

Do Vision-Language Models Learn In Context? Not So Fast

Anonymous ACL submission

Abstract

In-context learning enables Large Language Models (LLMs) to learn tasks from demonstration examples without parameter updates. While this ability has been extensively studied in LLMs, its effectiveness in Vision-Language Models (VLMs) remains underexplored. Existing research primarily focuses on a few models trained on interleaved image-text datasets and often overlooks image captioning in their analysis. In this work, we systematically analyze in-context learning in VLMs, evaluating six models across four architectures on three image captioning and four visual question answering benchmarks. We investigate the influence of prompt design, demonstration selection, model architecture, and training strategies. We also extend our analysis beyond models trained on interleaved datasets to include those trained on image-text pairs, often considered incapable of in-context learning. Our findings show that VLMs still struggle to leverage contextual information to adapt their outputs. However, detailed prompts specifying the task and structure of demonstrations improve performance more than simply concatenating examples. Additionally, while instruction-tuning enhances comprehension of detailed instructions, it reduces reliance on contextual examples and may hinder models' in-context learning capacity. Moreover, VLMs with advanced modality projectors can achieve competitive in-context learning performance even trained on image-text pairs.

1 Introduction

In recent years, Large Language Models (LLMs) have attracted significant attention for their notable performance across a wide range of Natural Language Processing tasks. As these models scale, in-context learning emerges as a new ability that allows LLMs to learn tasks given only a few examples through demonstrations (Brown et al., 2020; Wei et al., 2022). In this paradigm, before being asked to perform a task, the model is given a set

of demonstrations, i.e., input-output examples, illustrating how to do it. Unlike supervised learning, in-context learning does not involve further parameter updates. Instead, the model should learn from analogy (Dong et al., 2024).

Despite the advancements, LLMs remain restricted to processing text-based data. They cannot handle other modalities such as image, audio, or video directly. However, the capacity to handle multimodal information contributes to knowledge acquisition and interaction with the real world. To bridge this gap, Vision-Language Models (VLMs) arise as a proposal to extend LLMs' capabilities to process visual information. Although in-context learning has been extensively studied in LLMs from various perspectives (Dong et al., 2024), relatively few works have explored this ability in VLMs (Baldassini et al., 2024; Qin et al., 2024; Yang et al., 2024). Moreover, they primarily evaluate a limited number of models trained on interleaved image-text datasets and focus predominantly on tasks such as Visual Question Answering (VQA) and image classification, often overlooking the task of image captioning.

In this paper, we systematically analyze in-context learning in VLMs, evaluating six models from four architectures across three image captioning and four VQA benchmarks. Specifically, we investigate how prompt construction, demonstration selection, and design decisions on model architecture and training impact in-context learning ability. Also, besides models trained on interleaved image-text datasets (OpenFlamingo (Awadalla et al., 2023) and Idefics2 (Laurençon et al., 2024)), we extend our analysis to include InstructBLIP (Dai et al., 2024) and LLaVA v1.5 (Liu et al., 2023), both originally designed to process a single image-text pair. To do so, we adapted their modality alignment method for multiple input images. We conduct all experiments in a controlled environment for fair comparisons, evaluating models under iden-

084 tical conditions.

085 Our main findings are as follows: (1) Overall, 135
086 evaluated VLMs struggle to leverage the contextual 136
087 information to adapt the output. However, using 137
088 detailed prompts that explicitly define the task 138
089 and the structure of demonstration examples proves 139
090 more effective than simply concatenating examples. 140
091 Additionally, increasing the number of demonstra- 141
092 tions does not necessarily improve performance. 142
093 (2) While instruction-tuning enhances the model’s 143
094 ability to comprehend detailed instructions, it may 144
095 reduce its reliance on contextual examples. Con- 145
096 versely, training on interleaved image-text datasets 146
097 improves the model’s use of contextual informa- 147
098 tion. (3) VLMs with advanced modality projectors 148
099 achieve competitive in-context learning abilities 149
100 even when trained on single image-text pairs, offer- 150
101 ing a cost-efficient alternative to models trained on 151
102 large-scale interleaved datasets. In contrast, mod- 152
103 els with poor visual-text alignment – relying on 153
104 long token sequences to represent images – show 154
105 weaker in-context learning capabilities. These find- 155
106 ings highlight crucial limitations in current VLMs 156
107 that should be addressed to enhance their in-context
108 learning ability.

109 2 Related Work

110 **Vision-Language Models.** VLMs excel in vi- 161
111 sion-language tasks due to pre-trained visual en- 162
112 coders and LLMs (Yin et al., 2024; Zhang et al., 163
113 2024). They comprise three key components: a 164
114 visual encoder for image features, an LLM for text 165
115 generation, and a modality projector to align visual 166
116 and textual data, bridging the modality gap. 167

117 Various approaches have been explored for the 168
118 modality projector, including linear layers and 169
119 multi-layer perceptrons (MLPs) (Koh et al., 2023; 170
120 Liu et al., 2023; Shukor et al., 2023; Su et al., 2023; 171
121 Lin et al., 2024; Liu et al., 2024), which, despite 172
122 the low training costs, can lead to long sequences 173
123 of tokens thereby increasing the inference costs. 174
124 Pooling strategies help mitigate this issue (Cha 175
125 et al., 2024; Sun et al., 2024; Hu et al., 2024). Ad- 176
126 vanced methods like Q-Former (Li et al., 2023) 177
127 improve alignment between frozen visual encoders 178
128 and LLMs (Zhu et al., 2024a; Dai et al., 2024; Gei- 179
129 gle et al., 2024). Another alternative is interleaved 180
130 cross-attention layers (Alayrac et al., 2022; Lau- 181
131 rençon et al., 2023; Xue et al., 2024), where the 182
132 LLM directly attends to visual features but signifi- 183
133 cantly increases the number of trainable parameters,

as pointed out by Laurençon et al. (2024).

134 Training these models typically involves pre- 135
136 training the modality projector on large-scale 137
138 image-text datasets while keeping the visual en- 139
140 coder and LLM frozen for feature alignment. Sub- 141
142 sequently, the LLM can be fine-tuned alongside 143
144 the modality projector on instruction-following 144
145 datasets to improve zero-shot generalization. Most 145
146 works (Dai et al., 2024; Liu et al., 2024, 2023; 146
147 Zhu et al., 2024a; Hu et al., 2024) train on a 147
148 mixture of image captioning (Lin et al., 2014; Li 148
149 et al., 2022; Sharma et al., 2018), VQA (Goyal 149
150 et al., 2017; Schwenk et al., 2022; Marino et al., 150
151 2019), and instruction-following (Liu et al., 2024) 151
152 datasets. Some models, such as Flamingo (Alayrac 152
153 et al., 2022), Idefics (Laurençon et al., 2023; 153
154 Laurençon et al., 2024; Laurençon et al., 2024), 154
155 VILA (Lin et al., 2024), MMICL (Zhao et al., 155
156 2024), MM1 (McKinzie et al., 2025), and xGen- 156
157 MM (BLIP-3) (Xue et al., 2024), are trained on 157
158 interleaved image-text datasets (Laurençon et al., 158
159 2023; Zhu et al., 2024b) to further enhance multi- 159
160 modal reasoning capabilities. 160

161 **In-Context Learning in VLMs.** In-context 161
162 learning has been extensively studied in LLMs, 162
163 but this paradigm remains underexplored in VLMs. 163
164 Recent studies investigate different factors that af- 164
165 fect the in-context learning ability of VLMs, in- 165
166 cluding modality importance, recency bias, demon- 166
167 stration retrieval, and ordering strategies. However, 167
168 these studies primarily evaluate a limited number of 168
169 models trained on interleaved image-text datasets, 169
170 mainly in VQA and image classification tasks, of- 170
171 ten neglecting image captioning. 171

172 Yang et al. (2024) investigated in-context learn- 172
173 ing for image captioning, analyzing different 173
174 demonstration retrieval and caption assignment 174
175 methods. Their findings suggest that when demon- 175
176 stration images are similar to the query image, 176
177 VLMs may leverage in-context captions as short- 177
178 cuts to generate a new one rather than learning the 178
179 captioning task. 179

180 Chen et al. (2024) and Baldassini et al. (2024) 180
181 showed that textual information is more critical 181
182 than visual information in the demonstrations for 182
183 in-context learning in VLMs. Removing images 183
184 causes a minor performance drop, while corrupt- 184
185 ing textual descriptions leads to a significant per- 185
186 formance decline, indicating that VLMs heavily 186
187 rely on textual cues even when processing multi- 187
188 modal demonstrations. 188

Beyond modality importance, Baldassini et al. (2024) explored recency bias in VLMs. They showed that models tend to replicate outputs from the most recent demonstrations, even when earlier demonstrations are more semantically relevant. Qin et al. (2024) further studied demonstration retrieval and ordering, revealing that multi-modal retrieval methods outperform single-modal approaches. They showed that the order of modalities within each demonstration can significantly influence model performance more than the arrangement of demonstrations themselves. Also, unlike traditional text-based in-context learning, where increasing the number of demonstrations improves performance, they found no significant performance gains when providing more demonstrations.

In contrast to previous studies, we systematically analyze the in-context learning ability of six models from four distinct architectures across three image captioning and four VQA benchmarks. We investigate the impact of prompt construction, demonstration selection, model architecture, and training choices. Additionally, previous works have explored models that support interleaved image-text inputs, in contrast, we modify InstructBLIP (Dai et al., 2024) and LLaVA v1.5 (Liu et al., 2023) to extend our analysis to models that originally designed for single image-text pairs.

3 Methodology

3.1 Experimental Setup

Models. We analyze four distinct families of VLMs: InstructBLIP (Dai et al., 2024), LLaVA v1.5 (Liu et al., 2023), OpenFlamingo (Awadalla et al., 2023), and Idefics2 (Laurençon et al., 2024). These families were selected to systematically explore how various design choices – such as bridging the modality gap and different training methods – affect the in-context learning capabilities of VLMs.

We use model checkpoints with parameter sizes ranging from 4B to 9B for a fair comparison across similar scenarios. Specifically, for InstructBLIP, we evaluate two checkpoints with different LLMs: InstructBLIP FlanT5-XL and InstructBLIP Vicuna 7B. For the other families, we assess LLaVA v1.5 7B, OpenFlamingo 9B, and two checkpoints of Idefics2 – before and after the instruction-tuning phase – namely, Idefics2 (Base) and Idefics2 (IT)¹.

¹Salesforce/instructblip-flan-t5-xl
Salesforce/instructblip-vicuna-7b
llava-hf/llava-1.5-7b-hf

Datasets & Metrics. We evaluate the models using different benchmarks proposed for image captioning and VQA. For image captioning, we use MS COCO (Lin et al., 2014), Flickr30K (Young et al., 2014) and NoCaps (Agrawal et al., 2019) datasets. We conduct our evaluation on the validation sets of each dataset. In image captioning experiments involving in-context learning, we utilize the MS COCO training set as the knowledge base from which we retrieve similar examples to construct the context. Each demonstration example comprises an image-text pair, where we randomly sample one of the human-annotated captions per image. We employ the CIDEr-D (Vedantam et al., 2015) and CIDEr-R (dos Santos et al., 2021), which are n-gram-based evaluation metrics, with CIDEr-R being less sensitive to variations in caption length.

For the VQA task, we utilize the VizWiz (Gurari et al., 2018), GQA (Hudson and Manning, 2019), TextVQA (Singh et al., 2019), and OKVQA (Marino et al., 2019) datasets, each designed to evaluate different model capabilities. VizWiz involves real-world images taken by visually impaired users with user-specific questions, while GQA assesses reasoning and compositional skills. TextVQA focuses on optical character recognition; thus, models should recognize text in images to answer the questions. OKVQA is designed to test models’ ability to answer questions about images using external resources or commonsense knowledge. Unlike image captioning, we use each dataset’s training set as the knowledge base. Performance is evaluated using the VQA accuracy metric (Antol et al., 2015), suitable for the open-ended nature of the questions.

3.2 Evaluation Protocol

Demonstrations Retrieval. Inspired by Yang et al. (2023), we retrieve demonstration examples employing a k-Nearest Neighbor approach based on the similarity distance in the visual feature space. We construct a knowledge base $\mathcal{D} = \{(i_1, t_1), \dots, (i_n, t_n)\}$, consisting of images i paired with their corresponding texts t different from those in the evaluation sets. Then, for each query image I , we extract its features $f(I)$ and we retrieve the top- k most similar image-text pairs based on the cosine similarity between visual features. Formally, the retrieved set

openflamingo/OpenFlamingo-9B-vitl-mpt7b
HuggingFaceM4/idefics2-8b-base
HuggingFaceM4/idefics2-8b

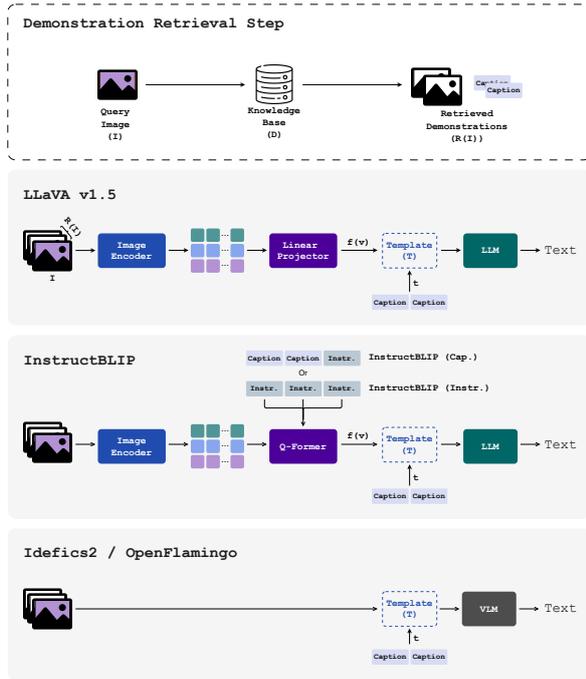


Figure 1: Evaluation pipeline for assessing the in-context learning capability of each analyzed model architecture. We illustrate the modifications made to the original LLaVA v1.5 and InstructBLIP pipelines to support interleaved image-text inputs.

$\mathcal{R}(I)$ of image-text pairs is defined as $\mathcal{R}(I) = \{(i, t) \mid \text{top-}k_{(i,t)} \in \mathcal{D} \text{ sim}(f_I, f_i)\}^2$, where $\text{sim}(\cdot)$ denotes the cosine similarity. We use a ViT (Dosovitskiy et al., 2021)³ to encode the images. To investigate the impact of including multiple demonstration examples, we evaluate prompts containing 0, 1, 3, and 5 demonstrations.

In-Context Learning. We assess the in-context learning capabilities of the InstructBLIP, LLaVA, Idefics2, and OpenFlamingo architectures across various scenarios. Although in-context learning is straightforward for Idefics2 and OpenFlamingo, as they were trained with multiple interleaved image-text instances, implementing a similar pipeline for InstructBLIP and LLaVA poses some challenges. In Figure 1, we illustrate the pipeline adopted for each model architecture.

Since InstructBLIP and LLaVA were trained on image-text pairs, we adapted these models to handle multiple images per sample. Regarding LLaVA, for each sample, comprising multiple images interleaved with texts, we pass the images through the visual encoder and extract the visual features,

²For simplicity, we denote $f(i)$ as f_i and $f(I)$ as f_I .

³<https://huggingface.co/google/vit-large-patch16-224-in21k>



Figure 2: Investigated prompt templates.

which are then projected into the LLM token embedding space using an MLP block. Similarly, token embeddings are extracted for the texts. The projected visual features $f(v)$ and text embeddings t are concatenated into a single sequence and passed as input to the LLM.

In the case of InstructBLIP, we first extract the visual features for all images in the sample. However, unlike LLaVA, InstructBLIP employs an instruction-aware Q-Former to bridge modalities, which takes an image-text pair as input. This way, for the image captioning task, we explore two different approaches: (InstructBLIP Cap.) passing to the Q-Former the image-caption pairs for the demonstration examples, and the query image – for which we aim to generate the caption – alongside an instruction; and, (InstructBLIP Instr.) feeding Q-Former with image-instruction pairs for each image in the sample, including the query image. The output of the Q-Former is a set of query embeddings $f(v)$ that represent the visual information, with dimensions matching those of the LLM’s input token embeddings. These query embeddings are, then, inserted into the template textual embeddings and fed into the LLM. For VQA, each demonstration example consists of an image and a corresponding question-answer pair, which are passed to the Q-Former. For the query image, we provide the image along with the question.

Prompt Construction. To evaluate the models’ ability to adapt at inference time, we construct a

prompt by inserting the visual information $f(v)$ ⁴ into a natural language template \mathcal{T} . We investigate scenarios using prompts with three different levels of detail, as illustrated in Figure 2. First, we use prompts containing only an instruction. Note that we do not necessarily use the same instructions as those reported in the original works. Instead, we choose to evaluate the different models under the same conditions. Next, we test straightforward prompts that include demonstration examples $\mathcal{R}(I)$ – image-caption pairs for image captioning and image-question-answer triplets for VQA – directly into the template \mathcal{T} . These examples are concatenated and followed by an instruction. Finally, building upon the Socratic Models (Zeng et al., 2023), we further explore detailed prompts based on Socratic templates (Zeng et al., 2023; Ramos et al., 2023) that specify the task and detail the format in which the demonstration examples are presented. In this case, the demonstrations are inserted at predefined positions within the template. We also experiment with minor variations of these templates to assess their impact. In all experiments involving demonstration examples, we follow the approach proposed by Baldassini et al. (2024), presenting examples in increasing similarity order relative to the query image as models tend to give more relevance to the last demonstrations. Specifically, we select the top- k examples, sorting them so that the most similar example is presented last.

4 Results and Discussions

Instruction-only Scenario. To establish a baseline and analyze the in-context learning capabilities of VLMs, we first conduct inference using instruction-only prompts without demonstration examples and investigate the VLMs’ sensitivity to minor prompt variations. For this, we evaluate models on the image captioning task using four similar instructions, three sourced from (Dai et al., 2024): “Write a short description for the image.”, “A short image caption.”, and “A short image description:” along with a fourth variant, “A short image description.” where the colon is replaced with a period. Results in Table A1 show that InstructBLIP models (with Vicuna-7B and FlanT5-XL) exhibit consistent performance with minimal fluctuations, unlike

⁴The visual information $f(v)$ can consist solely of the query image, as in the instruction-only scenario, or also include the demonstrations $\mathcal{R}(I)$, which is the case of the in-context learning.

other models. LLaVA demonstrates high sensitivity, with its performance on the MS COCO dataset declining significantly when the period in “A short image description.” is replaced with a colon, while remaining stable on other datasets. This suggests a potential memorization of MS COCO’s content, as this dataset is used to generate instruction-following training data. In contrast, Idefics2 and OpenFlamingo perform best with “A short image description:” and show reduced performance when the colon is replaced with a period. Idefics2 (Base) exhibits greater variation before instruction-tuning, indicating that this phase enhances robustness to prompt variations.

Impact of Prompt on In-Context Learning. To investigate the influence of prompt structure on in-context learning, we evaluate models on the image captioning task using prompts designed with two levels of detail. The first prompt follows a straightforward template, where demonstration image-caption pairs are directly concatenated with an instruction. In contrast, the second prompt is more detailed, explicitly specifying the format in which the demonstration examples are presented and including the phrase “I am an intelligent image captioning bot.” (Section 3.2). For this experiment, we use the MS COCO training set as the knowledge base, and each sample includes three demonstration examples retrieved as context. The results of this evaluation, along with the best performance in the instruction-only scenario, are reported in Figure 3.

One can observe that all models, except instruction-tuned Idefics2, perform better in the instruction-only scenario than when provided with in-context demonstrations. These results indicate that these VLMs struggle to effectively utilize contextual information to adapt their outputs, thus exhibiting weak in-context learning abilities. Particularly, OpenFlamingo performs poorly with straightforward prompts, demonstrating a sharp decline in performance in this scenario. Furthermore, OpenFlamingo and instruction-tuned Idefics2, both of which are trained on interleaved image-text datasets, are the models least affected by the shift from instruction-only to in-context learning scenarios. It is worth noting that Idefics2 (Base) performs better with the straightforward prompt than with the detailed one. However, after instruction-tuning, its performance with the detailed prompt improves significantly, outperforming even the instruction-only

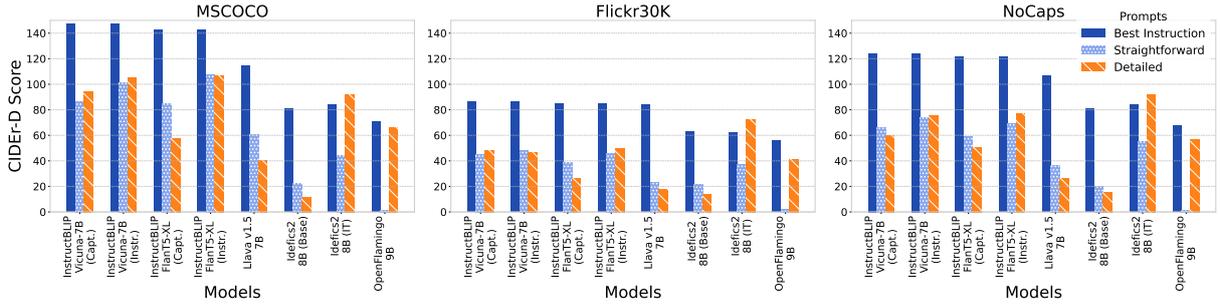


Figure 3: **Comparison of Instruction-only and In-Context Learning Scenarios.** Evaluation results for image captioning task under in-context learning using straightforward and detailed prompts. “Idefics2 8B (IT)” stands for the instruction-tuned checkpoint of the Idefics2 architecture.

432 setup, where its performance remains relatively stable. This result indicates that instruction-tuning
 433 enhances the model’s ability to comprehend the detailed instruction, while training on interleaved
 434 image-text datasets helps the model better leverage contextual information.
 435
 436
 437

438 Regarding InstructBLIP models, performance is further influenced by the type of input provided to
 439 the Q-Former. Specifically, using image-caption pairs from demonstration examples leads to lower
 440 performance than image-instruction pairs. Possibly, this is because InstructBLIP’s Q-Former is pri-
 441 marily exposed to instructions rather than captions during instruction-tuning. Additionally, Instruct-
 442 BLIP FlanT5-XL performs better with straight-
 443 forward prompts, whereas InstructBLIP Vicuna-
 444 7B achieves higher results with detailed prompts.
 445 This discrepancy is likely due to FlanT5-XL’s
 446 fine-tuning on datasets containing few-shot exam-
 447 ples, whose format is similar to the straightfor-
 448 ward prompt.
 449
 450
 451
 452

453 Although there is a notable performance drop
 454 when shifting from instruction-only to in-context
 455 learning setup, InstructBLIP models remain com-
 456 petitive with Idefics2 and OpenFlamingo, despite
 457 not being trained on interleaved image-text datasets.
 458 In contrast, LLaVA struggles significantly in the
 459 in-context learning scenario. We hypothesize that
 460 Q-Former can compress the visual information into
 461 a small set of tokens, allowing InstructBLIP to bet-
 462 ter leverage the LLM’s in-context learning ability.
 463 Conversely, LLaVA maps each visual patch into
 464 one input token embedding using a linear layer, re-
 465 quiring a long sequence of tokens to represent all in-
 466 put images (demonstrations and query), which may
 467 confound its LLM block. This hypothesis aligns
 468 with the findings of [Laurençon et al. \(2024\)](#), which
 469 suggest that reducing the number of visual tokens

470 can improve performance on downstream tasks.

471 Overall, these results indicate that the evaluated
 472 VLMs struggle to leverage the contextual infor-
 473 mation and underscore the impact of prompt de-
 474 sign on in-context learning performance. Detailed
 475 prompts that specify both the task and the structure
 476 of demonstration examples proved to be more ef-
 477 fective than simply concatenating demonstrations.
 478 Also, our findings indicate that both instruction-
 479 tuning and training on interleaved image-text
 480 datasets enhance in-context learning ability. No-
 481 tably, models with advanced modality projectors
 482 can achieve competitive performance even when
 483 trained on datasets containing only single image-
 484 text pairs per sample, offering a more cost-efficient
 485 alternative to training on interleaved datasets.

486 **Influence of the Number of Demonstrations.** In
 487 our previous experiments (Section 4), we fixed the
 488 number of demonstrations at three per sample. We
 489 observed that the detailed prompt generally im-
 490 proves performance. Building on this finding, we
 491 now investigate whether increasing the number of
 492 demonstrations (shots) in context further enhances
 493 model performance. To test this hypothesis, we
 494 evaluate the models on image captioning, using the
 495 previously defined detailed prompt, and on four
 496 VQA datasets. In this experiment, we vary the
 497 number of shots among 0, 1, 3, and 5. In the 0-
 498 shot setting, the prompt consists only of the tem-
 499 plate, without any demonstrations. We emphasize
 500 that this differs from the instruction-only scenario,
 501 as the 0-shot prompt signals a demonstration will
 502 be provided, but no actual demonstration is given.
 503 This setup allows us to evaluate the performance
 504 gains achieved by incorporating more demon-
 505 strations. The image captioning and VQA results are
 506 summarized in Figure 4 (the numeric results for
 507 image captioning and VQA can also be found in

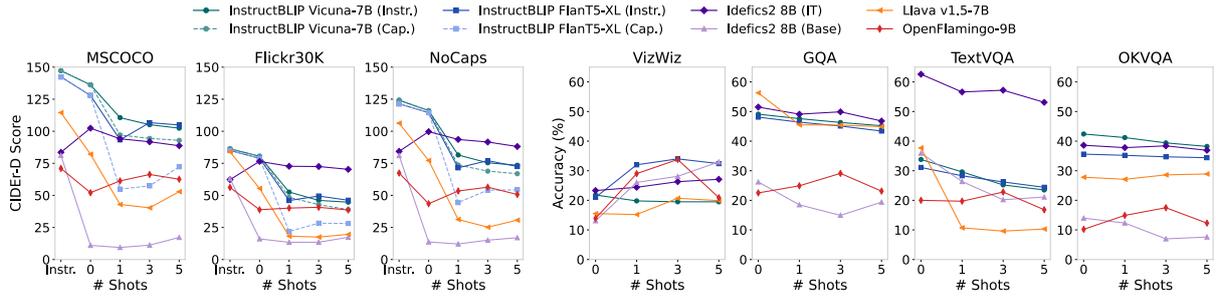


Figure 4: **Influence of the number of demonstration examples on performance.** We evaluate the impact of varying the number of demonstration examples (shots) in the context. For image captioning, we use a detailed prompt and employ the MS COCO training set as the knowledge base, plotting the CIDEr-D score. “Instr.” in the x-axis of charts with image captioning results stands for the best results in the instruction-only scenario. For VQA, we utilize the corresponding training set for each dataset as the knowledge base and report the VQA accuracy.

Tables A3 and A4, respectively).

For image captioning, our results reveal that most models perform better in the instruction-only and 0-shot scenarios than when demonstrations are provided. Furthermore, we do not observe consistent improvements as demonstrations increase. In fact, incorporating more demonstrations often degrades performance relative to the 0-shot setup. However, consistent with prior observations, the Idefics2 and OpenFlamingo models appear to be the least affected by the demonstrations in the in-context learning setting. Specifically, Idefics2 (Base) and OpenFlamingo show slight performance gains as the number of shots increases, while the instruction-tuned Idefics2 model maintains a relatively stable performance. Note that InstructBLIP models achieve the highest performance on MS COCO but experience significant drops on Flickr30K and NoCaps, where the instruction-tuned Idefics2 model outperforms them. LLaVA is the most hampered by the demonstrations, it faces a notable decline on Flickr30K and NoCaps when demonstrations are included. This result corroborates our hypothesis that the long sequence of tokens required to represent the input images may confound the LLM.

Similar to image captioning, in VQA, we observe that models generally perform better across most datasets without in-context demonstrations. However, an opposite trend is observed for Vizwiz, where the inclusion of demonstrations appears beneficial. A detailed analysis (Section A.5.1) reveals that this effect is due to a strong dataset imbalance: the answer “unanswerable” appears more than a thousand times, while most other answers occur only once. Likewise, many of the provided demon-

strations are also annotated as “unanswerable” leading models to favor this response. Additionally, in the TextVQA dataset, models’ performance declines consistently as more demonstrations are introduced. This drop aligns with expectations, as answering questions in TextVQA requires recognizing text within images, and, in this case, similar examples in the context may confound the models.

Furthermore, for the GQA and OKVQA datasets, the performance of InstructBLIP models, LLaVA, and instruction-tuned Idefics2 remains relatively unchanged as the number of shots increases. This suggests that these models overlook in-context demonstrations for reasoning-based tasks. Nevertheless, it is interesting to note that they significantly outperform Idefics2 (Base) and OpenFlamingo on these datasets, underscoring the importance of instruction-tuning for VQA tasks requiring reasoning.

Our results suggest that increasing the number of demonstrations in the context does not necessarily enhance model performance. Instead, refining model architectures or training strategies may be necessary to leverage contextual information better. Particularly, instruction-tuned models achieve better results on reasoning-intensive VQA tasks, while models trained on interleaved image-text datasets exhibit better in-context learning ability. Due to computational constraints, our evaluation is limited to up to 5 demonstrations. However, our results show fluctuations in scores across 1, 3, and 5 shots. Therefore, further large-scale exploration is needed to fully understand the impact of number of demonstrations on performance.

Similar vs. Random Demonstrations. To investigate the impact of similar demonstrations on final

580 results, we conduct a comparative analysis under
 581 two scenarios: one where demonstrations are similar
 582 to the query image and another with demon-
 583 strations from the same task but randomly chosen,
 584 either related or unrelated to the query image,
 585 referred to as random demonstrations. We hypothe-
 586 size is that providing examples with content similar
 587 to the query image leads to better performance
 588 than using random demonstrations. To validate
 589 this, we fix the number of demonstrations at three
 590 and conduct experiments using both similar (as
 591 described in Section 3.2) and random demon-
 592 strations for image captioning and VQA tasks. For
 593 image captioning, we employ a detailed prompt
 594 to maintain consistency with previous experiments.
 595 Figure 5 illustrates the difference in scores between
 596 similar and random demonstrations across image
 597 captioning and VQA datasets.

598 Our experimental results highlight distinct be-
 599 haviors across models in both image captioning and
 600 VQA tasks when exposed to similar and random
 601 demonstrations. In image captioning, InstructBLIP
 602 Vicuna-7B and LLaVA 7B demonstrate the most
 603 substantial performance gains with similar demon-
 604 strations, particularly on MS COCO and NoCaps
 605 datasets. In contrast, OpenFlamingo 9B exhibits a
 606 sharp performance drop, indicating that this model
 607 struggles to effectively leverage visual elements
 608 similar to the query image.

609 In VQA, most models benefit more from simi-
 610 lar demonstrations than from random ones, with
 611 notable improvements on the OKVQA dataset.
 612 OKVQA consists of images and general questions
 613 that require commonsense knowledge. Then, simi-
 614 lar demonstrations help models generate more ac-
 615 curate responses, whereas random demonstrations
 616 can confound them. In contrast, in TextVQA, mod-
 617 els exhibit the greatest performance drop when
 618 using similar demonstrations. That is, models
 619 perform better with random demonstrations than
 620 with similar ones. We hypothesize that, as answer-
 621 ing TextVQA questions requires recognizing text
 622 within images, showing random task-related exam-
 623 ples might help models focus on the task itself. On
 624 the other hand, similar demonstrations could intro-
 625 duce visual distractions and lead to answer copying
 626 from provided examples.

627 5 Conclusion

628 In this paper, we systematically analyze in-context
 629 learning in VLMs, evaluating six models from four

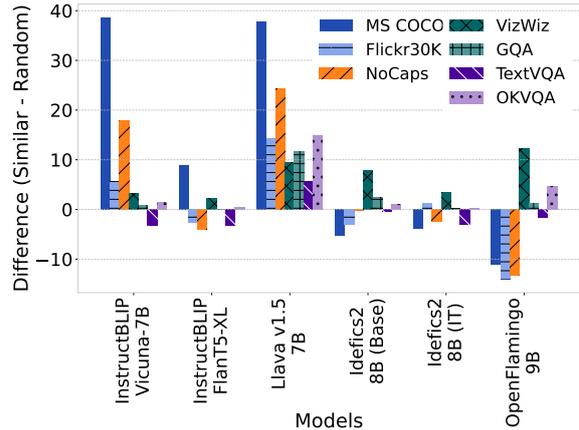


Figure 5: Difference in scores between in-context learning using similar demonstrations and random ones across image captioning and VQA datasets. For the image captioning datasets, we consider the detailed prompt. We plot the difference in CIDEr-D score for image captioning and VQA accuracy for VQA datasets.

630 distinct architectures across multiple image cap-
 631 tioning and VQA benchmarks. We investigate the
 632 impact of prompt construction, demonstration se-
 633 lection, and model design on in-context learning.
 634 Unlike previous work, we analyze models beyond
 635 those trained on interleaved image-text datasets.
 636 Our findings reveal that the evaluated models strug-
 637 gle to utilize contextual information to refine their
 638 outputs. However, detailed prompts, explicitly
 639 defining both the task and the structure of demon-
 640 stration examples, significantly enhance this ability
 641 compared to simply concatenating examples. In-
 642 creasing the number of demonstrations does not
 643 necessarily yield better results. While instruction-
 644 tuning helps models comprehend detailed instruc-
 645 tions, it may reduce their in-context learning ca-
 646 pacity. In contrast, training on interleaved image-text
 647 datasets enhances such ability. Additionally, we
 648 show that models with advanced modality projec-
 649 tors can achieve competitive in-context learning
 650 performance even when trained on single image-
 651 text pairs, offering a cost-efficient alternative.

652 This work sheds light on key limitations in cur-
 653 rent VLMs regarding their in-context learning abil-
 654 ity. Future research should explore modality pro-
 655 jectors to better integrate LLMs’ in-context learn-
 656 ing abilities into VLMs, as well as a combined
 657 approach using instruction-tuning and interleaved
 658 image-text training. Another promising direction is
 659 the inclusion of both positive and negative demon-
 660 strations, which could help models better distin-
 661 guish between correct and incorrect responses.

662 Limitations

663 Although our analysis focuses on VLMs with up to
664 9B parameters and a maximum of 5 demonstrations
665 per query due to computational constraints, study-
666 ing larger models and increasing the number of
667 shots would be important to determine whether our
668 conclusions hold at a greater scale. Furthermore,
669 to better understand the role of instruction-tuning
670 and training of interleaved image-text datasets, it
671 would be interesting to extend our analysis to a
672 broader range of model architectures evaluating
673 models before and after instruction-tuning. Finally,
674 our analysis is limited to VLMs trained in English-
675 language texts. However, evaluating the in-context
676 learning capacity of multilingual models is essen-
677 tial. It would be necessary to study whether in-
678 context learning can improve VLMs performance
679 on low-resource languages.

680 Ethics Statement

681 This study systematically analyzes the in-context
682 learning capabilities of publicly available VLMs.
683 Our analysis is based solely on publicly available
684 image captioning and VQA datasets, and we fully
685 comply with the terms of use and licensing agree-
686 ments associated with each model and dataset. We
687 do not conduct any fine-tuning or modifications in
688 the models that could introduce unintended risks.
689 However, we recognize that our work reflects the
690 existing limitations and potential risks of the eval-
691 uated models, including but not limited to gender,
692 racial, and cultural biases, as well as the potential
693 for generating misinformation or disinformation.

694 References

695 Harsh Agrawal, Karan Desai, Yufei Wang, Xinlei Chen,
696 Rishabh Jain, Mark Johnson, Dhruv Batra, Devi
697 Parikh, Stefan Lee, and Peter Anderson. 2019. No-
698 Caps: Novel object captioning at scale. In *IEEE/CVF*
699 *International Conference on Computer Vision*, pages
700 8948–8957.

701 Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc,
702 Antoine Miech, Iain Barr, Yana Hasson, Karel
703 Lenc, Arthur Mensch, Katherine Millican, Malcolm
704 Reynolds, Roman Ring, Eliza Rutherford, Serkan
705 Cabi, Tengda Han, Zhitao Gong, Sina Samangooei,
706 Marianne Monteiro, Jacob L Menick, Sebastian
707 Borgeaud, Andy Brock, Aida Nematzadeh, Sahand
708 Sharifzadeh, Miłojaj Bińkowski, Ricardo Barreira,
709 Oriol Vinyals, Andrew Zisserman, and Karén Si-
710 monyan. 2022. *Flamingo: a visual language model*
711 *for few-shot learning*. In *Advances in Neural Infor-*

mation Processing Systems, volume 35, pages 23716–
23736. Curran Associates, Inc. 712 713

Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Mar-
garet Mitchell, Dhruv Batra, C Lawrence Zitnick, and
Devi Parikh. 2015. VQA: Visual Question Answer-
ing. In *IEEE International Conference on Computer*
Vision, pages 2425–2433. 714 715 716 717 718

Anas Awadalla, Irena Gao, Josh Gardner, Jack Hes-
sel, Yusuf Hanafy, Wanrong Zhu, Kalyani Marathe,
Yonatan Bitton, Samir Gadre, Shiori Sagawa, et al.
2023. OpenFlamingo: An open-source framework
for training large autoregressive vision-language
models. *arXiv preprint arXiv:2308.01390*. 719 720 721 722 723 724

Folco Bertini Baldassini, Mustafa Shukor, Matthieu
Cord, Laure Soulier, and Benjamin Piwowarski. 2024.
What makes multimodal in-context learning work?
In *IEEE/CVF Conference on Computer Vision and*
Pattern Recognition, pages 1539–1550. 725 726 727 728 729

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. 730 731 732 733 734 735

Junbum Cha, Wooyoung Kang, Jonghwan Mun, and
Byungseok Roh. 2024. Honeybee: Locality-
enhanced projector for multimodal llm. In
IEEE/CVF Conference on Computer Vision and Pat-
tern Recognition (CVPR), pages 13817–13827. 736 737 738 739 740

Shuo Chen, Zhen Han, Bailan He, Mark Buckley, Philip
Torr, Volker Tresp, and Jindong Gu. 2024. *Under-*
standing and improving in-context learning on vision-
language models. In *Workshop on Mathematical and*
Empirical Understanding of Foundation Models. 741 742 743 744 745

Wenliang Dai, Junnan Li, Dongxu Li, Anthony
Meng Huat Tiong, Junqi Zhao, Weisheng Wang,
Boyang Li, Pascale N Fung, and Steven Hoi.
2024. InstructBLIP: Towards general-purpose vision-
language models with instruction tuning. *Advances*
in Neural Information Processing Systems, 36. 746 747 748 749 750 751

Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Jingyuan
Ma, Rui Li, Heming Xia, Jingjing Xu, Zhiyong Wu,
Baobao Chang, et al. 2024. A survey on in-context
learning. In *Conference on Empirical Methods in*
Natural Language Processing, pages 1107–1128. 752 753 754 755 756

Gabriel Oliveira dos Santos, Esther Luna Colom-
bini, and Sandra Avila. 2021. CIDEr-R: Robust
Consensus-based Image Description Evaluation. In
Seventh Workshop on Noisy User-generated Text (W-
NUT 2021), pages 351–360. 757 758 759 760 761

Alexey Dosovitskiy, Lucas Beyer, Alexander
Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,
Thomas Unterthiner, Mostafa Dehghani, Matthias
Minderer, Georg Heigold, Sylvain Gelly, Jakob
Uszkoreit, and Neil Houlsby. 2021. *An image*
is worth 16x16 words: Transformers for image 762 763 764 765 766 767

768	recognition at scale . In <i>International Conference on Learning Representations</i> .	825
769		826
770	Gregor Geigle, Abhay Jain, Radu Timofte, and Goran Glavaš. 2024. mBLIP: Efficient Bootstrapping of Multilingual Vision-LLMs . In <i>3rd Workshop on Advances in Language and Vision Research (ALVR)</i> , pages 7–25, Bangkok, Thailand. Association for Computational Linguistics.	827
771		828
772		829
773		830
774		831
775		
776	Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. 2017. Making the v in VQA matter: Elevating the role of image understanding in visual question answering. In <i>IEEE Conference on Computer Vision and Pattern Recognition</i> , pages 6904–6913.	832
777		833
778		834
779		835
780		836
781		
782	Danna Gurari, Qing Li, Abigale J Stangl, Anhong Guo, Chi Lin, Kristen Grauman, Jiebo Luo, and Jeffrey P Bigham. 2018. VizWiz grand challenge: Answering visual questions from blind people. In <i>IEEE Conference on Computer Vision and Pattern Recognition</i> , pages 3608–3617.	837
783		838
784		839
785		840
786		841
787		
788	Anwen Hu, Haiyang Xu, Jiabo Ye, Ming Yan, Liang Zhang, Bo Zhang, Ji Zhang, Qin Jin, Fei Huang, and Jingren Zhou. 2024. mPLUG-DocOwl 1.5: Unified structure learning for OCR-free document understanding . In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 3096–3120, Miami, Florida, USA. Association for Computational Linguistics.	842
789		843
790		844
791		845
792		846
793		
794		
795		
796	Drew A Hudson and Christopher D Manning. 2019. GQA: A new dataset for real-world visual reasoning and compositional question answering. In <i>IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 6700–6709.	847
797		848
798		849
799		850
800		
801	Jing Yu Koh, Ruslan Salakhutdinov, and Daniel Fried. 2023. Grounding language models to images for multimodal inputs and outputs . In <i>40th International Conference on Machine Learning</i> , volume 202 of <i>Proceedings of Machine Learning Research</i> , pages 17283–17300. PMLR.	851
802		852
803		853
804		
805		
806		
807	Hugo Laurençon, Lucile Saulnier, Leo Tronchon, Stas Bekman, Amanpreet Singh, Anton Lozhkov, Thomas Wang, Siddharth Karamcheti, Alexander Rush, Douwe Kiela, Matthieu Cord, and Victor Sanh. 2023. OBELICS: An Open Web-Scale Filtered Dataset of Interleaved Image-Text Documents . In <i>Advances in Neural Information Processing Systems</i> , volume 36, pages 71683–71702. Curran Associates, Inc.	854
808		855
809		856
810		857
811		858
812		859
813		860
814		861
815		862
816	Hugo Laurençon, Leo Tronchon, Matthieu Cord, and Victor Sanh. 2024. What matters when building vision-language models? In <i>The Thirty-eighth Annual Conference on Neural Information Processing Systems</i> .	863
817		864
818		865
819		866
820		867
821	Hugo Laurençon, Andrés Marafioti, Victor Sanh, and Léo Tronchon. 2024. Building and better understanding vision-language models: insights and future directions. <i>arXiv preprint arXiv:2408.12637</i> .	868
822		869
823		870
824		
	Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023. BLIP-2: Bootstrapping language-image pre-training with frozen image encoders and large language models . In <i>40th International Conference on Machine Learning</i> , volume 202 of <i>Proceedings of Machine Learning Research</i> , pages 19730–19742. PMLR.	871
		872
		873
		874
		875
	Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. 2022. BLIP: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In <i>International Conference on Machine Learning</i> , pages 12888–12900. PMLR.	876
		877
		878
		879
		880
	Ji Lin, Hongxu Yin, Wei Ping, Pavlo Molchanov, Mohammad Shoeybi, and Song Han. 2024. VILA: On Pre-training for Visual Language Models . In <i>IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 26689–26699.	871
		872
		873
		874
		875
	Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 2014. Microsoft coco: Common objects in context. In <i>European Conference on Computer Vision</i> , pages 740–755.	876
		877
		878
		879
		880
	Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2023. Improved baselines with visual instruction tuning. In <i>NeurIPS 2023 Workshop on Instruction Tuning and Instruction Following</i> .	871
		872
		873
		874
		875
	Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2024. Visual instruction tuning. <i>Advances in Neural Information Processing Systems</i> , 36.	871
		872
		873
		874
		875
	Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. 2019. OK-VQA: A visual question answering benchmark requiring external knowledge. In <i>IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 3195–3204.	876
		877
		878
		879
		880
	Brandon McKinzie, Zhe Gan, Jean-Philippe Fauconnier, Sam Dodge, Bowen Zhang, Philipp Dufter, Dhruvi Shah, Xianzhi Du, Futang Peng, Anton Belyi, Haotian Zhang, Karanjeet Singh, Doug Kang, Hongyu Hè, Max Schwarzer, Tom Gunter, Xiang Kong, Aonan Zhang, Jianyu Wang, Chong Wang, Nan Du, Tao Lei, Sam Wiseman, Mark Lee, Zirui Wang, Ruoming Pang, Peter Grasch, Alexander Toshev, and Yinfei Yang. 2025. MM1: Methods, Analysis and Insights from Multimodal LLM Pre-training. In <i>Computer Vision – ECCV 2024</i> , pages 304–323, Cham. Springer Nature Switzerland.	876
		877
		878
		879
		880
	Libo Qin, Qiguang Chen, Hao Fei, Zhi Chen, Min Li, and Wanxiang Che. 2024. What factors affect multimodal in-context learning? an in-depth exploration . In <i>The Thirty-Eighth Annual Conference on Neural Information Processing Systems</i> .	876
		877
		878
		879
		880
	Rita Ramos, Bruno Martins, and Desmond Elliott. 2023. Lmcap: Few-shot multilingual image captioning by retrieval augmented language model prompting. In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 1635–1651.	876
		877
		878
		879
		880

881	Dustin Schwenk, Apoorv Khandelwal, Christopher Clark, Kenneth Marino, and Roozbeh Mottaghi. 2022. A-OKVQA: A benchmark for visual question answering using world knowledge. In <i>European Conference on Computer Vision</i> , pages 146–162. Springer.	Xu Yang, Yongliang Wu, Mingzhuo Yang, Haokun Chen, and Xin Geng. 2024. Exploring diverse in-context configurations for image captioning. <i>Advances in Neural Information Processing Systems</i> , 36.	936 937 938 939 940
886	Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. 2018. Conceptual captions: A cleaned, hypernamed, image alt-text dataset for automatic image captioning . In <i>56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 2556–2565, Melbourne, Australia. Association for Computational Linguistics.	Zhuolin Yang, Wei Ping, Zihan Liu, Vijay Korthikanti, Weili Nie, De-An Huang, Linxi Fan, Zhiding Yu, Shiyi Lan, Bo Li, et al. 2023. Re-ViLM: Retrieval-Augmented Visual Language Model for Zero and Few-Shot Image Captioning. In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 11844–11857.	941 942 943 944 945 946 947
893	Mustafa Shukor, Corentin Dancette, and Matthieu Cord. 2023. eP-ALM: Efficient Perceptual Augmentation of Language Models. In <i>IEEE/CVF International Conference on Computer Vision (ICCV)</i> , pages 22056–22069.	Shukang Yin, Chaoyou Fu, Sirui Zhao, Ke Li, Xing Sun, Tong Xu, and Enhong Chen. 2024. A survey on multimodal large language models . <i>National Science Review</i> , pages 1–18.	948 949 950 951
898	Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh, and Marcus Rohrbach. 2019. Towards VQA models that can read. In <i>IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 8317–8326.	Peter Young, Alice Lai, Micah Hodosh, and Julia Hockenmaier. 2014. From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions. <i>Transactions of the Association for Computational Linguistics</i> , 2:67–78.	952 953 954 955 956
903	Yixuan Su, Tian Lan, Huayang Li, Jialu Xu, Yan Wang, and Deng Cai. 2023. PandaGPT: One model to instruction-follow them all . In <i>1st Workshop on Taming Large Language Models: Controllability in the Era of Interactive Assistants!</i> , pages 11–23, Prague, Czech Republic. Association for Computational Linguistics.	Andy Zeng, Maria Attarian, brian ichter, Krzysztof Marcin Choromanski, Adrian Wong, Stefan Welker, Federico Tombari, Aveek Purohit, Michael S Ryoo, Vikas Sindhwani, Johnny Lee, Vincent Vanhoucke, and Pete Florence. 2023. Socratic models: Composing zero-shot multimodal reasoning with language . In <i>The Eleventh International Conference on Learning Representations</i> .	957 958 959 960 961 962 963 964
910	Quan Sun, Yufeng Cui, Xiaosong Zhang, Fan Zhang, Qiyang Yu, Yueze Wang, Yongming Rao, Jingjing Liu, Tiejun Huang, and Xinlong Wang. 2024. Generative multimodal models are in-context learners. In <i>IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 14398–14409.	Duzhen Zhang, Yahan Yu, Jiahua Dong, Chenxing Li, Dan Su, Chenhui Chu, and Dong Yu. 2024. MM-LLMs: Recent advances in MultiModal large language models . In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 12401–12430, Bangkok, Thailand. Association for Computational Linguistics.	965 966 967 968 969 970 971
916	Ramakrishna Vedantam, C Lawrence Zitnick, and Devi Parikh. 2015. CIDEr: Consensus-based Image Description Evaluation. In <i>IEEE Conference on Computer Vision and Pattern Recognition</i> , pages 4566–4575.	Haozhe Zhao, Zefan Cai, Shuzheng Si, Xiaojian Ma, Kaikai An, Liang Chen, Zixuan Liu, Sheng Wang, Wenjuan Han, and Baobao Chang. 2024. MMICL: Empowering Vision-language Model with Multi-Modal In-Context Learning . In <i>The Twelfth International Conference on Learning Representations</i> .	972 973 974 975 976 977
921	Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. 2022. Emergent abilities of large language models. <i>Transactions on Machine Learning Research</i> .	Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. 2024a. MiniGPT-4: Enhancing Vision-Language Understanding with Advanced Large Language Models . In <i>The Twelfth International Conference on Learning Representations</i> .	978 979 980 981 982
926	Le Xue, Manli Shu, Anas Awadalla, Jun Wang, An Yan, Senthil Purushwalkam, Honglu Zhou, Viraj Prabhu, Yutong Dai, Michael S. Ryoo, Shrikant Kendre, Jieyu Zhang, Can Qin, Shu Zhang, Chia-Chih Chen, Ning Yu, Juntao Tan, Tulika Manoj Awalganekar, Shelby Heinecke, Huan Wang, Yejin Choi, Ludwig Schmidt, Zeyuan Chen, Silvio Savarese, Juan Carlos Niebles, Caiming Xiong, and Ran Xu. 2024. xGen-MM (BLIP-3): A Family of Open Large Multimodal Models . <i>CoRR</i> , abs/2408.08872.	Wanrong Zhu, Jack Hessel, Anas Awadalla, Samir Yitzhak Gadre, Jesse Dodge, Alex Fang, Youngjae Yu, Ludwig Schmidt, William Yang Wang, and Yejin Choi. 2024b. Multimodal c4: An open, billion-scale corpus of images interleaved with text. <i>Advances in Neural Information Processing Systems</i> , 36.	983 984 985 986 987 988 989

A Appendix

A.1 Results on Instruction-only Scenario

As detailed in Section 4, we first conduct inference using instruction-only prompts, i.e., without including any demonstration examples, to establish a baseline for our in-context learning experiments. To do so, we test four similar instructions, three of which are selected from (Dai et al., 2024): (I1) “Write a short description for the image.”, (I2) “A short image caption.”, and (I3) “A short image description:”. We also create a fourth instruction, (I4) “A short image description.”, by replacing the colon in the latter instruction with a period. Table A1 summarizes the results of these experiments.

InstructBLIP models (with Vicuna-7B and FlanT5-XL) exhibit consistent performance, with only minor fluctuations across the different instructions. Interestingly, this consistency does not extend to the other models. LLaVA shows one of the greatest sensitivity to instruction variations, performing best with the instruction “A short image caption.” and worst with “Write a short description for the image.”. Notably, its performance on the MS COCO dataset declines significantly when the period in “A short image description.” is replaced with a colon, while remaining stable on the other datasets. This drop in results on MS COCO suggests that LLaVA may be memorizing the content of MS COCO, as this dataset is used to generate instruction-following training data. In contrast, Idefics2 models and OpenFlamingo perform best with the instruction “A short image description:” and show reduced performance when the colon is replaced with a period. Also, the difference between the highest and lowest scores is more pronounced in Idefics2 before the instruction-tuning phase (Idefics2 (Base)), possibly because this phase enhances the model’s robustness to minor prompt variations. A similar trend is observed in OpenFlamingo, which also does not undergo an instruction-tuning phase during training.

A.2 Experimental Results in Numbers

We provide the numerical results of the experiments regarding the impact of prompt in the in-context learning ability (Section 4) and the influence of the number of demonstrations in the context on the performance (Section 4). The results are divided into three tables. Table A2 presents

the results for the image captioning task under instruction-only and in-context learning scenarios; it shows the best performance in the instruction-only scenario alongside the results of in-context learning with straightforward and detailed prompts. Tables A3 and A4 show the results for image captioning and VQA, respectively, varying the number of demonstration examples in the context.

A.3 Ablation on Detailed Prompt

Building on the findings from the instruction-only scenario (Section 4), we investigate the impact of small modifications to the detailed prompt used to evaluate the in-context learning capabilities of models in the image captioning task. We use (I) *Base detailed prompt* for our experiments and test various small changes to this template (Figure A1).

The results, summarized in Table A5, reveal interesting insights. First, removing the initial phrase (prompt 2) significantly hampers the performance of most models. Second, models generally perform better when the word “creative” is removed from the prompt (prompt 3). However, removing both the initial phrase and the word “creative” (prompt 4) produces intermediate results, suggesting that the effects of these changes are combined. The best prompt in most cases is to keep the initial phrase while removing the word “creative” (prompt 3), which leads to the highest performance.

These changes in the prompt can result in substantial performance differences, with variations of up to 20 points in CIDEr scores. Among the experimented models, InstructBLIP FlanT5-XL demonstrates major sensitivity to prompt modifications. Notably, it fails to generate captions when the word “creative” is included in the prompt, underscoring its dependence on precise prompt phrasing. Finally, as expected, altering the name of the modality projector (prompts 5, 6, and 7, Figure A1) has no impact on model performance, indicating that the models simply ignore this detail.

A.4 Impact of Training Size on In-Context Learning

To further explore the impacts of design decisions on in-context learning, we investigate the impact of the training set size on the model performance in both instruction-only and in-context learning scenarios. Figure A2 illustrates model performance across image captioning datasets as a function of training set size.

Table A1: **Instruction-only scenario.** We evaluate the VLMs on image captioning datasets with different instructions and report the CIDEr-D (\uparrow) and CIDEr-R (\uparrow) scores. The numbers in bold are at least 1 point better than the others. The evaluated instructions are: (I1) “Write a short description for the image.”, (I2) “A short image description.”, (I3) “A short image description:” and (I4) “A short image caption.”.

Model	Instruction	MS COCO		Flickr30K		NoCaps	
		CIDEr-D	CIDEr-R	CIDEr-D	CIDEr-R	CIDEr-D	CIDEr-R
InstructBLIP Vicuna-7B	I1	147.4	149.5	85.1	97.0	123.7	130.2
	I2	146.7	148.4	85.9	97.8	124.0	130.0
	I3	146.8	149.0	86.3	98.4	123.7	130.5
	I4	147.2	149.0	86.3	98.2	124.2	130.5
InstructBLIP FlanT5-XL	I1	142.5	144.5	85.1	96.9	121.5	128.2
	I2	142.4	144.3	85.4	97.2	121.6	128.2
	I3	142.4	144.4	85.1	96.8	121.4	128.1
	I4	142.4	144.4	85.0	97.0	121.4	127.9
LLaVA v1.5-7B	I1	64.9	88.9	47.3	71.8	72.2	93.0
	I2	101.3	113.9	69.6	88.4	99.0	113.1
	I3	78.3	90.6	69.7	87.4	96.5	111.7
	I4	114.5	122.7	83.9	99.2	106.3	117.5
Idefics2-8B (Base)	I1	0.1	3.1	0.0	1.0	0.3	4.6
	I2	9.9	64.6	9.7	61.3	19.0	72.0
	I3	81.2	94.6	63.0	79.7	81.0	95.3
	I4	0.7	1.5	0.7	2.0	0.4	0.9
Idefics2-8B (Instruction-Tuned)	I1	57.5	70.1	51.8	66.5	69.1	80.9
	I2	49.1	59.6	47.5	61.8	67.7	79.5
	I3	83.6	90.1	62.3	74.7	84.3	93.0
	I4	53.5	65.3	41.9	55.9	63.2	75.8
OpenFlamingo-9B	I1	36.1	50.9	31.4	43.4	29.8	49.8
	I2	60.9	72.8	49.4	62.8	63.1	75.3
	I3	71.0	82.0	56.2	69.8	67.4	81.5
	I4	58.7	66.7	47.2	58.8	53.0	63.7

Table A2: **Comparison between instruction-only and in-context learning scenarios.** Evaluation results for image captioning task under in-context learning using straightforward and detailed prompts. “Instruction” refers to the best performance in the instruction-only scenario. Bold numbers highlight the best performance for each model.

Model	Prompt	MS COCO		Flickr30K		NoCaps	
		CIDEr-D (\uparrow)	CIDEr-R (\uparrow)	CIDEr-D (\uparrow)	CIDEr-R (\uparrow)	CIDEr-D (\uparrow)	CIDEr-R (\uparrow)
InstructBLIP Vicuna-7B (Q-Former fed with Caption)	Instruction	147.2	149.0	86.3	98.2	124.2	130.5
	Straightforward	86.2	96.0	45.0	52.1	66.1	73.3
	Detailed	94.4	97.9	47.8	55.8	60.0	71.3
InstructBLIP Vicuna-7B (Q-Former fed with Instruction)	Instruction	147.2	149.0	86.3	98.2	124.2	130.5
	Straightforward	100.9	106.3	47.7	53.6	74.2	78.5
	Detailed	105.1	107.5	46.1	52.3	75.4	79.0
InstructBLIP FlanT5-XL (Q-Former fed with Caption)	Instruction	142.4	144.4	85.0	97.0	121.4	127.9
	Straightforward	84.9	87.1	39.0	44.1	58.6	61.6
	Detailed	57.4	59.1	26.5	30.5	50.0	52.9
InstructBLIP FlanT5-XL (Q-Former fed with Instruction)	Instruction	142.4	144.4	85.0	97.0	121.4	127.9
	Straightforward	107.2	108.2	45.4	51.6	69.2	72.4
	Detailed	106.7	108.6	49.5	56.7	77.1	81.2
LLaVA v1.5-7B	Instruction	114.5	122.7	83.9	99.2	106.3	117.5
	Straightforward	60.3	65.1	23.4	28.0	36.1	40.4
	Detailed	40.6	44.1	17.6	21.1	26.2	29.4
Idefics2-8B (Base)	Instruction	81.2	94.6	63.0	79.7	81.0	95.3
	Straightforward	21.9	35.6	21.2	32.6	19.8	31.4
	Detailed	11.2	33.2	13.5	29.4	15.1	37.6
Idefics2-8B (Instruction-Tuned)	Instruction	83.6	90.1	62.3	74.7	84.3	93.0
	Straightforward	44.1	57.1	37.0	51.1	55.2	68.0
	Detailed	91.8	99.6	72.5	85.5	91.6	101.0
OpenFlamingo-9B	Instruction	71.0	82.0	56.2	69.8	67.4	81.5
	Straightforward	1.2	18.4	1.8	10.6	1.3	13.9
	Detailed	66.3	70.6	40.7	50.3	56.4	65.2

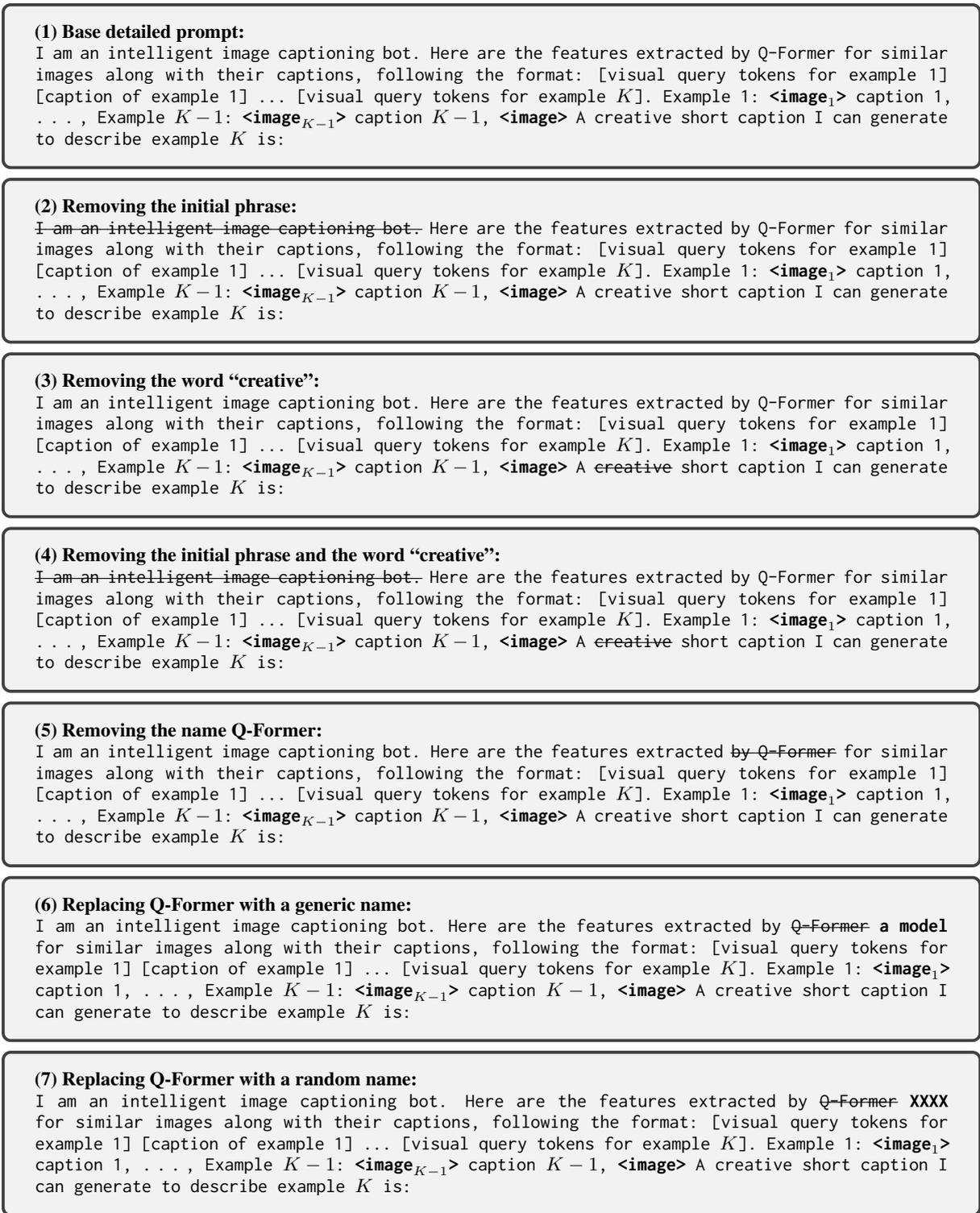


Figure A1: Detailed prompts to evaluate the in-context learning capabilities of models in the image captioning task.

Table A3: **Influence of the number of demonstration examples on image captioning performance.** Evaluation results for image captioning task with the detailed prompt varying the number of demonstration examples (Shot) in the context. Bold numbers highlight the best performance for each model. We use the MS COCO training set as the knowledge base.

Model	Shot	MS COCO		Flickr30K		NoCaps	
		CIDEr-D (↑)	CIDEr-R (↑)	CIDEr-D (↑)	CIDEr-R (↑)	CIDEr-D (↑)	CIDEr-R (↑)
InstructBLIP Vicuna-7B (Q-Former fed with Caption)	Instruction	147.2	149.0	86.3	98.2	124.2	130.5
	0	136.2	139.1	80.5	92.7	116.1	123.2
	1	97.1	100.1	48.5	55.7	73.7	78.3
	3	94.4	97.9	42.8	50.3	68.9	74.3
	5	92.9	96.6	38.9	46.3	66.9	72.7
InstructBLIP Vicuna-7B (Q-Former fed with Instruction)	Instruction	147.2	149.0	86.3	98.2	124.2	130.5
	0	136.2	139.1	80.5	92.7	116.1	123.2
	1	110.6	112.7	52.6	59.8	81.5	85.4
	3	105.1	107.5	46.1	52.3	75.4	79.0
	5	102.4	105.2	44.7	50.3	73.5	77.1
InstructBLIP FlanT5-XL (Q-Former fed with Caption)	Instruction	142.4	144.4	85.0	97.0	121.4	127.9
	0	127.9	131.2	78.9	90.9	114.5	122.0
	1	54.9	56.3	21.9	25.1	44.4	47.1
	3	57.4	59.1	28.3	32.5	54.1	57.5
	5	72.3	74.8	28.1	32.3	54.3	57.5
InstructBLIP FlanT5-XL (Q-Former fed with Instruction)	Instruction	142.4	144.4	85.0	97.0	121.4	127.9
	0	127.9	131.2	78.9	90.9	114.5	122.0
	1	93.3	95.3	46.0	52.5	71.5	75.5
	3	106.7	108.6	49.5	56.7	77.1	81.2
	5	104.9	106.6	46.2	53.0	72.5	76.2
LLaVA v1.5-7B	Instruction	114.5	122.7	83.9	99.2	106.3	117.5
	0	82.1	87.2	55.4	65.4	77.2	83.9
	1	42.9	45.8	18.2	21.5	31.2	34.4
	3	40.3	43.8	17.5	21.3	25.2	28.1
	5	52.9	56.6	19.6	23.2	30.7	33.9
Idefics2-8B (Base)	Instruction	81.2	94.6	63.0	79.7	81.0	95.3
	0	11.0	36.9	16.1	36.8	13.7	37.0
	1	9.4	31.8	13.4	30.2	12.1	32.6
	3	11.2	33.2	13.5	29.4	15.1	37.6
	5	17.2	38.1	17.3	33.3	17.0	36.9
Idefics2-8B (Instruction-Tuned)	Instruction	83.6	90.1	62.3	74.7	84.3	93.0
	0	102.4	110.4	76.5	91.5	99.7	110.1
	1	94.1	102.9	72.7	88.1	93.6	104.4
	3	91.8	99.6	72.5	85.5	91.6	101.0
	5	88.7	96.1	70.3	82.6	88.1	97.0
OpenFlamingo-9B	Instruction	71.0	82.0	56.2	69.8	67.4	81.5
	0	52.1	66.0	38.8	49.4	43.5	61.4
	1	61.3	65.1	40.1	48.2	53.4	61.6
	3	66.3	70.6	40.7	50.3	56.4	65.2
	5	62.6	66.1	38.6	46.2	50.7	56.8

One can observe that, in the instruction-only scenario, Idefics2 and OpenFlamingo exhibit lower efficiency compared to InstructBLIP and LLaVA models. However, in the in-context learning setting, instruction-tuned Idefics2 follows the same scaling trend as InstructBLIP and LLaVA on Flickr30K and NoCaps. This indicates that Idefics2 (IT) benefits from additional contextual information as efficiently as InstructBLIP models and LLaVA with respect to the training data volume.

In contrast, OpenFlamingo consistently underperforms across all datasets. This finding aligns with Qin et al. (2024) and suggests that the fully autoregressive approach – where visual information is passed as input soft tokens to the LLM – is

more data-efficient than OpenFlamingo’s strategy of integrating visual information directly within the LLM’s layers.

A.5 Qualitative Analysis

We qualitatively analyze how models’ outputs vary between instruction-only scenarios and those with context, using both similar and random demonstrations. Specifically, we investigate whether the models effectively leverage contextual information. To do so, we select four examples from Flickr30K, as well as the demonstrations for these examples, as shown in Figure A3.

Consistent with our quantitative analysis presented in Section 4, InstructBLIP models can generate captions correctly describing the visual con-

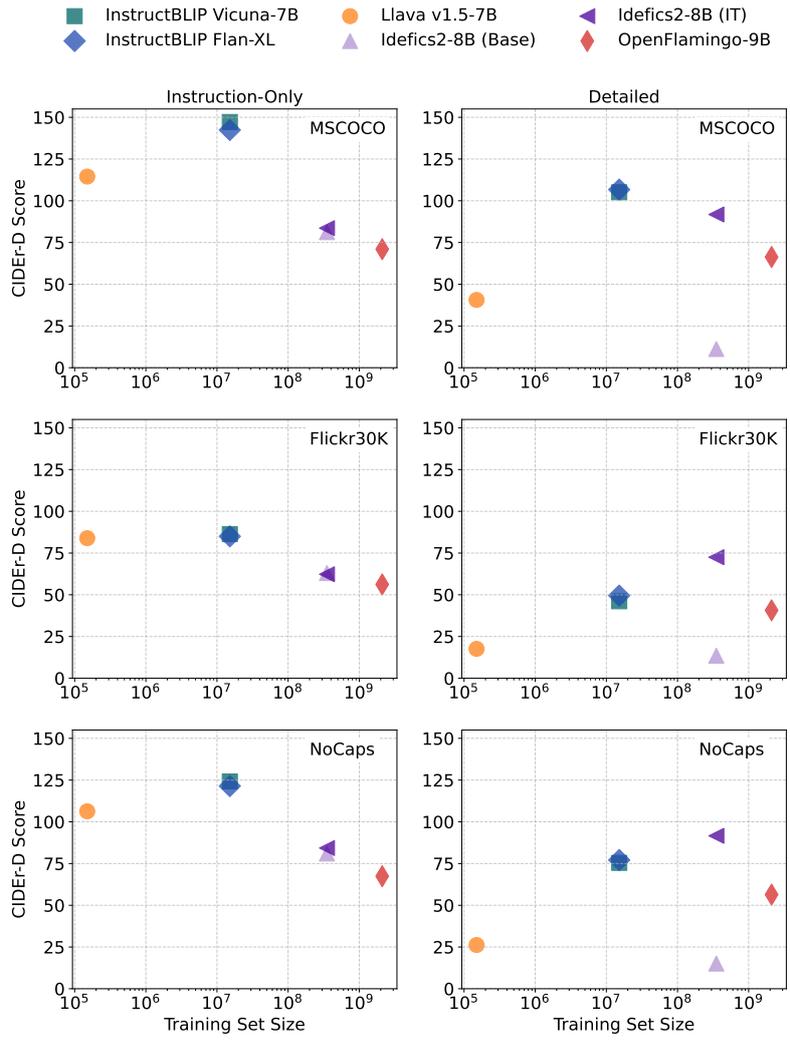


Figure A2: Influence of training dataset size on performance on instruction-only and in-context learning scenarios. Note that the training set size is in the log scale.

Table A4: **Influence of the number of demonstration examples on VQA performance.** Evaluation results for the VQA task varying the number of demonstration examples (Shot) in the context. We use the corresponding training set of each dataset as the knowledge base. We report the VQA accuracy, bold numbers highlight the best performance for each model.

Model	Shot	VizWiz (\uparrow)	GQA (\uparrow)	TextVQA (\uparrow)	OKVQA (\uparrow)
InstructBLIP Vicuna-7B	0	21.8	49.1	33.8	42.4
	1	19.8	47.6	29.6	41.2
	3	19.5	46.3	25.2	39.4
	5	19.5	45.2	23.5	38.2
InstructBLIP FlanT5-XL	0	21.0	48.1	31.1	35.6
	1	32.0	46.4	28.3	35.2
	3	34.0	45.1	26.2	34.7
	5	32.4	43.4	24.4	34.4
LLaVA v1.5	0	15.5	56.3	37.7	27.8
	1	15.2	45.5	10.7	27.1
	3	20.7	45.4	9.6	28.6
	5	19.9	44.9	10.3	28.9
Idefics2-8B (Base)	0	13.08	26.20	35.94	13.98
	1	26.05	18.49	26.34	12.33
	3	28.11	14.91	20.2	6.98
	5	32.98	19.38	21.1	7.6
Idefics2-8B (Instruction-Tuned)	0	23.3	51.5	62.6	38.6
	1	24.4	49.1	56.6	37.8
	3	26.3	49.9	57.2	38.4
	5	27.1	46.8	53.1	36.9
OpenFlamingo-9B	0	13.9	22.5	20.0	10.2
	1	29.0	24.9	19.7	14.9
	3	33.8	29.1	22.8	17.5
	5	20.9	23.1	16.8	12.3

1120 tent of the query image in instruction-only and
1121 in-context learning with similar demonstrations.
1122 However, random demonstrations can confound
1123 the models, leading to captions unrelated to query
1124 images. Then, it confirms the performance gain
1125 shown in Figure 5. Furthermore, LLaVA proves
1126 to leverage visual cues provided within the demon-
1127 strations, often utilizing information from the first
1128 demonstration. Particularly, we can observe in ex-
1129 ample #1 that it incorporates the visual detail of
1130 the number 4 from the player’s jersey from the last
1131 similar demonstration. And, in examples #2 and
1132 #4, it ignores the query image and describes ex-
1133 actly the first similar demonstration. In contrast,
1134 when random demonstrations are used, the model
1135 generates captions unrelated to the query image.

1136 Idefics2 (IT) appears to ignore contextual infor-
1137 mation during caption generation. Interestingly,
1138 when given demonstrations, the model generates
1139 shorter captions that align with the average sen-
1140 tence length found in the MS COCO dataset. This
1141 behavior suggests that the model is attempting to
1142 replicate the structure of the provided examples
1143 rather than leveraging their semantic content. On
1144 the other hand, Idefics2 (Base) struggles with cap-

1145 tion generation under the instruction-only scenario,
1146 often returning the prompt or part of it or even gen-
1147 erating a text totally unrelated to the query image
1148 as in example #2. This indicates a reliance on con-
1149 textual grounding. Nevertheless, its performance
1150 improves when given demonstrations, generating
1151 outputs similar to those of Idefics2 (IT), albeit with
1152 repetition of prompt segments at the end. This
1153 repetition suggests difficulty in processing instruc-
1154 tions effectively. However, after instruction tun-
1155 ing, the model improves its handling of prompts
1156 and generates cleaner captions without extraneous
1157 prompt segments.

1158 OpenFlamingo can generate captions in both
1159 instruction-only and in-context learning settings.
1160 However, when given demonstrations, it sometimes
1161 mixes up elements from different examples, re-
1162 sulting in captions unrelated to the query image.
1163 This is particularly seen in examples #2 and #3,
1164 where random demonstrations lead to captions in-
1165 corporating unrelated elements. This suggests that
1166 while OpenFlamingo can utilize contextual infor-
1167 mation, it lacks a robust mechanism to differentiate
1168 between relevant and irrelevant context, compro-
1169 mising caption generation.

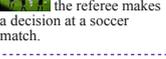
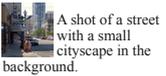
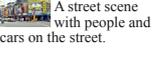
Context	Query	
	Example #1	Example #2
Instruction	<p>Similar Demonstrations</p>  <p>A player is down while others tend to him at a soccer game.</p>  <p>A group of children playing a game of soccer in a green grass covered soccer field.</p>  <p>Teams pause as the referee makes a decision at a soccer match.</p>	<p>Random Demonstrations</p>  <p>A Young child is throwing a big rock into a body of water.</p>  <p>Two people are climbing a rock wall with a rope.</p>  <p>Clown at a circus making a balloon toy for a woman.</p>
	<p>Instruction</p>  <p>A beautiful woman holding a brown dog in her arms near a refrigerator.</p>  <p>The dogs are playing with a yellow frisbee.</p>  <p>Some very cute brown dogs in the road with a toy.</p>	<p>Similar Demonstrations</p>  <p>A man with a tattoo on his arm cooking something in a frying pan.</p>  <p>A kid juggling balls in the middle of traffic.</p>  <p>Two football teams are in the middle of a play as white jerseys and the opposing green jerseys are in all positions of motion while scrambling for the football.</p>
Generated Captions	<p>InstructBLIP Vicuna: a soccer player standing on a field in front of a crowd</p> <p>InstructBLIP Flan-XL: a group of people watching a soccer game on a field</p> <p>LLaVA: A soccer player wearing a green shirt and black shorts.</p> <p>Idefics2 (Base): A short image caption...</p> <p>Idefics2 (IT): A soccer goalie runs on the field in front of thousands of spectators. The goalie is wearing a green jersey.</p> <p>OpenFlamingo: This is a picture of a goalkeeper.</p>	<p>InstructBLIP Vicuna: a group of men playing soccer on a field</p> <p>InstructBLIP Flan-XL: a group of soccer players on a field</p> <p>LLaVA: A man wearing a white jersey with the number 4 on it.</p> <p>Idefics2 (Base): "A soccer goalie warming up before a game." I am an intelligent image captioning bot. [...]</p> <p>Idefics2 (IT): A soccer goalie runs across the field during a game.</p> <p>OpenFlamingo: A goalkeeper stands in front of a soccer goal.</p>
	<p>InstructBLIP Vicuna: a woman playing a guitar in front of a microphone</p> <p>InstructBLIP Flan-XL: a girl in a spider-man costume riding a toy horse</p> <p>LLaVA: Little girl wearing a purple shirt and pink pajamas lying on the floor.</p> <p>Idefics2 (Base): "A soccer goalie warming up before a game." I am an intelligent image captioning bot. [...]</p> <p>Idefics2 (IT): A goalkeeper runs across the field during a soccer game.</p> <p>OpenFlamingo: A goalkeeper stands in front of the goal.</p>	<p>InstructBLIP Vicuna: two dogs playing with a frisbee in a field</p> <p>InstructBLIP Flan-XL: two dogs playing with a frisbee in a grassy field</p> <p>LLaVA: Two dogs fighting over a frisbee in a grassy field.</p> <p>Idefics2 (Base): Lorem Ipsum is simply dummy text of the printing and typesetting industry. Lorem Ipsum has been the industry's standard</p> <p>Idefics2 (IT): Dogs playing in the grass and one has a toy in its mouth.</p> <p>OpenFlamingo: Two dogs playing tug-of-war</p>
Context	<p>Similar Demonstrations</p>  <p>A busy intersection in the city is full of people and signs.</p>  <p>A shot of a street with a small cityscape in the background.</p>  <p>A street scene with people and cars on the street.</p>	<p>Random Demonstrations</p>  <p>A young dirty faced asian girl carrying a pillow.</p>  <p>A woman in a blue shirt and jeans is throwing something away in a dumpster.</p>  <p>A man in a gray shirt talks to a young woman in a white jacket in front of a church.</p>
	<p>Instruction</p>  <p>Guy yawns at the table of wine bottles, glasses and bowl of food.</p>  <p>A couple of men sits around a wooden table.</p>  <p>Some folks having an appetizer and some drinks.</p>	<p>Similar Demonstrations</p>  <p>A man in a black shirt sings into a microphone.</p>  <p>A construction worker working on a railroad wearing an orange shirt.</p>  <p>A man is playing a musical instrument outside on the sidewalk.</p>
Generated Captions	<p>InstructBLIP Vicuna: a couple of people hugging each other in the middle of a city</p> <p>InstructBLIP Flan-XL: a group of people standing in a city square</p> <p>LLaVA: Two people hugging in the middle of a city street.</p> <p>Idefics2 (Base): A short image caption.</p> <p>Idefics2 (IT): These people are in a city. Some are wearing costumes. There are cars on the street. There are stores in the background. The people are</p> <p>OpenFlamingo: A group of people dancing in the middle of a downtown street.</p>	<p>InstructBLIP Vicuna: there are a lot of people walking around on the street</p> <p>InstructBLIP Flan-XL: a busy city street with many people walking around</p> <p>LLaVA: A short caption I can generate to describe example 4 is:</p> <p>Idefics2 (Base): "A group of people standing in front of a building." I am an intelligent image captioning bot. [...]</p> <p>Idefics2 (IT): A group of people stand around in front of a large building.</p> <p>OpenFlamingo: A city street that has a large number of people on it.</p>
	<p>InstructBLIP Vicuna: a man working at a desk with three computer monitors in front of him</p> <p>InstructBLIP Flan-XL: a young boy running away from a bench.</p> <p>LLaVA: An older woman and a young girl are standing at a kitchen counter.</p> <p>Idefics2 (Base): "A group of people are standing in front of a building." I am an intelligent image captioning bot. [...]</p> <p>Idefics2 (IT): A group of people are standing around in front of a large building.</p> <p>OpenFlamingo: A group of people dancing in the street.</p>	<p>InstructBLIP Vicuna: a group of people sitting at a table in a bar</p> <p>InstructBLIP Flan-XL: a group of people sitting around a table in a bar</p> <p>LLaVA: A man and a woman are sitting at a table with wine bottles and glasses in front of them.</p> <p>Idefics2 (Base): A short image caption.</p> <p>Idefics2 (IT): This is as picture taken in a room, there are group of people sitting on chairs in front of the people there is a table on the table there</p> <p>OpenFlamingo: A group of people sitting at a table in a restaurant.</p>

Figure A3: Selected examples from Flickr30K illustrate how models' outputs vary between instruction-only scenarios and those with context, using similar and random demonstrations. The demonstration examples are retrieved from the MS COCO training set, consistent with our experimental pipeline.

Table A5: **Ablation on detailed prompt.** Evaluation results for image captioning task under in-context learning testing detailed prompts with minor changes. Each sample includes three demonstration examples as context, with the MS COCO training set serving as the knowledge base. Instruction refers to the best performance across instruction-only prompts. Bold numbers highlight the best performance for each model.

Model	Prompt	MS COCO		Flickr30K		NoCaps	
		CIDEr-D (↑)	CIDEr-R (↑)	CIDEr-D (↑)	CIDEr-R (↑)	CIDEr-D (↑)	CIDEr-R (↑)
InstructBLIP Vicuna-7B	(1) Base detailed prompt	101.6	104.2	44.3	51.8	72.1	77.0
	(2) Removing the initial phrase	98.5	101.4	42.6	50.4	70.3	75.2
	(3) Removing the word “creative”	106.2	108.5	46.8	53.0	73.6	77.2
	(4) Removing both the initial phrase and the word “creative”	103.1	105.9	46.2	52.5	73.0	76.8
	(5) Removing the name “q-former”	101.1	103.6	42.6	49.8	70.6	75.4
	(6) Replacing “q-former” with a generic name	101.4	103.9	43.7	50.8	70.7	75.5
	(7) Replacing “q-former” with a random name	100.4	102.9	43.2	50.0	71.1	76.0
InstructBLIP FlanT5-XL	(1) Base detailed prompt	0.4	0.4	0.0	0.0	0.5	0.6
	(2) Removing the initial phrase	0.7	0.7	0.1	0.1	2.0	2.2
	(3) Removing the word “creative”	106.7	108.6	49.6	56.7	76.1	80.2
	(4) Removing both the initial phrase and the word “creative”	105.0	106.8	48.6	55.7	75.3	79.3
	(5) Removing the name “q-former”	0.5	0.5	0.0	0.1	1.1	1.2
	(6) (6) Replacing “q-former” with a generic name	0.2	0.2	0.0	0.1	0.4	0.4
	(7) Replacing “q-former” with a random name	0.3	0.3	0.0	0.1	0.2	0.3
LLaVA v1.5	(1) Base detailed prompt	38.9	43.5	18.3	22.2	26.6	30.5
	(2) Removing the initial phrase	29.3	32.0	12.8	15.4	18.4	20.7
	(3) Removing the word “creative”	40.6	44.1	17.6	21.1	26.2	29.4
	(4) Removing both the initial phrase and the word “creative”	29.0	31.4	12.9	15.6	17.0	19.3
	(5) Removing the name “q-former”	39.5	44.5	17.2	21.0	26.0	29.5
	(6) Replacing “q-former” with a generic name	39.4	44.0	17.4	21.3	25.9	29.7
	(7) Replacing “q-former” with a random name	39.8	44.4	17.8	21.6	26.3	29.9
Idefics2-8B	(1) Base detailed prompt	66.5	74.4	53.0	64.2	65.6	74.0
	(2) Removing the initial phrase	69.5	76.7	54.0	64.5	69.7	77.8
	(3) Removing the word “creative”	91.8	99.6	72.5	85.5	91.6	101.0
	(4) Removing both the initial phrase and the word “creative”	87.7	95.4	66.9	80.3	88.4	98.8
	(5) Removing the name “q-former”	65.0	72.7	51.0	61.6	64.0	72.3
	(6) Replacing “q-former” with a generic name	64.0	72.3	51.5	62.6	62.5	70.9
	(7) Replacing “q-former” with a random name	67.6	75.9	54.4	65.4	67.3	76.0
OpenFlamingo-9B	(1) Base detailed prompt	66.9	72.5	43.9	53.8	57.5	66.9
	(2) Removing the initial phrase	67.2	72.9	44.4	53.8	58.8	68.0
	(3) Removing the word “creative”	66.3	70.6	40.7	50.3	56.4	65.2
	(4) Removing both the initial phrase and the word “creative”	66.2	70.8	43.2	52.5	58.3	66.9
	(5) Removing the name “q-former”	65.6	71.4	43.0	52.5	56.4	66.0
	(6) Replacing “q-former” with a generic name	66.3	72.1	44.4	54.1	56.9	67.1
	(7) Replacing “q-former” with a random name	66.5	72.3	42.7	51.8	56.6	65.9

A.5.1 VizWiz

To investigate why adding in-context demonstrations benefits performance on the VizWiz dataset, unlike other datasets, we first analyze the distribution of answers in this dataset. Figure A4 shows the answers that appear more than 10 times in the dataset. A clear imbalance can be seen, with “unanswerable” being by far the most frequent answer. It seems approximately 1,750 times, whereas all other answers occur fewer than 100 times, with most appearing only once. We hypothesize that this imbalance explains the performance gains observed when the number of demonstrations increases. To test this, we examine the top 10 most frequently generated answers from each model in the 0-shot setting, where only the question is provided, and the 5-shot setting, as shown in Figure A5.

In the 0-shot setting, the distribution of answers is more balanced, and in most cases, “unanswerable” does not even appear in the top 10. However, when demonstrations are included, “unanswerable” becomes the most frequent answer, with a significant margin over others for most models. Notably, InstructBLIP FlanT5-XL, Idefics2

(Base), and OpenFlamingo generate “unanswerable” most often, with InstructBLIP FlanT5-XL and OpenFlamingo outputting it incorrectly many times. This suggests that demonstrations labeled as “unanswerable” influence the models to replicate this response, leading to improved performance due to the dataset’s strong bias toward this answer.

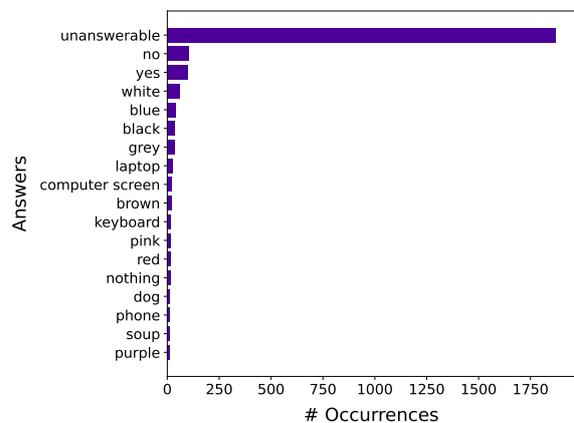


Figure A4: Answers that occur more than 10 times in the VizWiz dataset.

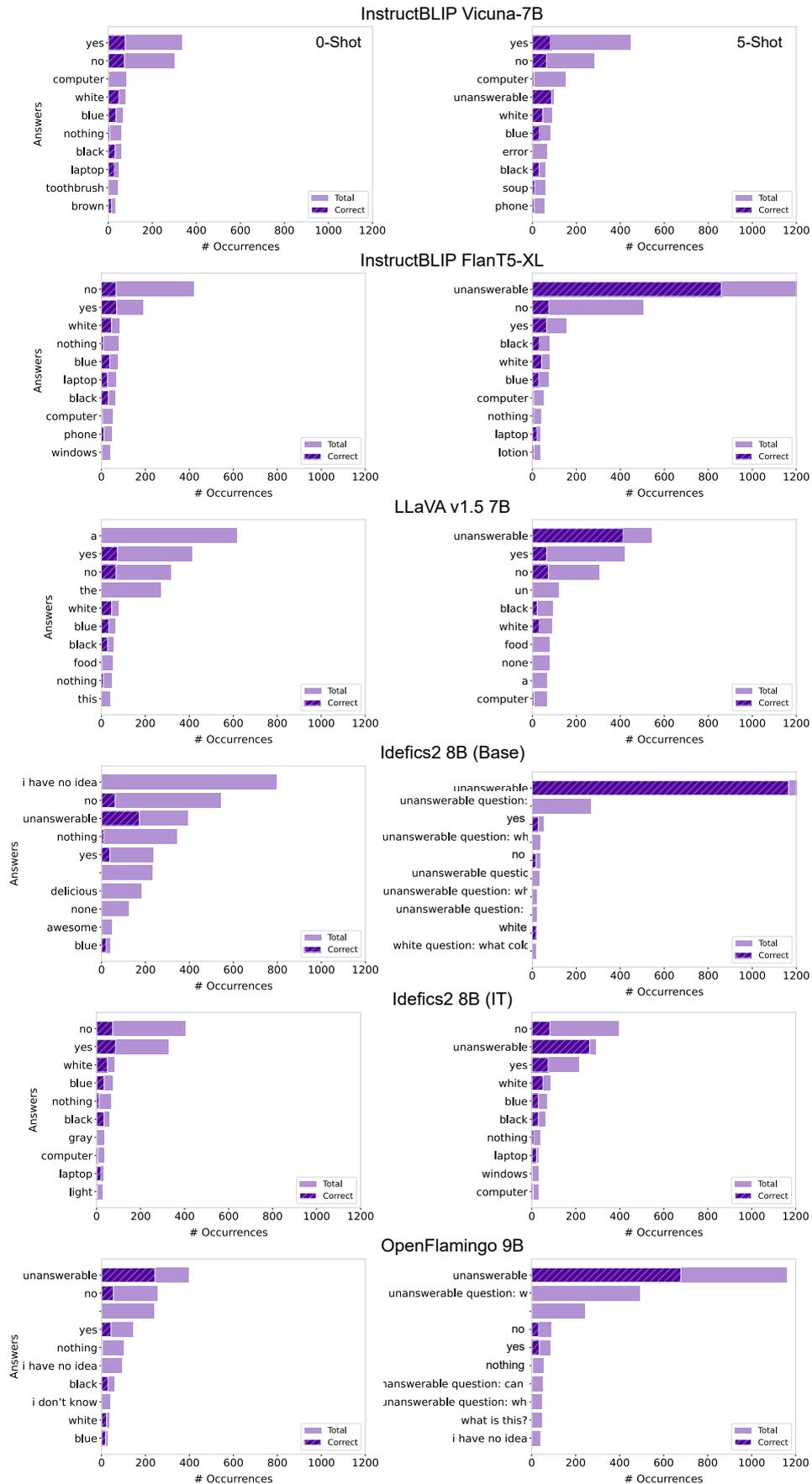


Figure A5: Top 10 most frequent answers for each model on the VizWiz dataset in both 0-shot and 5-shot scenarios. “Total” refers to the total number of occurrences of a given answer and “Correct” indicates the number of correct ones.

A.6 More Details on Experimental Setup

To facilitate the reproducibility of our work, we report in Table A6 the models we analyzed, along with details on their number of parameters and training set size. Table A7 shows the datasets used in our experiments, including basic statistics on their size. Additionally, we outline the main hyperparameters used in our experiments in Tables A8 and A9. Table A8 lists the hyperparameters specific to image captioning, while Table A9 includes those used for VQA. We conducted our experiments in a heterogeneous computing environment; however, the majority were performed on a single Quadro RTX 8000 GPU. Also, all experiments were conducted only once.

Table A6: VLMs investigated in this work. For each model, we report the number of parameters and the size of the dataset used for training.

Model	#Params (B)	Training set size (M)
Llava v1.5-7B	7.1	0.15
InstructBLIP Vicuna-7B	7.9	15.1
InstructBLIP Flan-XL	4.0	15.1
Idefics2-8B	8.4	351.2
OpenFlamingo-9B	8.1	2,101.0

Table A7: Datasets used in our experiments. For each dataset, we report the number of samples in each split and the specific task it is used for. Note that we do not use Flickr or NoCaps training sets, as we rely on the MS COCO training set as the knowledge base for these datasets. “Val.” stands for the validation dataset.

Dataset	Size	Task
MS COCO	Train: 118.2K/ Val: 5.0K	Image Captioning
Flickr30K	Val: 1.0K	Image Captioning
NoCaps	Val: 4.5K	Image Captioning
VizWiz	Train: 20.5K/ Val: 4.3K	VQA
GQA	Train: 943K/ Val: 12.5K	VQA
TextVQA	Train: 34.6K/ Val: 5K	VQA
OKVQA	Train: 9K/ Val: 5K	VQA

Table A8: Hyperparameters for image captioning.

Hyperparameters	Value
# Beams	5
Max. New Tokens	30
Min. Length	10
Repetition Penalty	1.5
Length Penalty	1.0
Temperature	1.0

Table A9: Hyperparameters for VQA.

Hyperparameters	Value
# Beams	5
Max. New Tokens	10
Min. Length	1
Length Penalty	-1.0