dit kunstwerk uit, {title}.
	Subsection	Leg de {subsection} over de {section} van dit kunstwerk uit, {title}.
	Sub subsection	Leg de {subsubsection} over de {section} van de {section} in dit werk uit, {title}.
	<b>Template 3</b>	
	Section	Hoe verduidelijkt {title} zijn {section}?
	Subsection	Hoe wordt in {title} de {subsection} van de {section} verduidelijkt?
	Sub subsection	Met betrekking tot {title}, hoe incorporeert de {section}'s {subsection} de {subsubsection}?
<b>Template 4</b>		
Section	Hoe wordt de {section} besproken in {title}?	
Subsection	Beschrijf de kenmerken van de {subsection} in de {section} van {title}.	
Sub subsection	Wanneer je kijkt naar de {section} van {title}, hoe bespreek je de {subsection}'s {subsubsection}?	

Table 19: Prompt Templates in Dutch

Language	Type	Template
Swedish	<b>Template 1</b>	
	Section	Fokusera på {title} och förklara {section}.
	Subsection	I samband med {title}, förklara {subsection} och {section}.
	Sub subsection	Med fokus på {section} i {title}, förklara {subsubsection} om {subsection}.
	<b>Template 2</b>	
	Section	Förklara {section} i detta konstverk, {title}.
	Subsection	Förklara {subsection} om {section} i detta konstverk, {title}.
	Sub subsection	Förklara {subsubsection} om {subsection} av {section} i detta verk, {title}.
	<b>Template 3</b>	
	Section	Hur förklarar {title} sitt {section}?
	Subsection	Hur förklaras {subsection} av {section} i {title}?
	Sub subsection	När det gäller {title}, hur innehåller {section}'s {subsection} {subsubsection}?
<b>Template 4</b>		
Section	I {title}, hur diskuteras {section}?	
Subsection	Beskriv egenskaperna hos {subsection} i {title}'s {section}.	
Sub subsection	När du tittar på {section} i {title}, hur diskuteras du dess {subsection}'s {subsubsection}?	

Table 20: Prompt Templates in Swedish

Language	Type	Template	
German	Section	Fokussiere dich auf {title} und erkunde erkläre die {section}.	
	Subsection	Im Kontext von {title}, erkunde erkläre die {subsection} und die {section}.	
	Sub subsection	Mit Fokus auf die {section} von {title}, erkunde erkläre die {subsubsection} über die {subsection}.	
	<b>Template 2</b>		
	Section	Erkunde Erkläre die {section} dieses Kunstwerks, {title}.	
	Subsection	Erkunde Erkläre die {subsection} über die {section} dieses Kunstwerks, {title}.	
	Sub subsection	Erkunde Erkläre die {subsubsection} über die {subsection} der {section} in diesem Werk, {title}.	
	<b>Template 3</b>		
	Section	Wie erläutert {title} seine {section}?	
	Subsection	In {title}, wie wird die {subsection} der {section} erläutert?	
	Sub subsection	Bezüglich {title}, wie integriert die {subsection} der {section} die {subsubsection}?	
	<b>Template 4</b>		
	Section	Wie wird die {section} in {title} diskutiert?	
	Subsection	Beschreibe die Merkmale der {subsection} in der {title}'s {section}.	
	Sub subsection	Wenn du die {section} von {title} betrachtest, wie diskutierst du die {subsection}'s {subsubsection} von der {subsection}?	

Table 21: Prompt Templates in German

Language	Type	Template	
Russian	Section	Сосредоточьтесь на {title} и объясните {section}.	
	Subsection	В контексте {title} объясните {subsection} и {section}.	
	Sub subsection	Сосредоточившись на {section} в {title}, объясните {subsubsection} о {subsection}.	
	<b>Template 2</b>		
	Section	Объясните {section} этого произведения искусства, {title}.	
	Subsection	Объясните {subsection} о {section} этого произведения искусства, {title}.	
	Sub subsection	Объясните {subsubsection} о {section} в {section} этого произведения, {title}.	
	<b>Template 3</b>		
	Section	Как {title} объясняет свой/свою {section}?	
	Subsection	Как объясняется в {title} {subsection} в {section}?	
	Sub subsection	Что касается {title}, как {section} в {subsection} включает {subsubsection}?	
	<b>Template 4</b>		
	Section	Как обсуждается {section} в {title}?	
	Subsection	Опишите черты {subsection} в {section} в {title}.	
	Sub subsection	Когда вы рассматриваете {section} в {title}, как вы обсуждаете {subsubsection} в {section}?	

Table 22: Prompt Templates in Russian