

Visual Program Distillation: Distilling Tools and Programmatic Reasoning into Vision-Language Models

Anonymous CVPR submission

Paper ID ***

Abstract

Solving complex visual tasks such as “Who invented the musical instrument on the right?” involves a composition of skills: understanding space, recognizing instruments, and also retrieving prior knowledge. Recent work shows promise by decomposing such tasks using a large language model (LLM) into an executable program that invokes specialized vision models. However, generated programs are error-prone: they omit necessary steps, include spurious ones, and are unable to recover when the specialized models give incorrect outputs. Moreover, they require loading multiple models, incurring high latency and computation costs. We propose **Visual Program Distillation (VPD)**, an instruction tuning framework that produces a vision-language model (VLM) capable of solving complex visual tasks with a single forward pass. VPD distills the reasoning ability of LLMs by using them to sample multiple candidate programs, which are then executed and verified to identify a correct one. It translates each correct program into a language description of the reasoning steps, which are then distilled into a VLM. Extensive experiments show that VPD improves the VLM’s ability to count, understand spatial relations, and reason compositionally. Our VPD-trained PaLI-X outperforms all prior VLMs, achieving state-of-the-art performance across complex vision tasks, including MMBench, OK-VQA, A-OKVQA, TallyQA, POPE, and Hateful Memes. An evaluation with human annotators also confirms that VPD improves model response factuality and consistency. Finally, experiments on content moderation demonstrate that VPD is also helpful for adaptation to real-world applications with limited data.

1. Introduction

Vision-language models (VLMs) have become the pre-trained backbone for many computer vision tasks [2, 4, 7, 9–11, 34, 39, 41, 59, 62, 72, 78]. Yet, all these models still fall short of solving numerous visual reasoning tasks expected of competent vision models. Even state-of-the-art (SOTA) proprietary vision-language models such as GPT-4V [49]

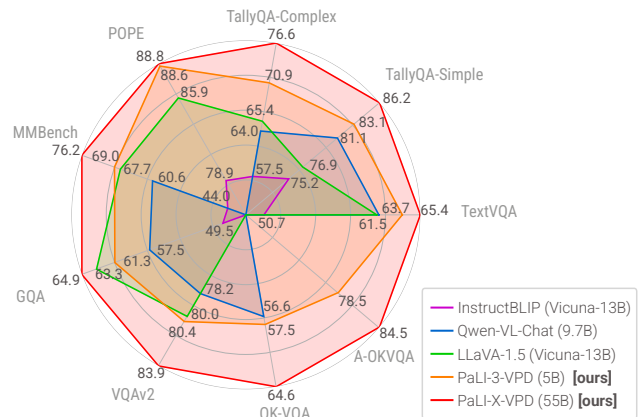


Figure 1. We introduce *Visual Program Distillation (VPD)*, a training framework which leverages LLM-generated programs to synthesis multimodal chain-of-thought training data for Vision-Language Models (VLMs). Our generalist models trained with VPD, PaLI-3-VPD (5B) and PaLI-X-VPD (55B), outperform prior VLMs on a broad range of tasks, while producing human-interpretable and faithful reasoning steps.

do not do well at tasks that involve counting and spatial reasoning [69]. They find it difficult to count (TallyQA [1]), to compositionally reason (GQA [26]), and to reason with external knowledge (OK-VQA [44], A-OKVQA [53]). Many of these tasks require VLMs to conduct compositional reasoning, which still remains an unsolved challenge. For instance, answering the question “Who invented the musical instrument on the right?” involves a composition of skills: identifying objects, applying spatial reasoning to locate the one on the right, recognizing the musical instrument, and accessing prior knowledge to retrieve the inventor.

In contrast, large language models (LLMs) have demonstrated remarkable performance at generating code that solves complex and compositional tasks [3, 8, 43, 49, 67]. Several recent papers [16, 20, 25, 56] capitalize on this by prompting LLMs to generate programs where each step corresponds to a reasoning step. The programs invoke specialized “tools” (or specialized vision models) to explicitly execute each reasoning step. For the question above, the

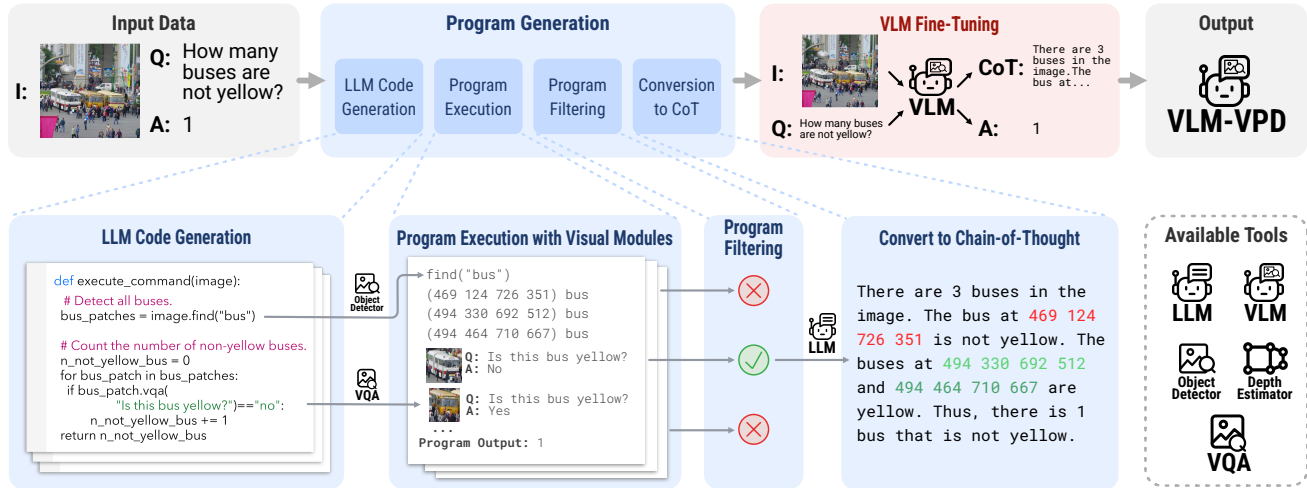


Figure 2. Overview of Visual Program Distillation (VPD). VPD uses an LLM and specialized vision tools to generate faithful chain-of-thought (CoT) training data for vision-language models (VLMs). Given a multimodal input, our 4-step data synthesis pipeline generates a CoT that answers the query. In the example above, our synthesized CoT contains a series of reasoning steps: find all buses, check if each bus is yellow, and aggregate the count into a final answer. The CoT also contains the grounding information given by object detection.

056 program would call an “object detector” tool to identify and
 057 isolate all the objects, a “fine-grained object classification”
 058 tool to recognize the musical instrument, and a “knowledge-
 059 based question answering” tool to retrieve its inventor.

060 Although innovative, generating explicit programs is compu-
 061 tationally expensive in practice, prone to errors, and still
 062 underperforms end-to-end models. Programs require load-
 063 ing and executing multiple tools, leading to high latency
 064 and computational cost. Moreover, generated programs may
 065 omit necessary steps or include spurious ones. Even when
 066 the program is correct, vision model invocations can produce
 067 incorrect outputs, from which the overall program cannot
 068 recover. Unfortunately, empirical results show that visual
 069 programs still fall short of end-to-end fine-tuned models [16].

070 Another line of work is visual instruction tuning [39],
 071 which tries to distill the instruction-following ability of
 072 LLMs into VLMs. They prompt LLMs with image captions
 073 and bounding box annotations in order to generate queries
 074 and answers that are used to fine-tune VLMs [7, 23, 39, 64].
 075 However, this approach has important limitations: image
 076 captions can miss fine-grained visual information, and LLM
 077 are prone to produce inconsistent outputs for custom vision
 078 representations like bounding boxes [18]. As a result, ex-
 079 isting instruction-tuned VLMs still struggle with tasks that
 080 requires complex visual reasoning [7, 14, 38, 40].

081 In this work, we present **Visual Program Distillation**
 082 (VPD), a novel distillation method that induces LLM-like
 083 complex reasoning capabilities into vision-language models
 084 (Fig. 2). As the name suggests, VPD combines two key
 085 insights to deliver a training paradigm that surpasses the sum
 086 of its parts: It relies on (1) advancements in visual programs
 087 that use tools [20] and (2) the recent breakthroughs in dis-

088 tillation through chain-of-thought reasoning [21]. Given a
 089 labeled training dataset of complex visual tasks, VPD gener-
 090 ates a correct program, and then distills its reasoning steps
 091 into vision-language models. To avoid using programs that
 092 give the wrong answer, VPD prompts an LLM to generate
 093 multiple candidate programs, and executes every one of them.
 094 When labeled data is available, it then filters for programs
 095 that produce the correct answer upon execution. Therefore,
 096 our programs comprise multiple vision tools, are executable,
 097 and yield the correct answer when executed. Next, VPD
 098 rewrites the correct programs as natural language chain-of-
 099 thought instructions and uses step-by-step distillation [21]
 100 to inject the reasoning abilities into VLMs.

101 Our best instruction-tuned model, PaLI-X-VPD
 102 (55B), sets a new SOTA result on 8 classical VQA tasks
 103 and 2 zero-shot multimodal benchmarks. Our models even
 104 outperform the recent SOTA established by PaLI-X [10].
 105 Importantly, we conduct a quality evaluation with human
 106 raters which shows that PaLI-X-VPD generates more con-
 107 sistent and faithful rationales compared to its counterpart
 108 trained using instruction-tuning data. In addition, we experi-
 109 ment with PaLI-3 (5B), and show that VPD also improves
 110 the performance of smaller-scale models. Finally, experi-
 111 ments on Hateful Memes [29] show that VPD is also helpful
 112 for adapting to new tasks, even when no labels are available.

113 2. Related work

114 VPD is a general method for improving any vision-language
 115 model, and includes automatic program generation and train-
 116 ing with chain-of-thought data as steps of the proposed
 117 framework. We discuss each of these research areas.

118 **Vision-language models (VLMs).** Most recent generative

VLMs share a common structure: a pre-trained visual encoder, a pre-trained LLM, and a connector between the two modalities [2, 4, 9, 34, 39, 41, 59, 62, 72, 78]. The models are trained on large-scale image-text pairs from various tasks to adapt to both modalities. They are also tuned on LLM-generated visual instructions to enable models to follow versatile instructions from users [39]. Some models also include bounding boxes in the pre-training data to improve the VLM on visually grounded tasks [6, 7, 10, 41, 50, 60, 77]. The bounding boxes are usually retrieved from in COCO [37] and Visual Genome [32]. Different from prior work, our method does not rely on provided dense annotations. We use LLM-generated code and specialized vision tools to generate our own visual instruction data.

Visual programming and agents. With the advancement of large language models (LLMs) [3, 5, 13, 49, 58], recent work has started using LLMs as an interface to solving complex reasoning tasks with tools [12, 42, 51, 54, 75], and also as an agent for vision tasks [23, 25, 68, 70]. From this line of work, the most relevant to us are VisProg [20] and ViperGPT [16], which leverage LLMs to generate executable programs with a sequence of invocations to specialized vision tools. They achieve SOTA zero-shot performance on various vision tasks, while being versatile and interpretable.

Training and inference with chain-of-thought. Chain-of-Thought (CoT) [65] has become a popular approach to improving LLM performance. Recent work such as Program-of-Thoughts (PoT) [8] and Faithful CoT [43] further improve this framework by splitting inference into two steps: first generate a program with LLM, then execute the program. This approach achieves better accuracy and reduces hallucinations in the reasoning steps. Moreover, CoT is also used to train language models. Distill step-by-step [21], PaD [79], and SCOTT [61] train smaller language models with CoT generated by LLMs, showing that this can improve model performance and reasoning consistency. Mammoth [74] trains an LLM with a hybrid of CoT and PoT rationales and achieves a SOTA model for math problems.

3. Visual Program Distillation (VPD)

We introduce VPD, a general model-agnostic framework that distills the reasoning power of LLM-generated programs together with the low-level image understanding abilities of vision tools into a single vision-language model (Fig. 2). At its core, VPD consists of two major steps:

1. **Program generation and verification:** Given a textual query q and a visual input i , VPD first generates a program that solves the query by making use of vision modules and tools, then converts the program execution trace into a chain-of-thought c (§3.1).
2. **Distilling step-by-step:** The visual input i , textual query q , and chain-of-thought c produced by the previous step

are distilled into a vision-language model (VLM) (§3.2). 170

3.1. Program generation and verification 171

Our training data synthesis pipeline is illustrated in the blue boxes of Fig. 2, and consists of four stages. Given a sample (i, q, y) consisting of a visual input i and a textual query q about its content and, when available, its ground truth answer y , we perform the following sequence of steps: 172

1. **Program generation with LLM:** Given q , we generate a list of k candidate programs $\pi(q) = \{z_1, z_2, \dots, z_k\}$, where π represents a program generation function. 173
174
175
176
177
178
179
2. **Program execution with vision modules:** We execute each program z_i with an execution engine ϕ to obtain its final result $\phi(i, z_i) = \hat{y}_i$. However, during program execution, we maintain the execution trace t_i recording all the intermediate function calls and outputs. At the end of this step, we produce a list of programs, results, and execution traces $\{(z_1, \hat{y}_1, t_1), \dots, (z_k, \hat{y}_k, t_k)\}$. 180
181
182
183
184
185
186
3. **Program filtering:** Among the k candidate programs from the previous step, we keep a single tuple (z, \hat{y}, t) with correct answer. 187
188
189
4. **Converting program execution traces into chains-of-thought:** We rewrite t into a CoT c using an LLM. 190
191

We now discuss in detail each of the steps above. 192

Program generation. We adopt a similar approach with recent work [16, 20] in our program generation step, and use PaLM-2 [3] as LLM to generate candidate programs for a given query q . We prompt PaLM-2 with the same kind of text prompt as used by ViperGPT, which contains a detailed description of the available vision modules, followed by the query q (prompt in Appendix §D). The LLM is expected to directly output a Python function definition that will be executed in the following steps. However, in our experiments, we find that the success rate of top- k programs is much higher than the top-1 program. Therefore, in contrast to prior work, which only samples one program z , when the ground truth answer y is available, we set a temperature T for LLM decoding and sample a list of top- k candidate programs $\{z_1, z_2, \dots, z_k\}$ from the LLM. Then, we filter out one correct program z in later steps. As shown in our ablation in §4.2, this becomes crucial to our performance gain. We use $T = 0.5$ and $k = 5$ for all our experiments. For unlabeled data, $k = 1$ and the filtering step can be skipped. 193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211

Program execution with vision modules. We use the same execution engine ϕ as ViperGPT [16]. An LLM-generated program z_i is a Python function that takes the visual input i as input. While the LLM outputs the program as a sequence of text, the execution engine ϕ is able to interpret this as a Python program and execute it, to obtain its return result \hat{y}_i . Additionally, ϕ also records the execution trace t_i of the program z_i , which keeps a record of all the vision module calls, their inputs, and outputs. We use the following tools for 212
213
214
215
216
217
218
219
220

221 the program: PaLI-X [10] for simple visual queries, PaLI-X
222 detection (distilled from OWLv2 [47]) for object detection;
223 Google Cloud Depth API [17] for depth estimation, and
224 PaLM-2 [3] for external knowledge.

225 **Program filtering.** As we discussed in §1, visual programs
226 are error-prone for a variety of reasons. The program might
227 be wrong, and the execution process can introduce additional
228 errors. To overcome these issues, we employ a program filter-
229 ing step. We start by sampling top- k programs and attempt
230 to execute each one. During this process, any program that
231 fails execution is instantly discarded. We further filter the re-
232 maining programs, the strategy depending on the availability
233 of labeled data. For tasks where human labels are available,
234 we select a single program per input sample whose answer
235 is correct—that is, when the program output \hat{y}_i matches the
236 human label y . In this case our visual program pipeline acts
237 like a CoT rationale annotator. One potential problem is that
238 some task answers are ambiguous. For example, for a ques-
239 tion like “Where are the horses?”, the answers “mountain”
240 and “mountains” are both correct. We adopt the method in
241 [28] and use an LLM to determine if the program output
242 is correct (details in §D). If there is more than one correct
243 program per question, we select the top scoring one, ac-
244 cording to the scores provided by the program-generating
245 LLM ϕ . If no programs pass the test, the rated sample is not
246 wasted—we simply use the correct answer y as supervision
247 in our fine-tuning stage, without an associated CoT. When
248 no human-rated answer is available, we directly use the top
249 scoring executable program¹.

250 **Converting program execution traces into chain-of-**
251 **thought rationales.** After the filtering step, for each visual
252 input i and query q , we will have selected at most one pro-
253 gram z together with its execution result \hat{y} and trace t . Since
254 most existing VLMs have been pre-trained with text in nat-
255 ural language and not code, we use an LLM to rewrite the
256 execution trace t into a natural language CoT c for our VLM
257 distillation. Some examples are shown in §A. Concretely,
258 similar to prior work [23, 24], we hand-craft 20 examples
259 of how an input (q, z, t) can be converted into a CoT, and
260 use these as few-shots for prompting PaLM-2 [3], which
261 performs in-context learning and generates CoTs for new
262 examples. We include a concrete example in §D.

263 3.2. Distilling Step-by-Step

264 In this step, we fine-tune a backbone VLM with the training
265 data generated in §3.1, distilling the knowledge and reason-
266 ing steps of our generated programs into a single end-to-end
267 model. We do so in a multitask fashion, where the VLM is
268 simultaneously fine-tuned on data synthesized for multiple
269 types tasks (e.g., free-form VQA, multiple choice).

¹Studying ways to approximate program correctness strategies for unlabeled data is an interesting direction for future work.

270 Let f represent our VLM model. While the same VLM ar-
271 chitecture could solve all these tasks, it needs to be prompted
272 differently to adapt to the task. Following prior work [39],
273 we manually design instructions for each task. For exam-
274 ple, for free-form VQA queries, the instruction is “Answer
275 with a single word or phrase”, while for multiple-choice
276 queries we use “Answer with the option letter from the given
277 choices directly”. During fine-tuning, we combine the train-
278 ing samples generated for all tasks into a single dataset, so
279 to account for the different types of tasks, we augment each
280 sample with its corresponding task-specific prompt p .

281 Using these instructions along with the data gener-
282 ated in §3.1, we put together the training dataset $\mathcal{D} =$
283 $\{(i_j, q_j, \hat{y}_j, c_j, p_j)\}_{j=1}^N$, where N is the total number of sam-
284 ples, i_j is the visual input, q_j is the textual query, \hat{y}_j is the
285 visual program output², and c_j is the CoT rationale.

286 We train f to minimize a loss for predicting both the label
287 and the rationale. As shown in the red box of Fig. 2, similar
288 to [21], we treat predicting the output label \hat{y}_j and the ratio-
289 nale c_j as two separate optimization goals. However, since
290 VLMs are open-ended text generation models, they need
291 additional prompting to indicate whether we want a short
292 answer or a long answer that includes a rationale. There-
293 fore, we append the suffix $s_c =$ “Explain the rationale to
294 answer the question” at the end of the prompt for generating
295 a CoT, and use the task instruction p_j for short answers. Our
296 optimization objective is:

$$297 \mathcal{L} = \mathcal{L}_{label} + \mathcal{L}_{rationale} \quad (1)$$

$$298 = \sum_{j=1}^N \ell(f(i_j, q_j, p_j), \hat{y}_j) + \ell(f(i_j, q_j, s_c), c_j) \quad (2)$$

299 Here ℓ is the cross entropy-loss normalized by sequence
300 length to ensure that the labels (typically short) and ratio-
301 nales (typically long) have similar weights. $\mathcal{L}_{rationale}$ both
302 teaches the VLM to generate faithful reasoning steps similar
303 to program execution traces, and carries more information
304 beyond the label \hat{y} that also helps the VLM in better predict-
305 ing the label. During test time, the rationale generation is
306 not required. We can adjust the task instruction to directly
307 get the short output label using p , and using s_c to get the
308 human-interpretable reasoning steps if needed.

309 4. Experiments

310 In this section, we demonstrate VPD’s effectiveness by using
311 it to train a generalist VLM. We attempt this for two VLMs
312 with different scales, PaLI-3 (5B) [11] and PaLI-X
313 (55B) [10]. Detailed experimental setups are given in §4.1.
314 Qualitatively, models fine-tuned with VPD exhibit the ability
315 to reason step-by-step like a program, as illustrated in Fig. 3,

²When labeled data is available, this is equivalent to the ground truth label y due to our filtering strategy.

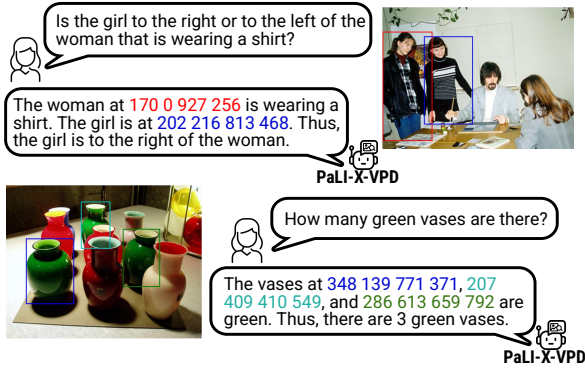


Figure 3. Outputs of PaLI-X trained with VPD. Training with execution traces of filtered programs improves the model’s ability to count, understand spatial relations, and reason compositionally. The images are from Visual Genome [32].

and supported by human evaluation results (§4.3). Quantitative results show that our PaLI-3-VPD and PaLI-X-VPD achieve new SOTA on a broad range of tasks, on both generalist and per-task fine-tuning settings (§4.2). We also conduct a detailed analysis of the source of performance gain (§4.2). Finally, we conduct human evaluation on the quality of the rationales generated by our models, and compare it with rationales generated by an instruction-tuned model trained without VPD (§4.3).

4.1. Experimental setup

Backbone models. We use two state-of-the-art VLMs, PaLI-3 (5B) [11] and PaLI-X (55B) [10] as our base models. Both take images and text as input, and generate text as output. For simplicity, we further refer to them as “PaLI model” when we discuss steps performed with each of them individually.

Data for Generalist Models. We fine-tune the pre-trained PaLI model on two types of datasets to make it a *generalist* VLM—that is, a model that performs relatively well on any task without further training on that specific task:

(1) *Multimodal Instruction-Tuning (MMIT) tasks*, created in the spirit of Self-Instruct [63]. An LLM is prompted with image captions, and generates task inputs and the desired outputs about the corresponding image. Details of MMIT tasks are covered in [64]. Note that these tasks cover a wide variety of common use cases, but do not include the specific in-domain tasks used in this work.

(2) *Academic task-oriented VQA tasks*. Since image captions contain only a coarse description of visual information in the image, and may miss details that are important for solving the task. Additionally, LLM are likely to hallucinate during this data curation process. To further boost the accuracy of our VLMs, we also fine-tune the PaLI models with academic task-oriented VQA tasks. The data mixture covers subsets of a wide variety of VQA tasks, including general VQA (VQAv2 [19]), optical character recognition

(OCRVQA [48]), compositional questions and reasoning (GQA [26]), counting (TallyQA [1]), and VQA that involves external knowledge (OK-VQA [45] and A-OKVQA [52]). The tasks contain a textual query and a short expected label. We use the pipeline in §3.1 to synthesize CoT reasoning steps for these labels, and tune the PaLI model with these the loss in Equation 1. Notice that sometimes the pipeline fails to find a program that generates the correct answer. In that case, we set $\mathcal{L}_{\text{rationale}}$ to 0 and only keep $\mathcal{L}_{\text{label}}$. §4.2 shows how many programs are kept after the filtering stage. More details of the training data mixture is in §B.

Data for specialist models. While fine-tuning the generalist model with VPD, we only use a subset of each task’s training data. To evaluate our model’s ability on each individual task, we continue fine-tuning PaLI-3-VPD and PaLI-X-VPD on each individual task on the training splits.

Training setup. Both PaLI-3 and PaLI-X follow an encoder-decoder architecture, where images are encoded into visual tokens via a visual encoder and then passed into a UL2 model. Due to resource constraint, we use LoRA [22] to fine-tune both PaLI models. Specifically, we add LoRA weights on each linear layer in the attention blocks and the MLP blocks for both the encoder and decoder in the UL2 transformer. More training details are in §B.

Evaluation setup. We evaluate our models on a wide range of tasks, including various VQA tasks and recent zero-shot VLM benchmarks. Noted that A-OKVQA contains two kinds of questions, multiple-choice (MC) and direct answer (DA). We report results on both. TallyQA [1] contains complex counting questions that involve object relationships, attribute identification, and reasoning to get the correct answer. It contains two evaluation partitions, simple and complex. TextVQA [55] focuses reading texts. We include it to evaluate our models on zero-shot tasks. In addition to VQA tasks, we also test our models on two popular VLM benchmarks. POPE [36] focuses on VLM hallucination, containing binary questions of whether or not an object exists in the image. MMBench [40] is a robust and comprehensive VLM benchmark testing a range of fine-grained abilities (e.g., object localization, attribute recognition, spatial relationship). Details of the evaluation sets and metrics are in §B.

Baselines. We refer to the two PaLI models fine-tuned with VPD as PaLI-3-VPD and PaLI-X-VPD, and compare them with various baselines. To evaluate the effectiveness of VPD, we experiment with removing the synthesized CoT from our data mixture, and train the PaLI models with the exact same hyper-parameters and steps. We call these models as PaLI-3 Instruct and PaLI-X Instruct. For a fair comparison, these models are also trained with the same supervised loss for predicting the ground truth answers, on the same images and textual queries as our VPD variant. Moreover, we also compare with the most recent SOTA

Generalist Models	VQAv2	GQA	OK-VQA	A-OKVQA		TallyQA		TextVQA	POPE	MMB
				MC	DA	Simp.	Comp.			
<i>Prior generalist VLMs</i>										
Flamingo (80B) [2]	82.0	-	57.8*	-	-	-	-	57.1	-	-
MiniGPT-4 (Vicuna-13B) [78]	-	43.5	-	67.2	-	-	-	-	-	42.3
InstructBLIP (Vicuna-13B) [14]	-	49.5*	-	-	-	75.2*	57.5*	50.7*	78.9	44.0
Shikra (Vicuna-13B) [7]	77.4	-	47.2	-	-	-	-	-	84.7	58.8
Qwen-VL (9.7B) [4]	78.8	59.3	58.6	-	-	82.6*	65.8*	63.8	-	38.2
Qwen-VL-Chat (9.7B) [4]	78.2	57.5	56.6	-	-	81.1*	64.0*	61.5	-	60.6
mPLUG-Owl2 (8.2B) [72]	79.4	56.1	57.7	-	-	-	-	54.3*	86.2	64.5
LLaVA-1.5 (Vicuna-13B) [38]	80.0	63.3	-	-	-	76.9*	65.4*	61.3*	85.9	67.7
<i>Our instruction-tuned baselines</i>										
PaLI-3-Instruct (5B)	79.9	59.7	56.7	78.3	57.6	81.9	70.4	63.3*	87.7	68.6
PaLI-X-Instruct (55B)	83.6	63.3	64.3	84.1	61.5	85.5	75.4	65.0*	88.9	75.0
<i>Our visual program distillation models</i>										
PaLI-3-VPD (5B)	<u>80.4</u>	<u>61.3</u>	<u>57.5</u>	<u>78.5</u>	56.5	<u>83.1</u>	<u>70.9</u>	<u>63.7*</u>	<u>88.6</u>	<u>69.0</u>
PaLI-X-VPD (55B)	83.9	64.9	64.6	84.5	62.7	86.2	76.6	65.4*	88.8	76.2

Table 1. Comparison with SOTA generalist VLMs. PaLI-X-VPD outperforms all prior models. VPD improves performance on 8/9 tasks for both PaLI-3 and PaLI-X, and is particularly effective on counting questions (TallyQA), compositional questions (GQA), and the comprehensive benchmark (MMBench). Underline indicates when PaLI-3-VPD outperforms the Instruct version. * marks tasks unseen during training. POPE and MMBench are zero-shot benchmarks for all models.

404 vision-language models. These VLMs are initialized with
405 pre-trained visual encoders and LLMs, and then trained with
406 image-text pairs, LLM-generated data, and academic tasks.

407 4.2. Quantitative results

408 **Generalist model.** Table 1 compares VPD with the base-
409 lines discussed in §4.1. We infer all answers by open-ended
410 generation with the prompt “Answer with a single word or
411 phrase.”, using greedy decoding without any constraint on
412 the model’s output space. For multiple-choice questions, we
413 run inference with the prompt “Answer with the option letter
414 from the given choices directly.” and generate the option letter.
415 As Table 1 shows, **PaLI-X-VPD sets the new state-of-**
416 **the-art on all benchmarks.** Compared with prior generalist
417 VLMs, it achieves significant improvement on MMBench
418 (+8.5), TallyQA complex (+9.8), and A-OKVQA(+9.5
419 compared with specialist SOTA[14]). PaLI-3-VPD, despite
420 its small architecture size (5B), outperforms all prior
421 models on VQAv2, TallyQA, POPE, and MMBench, including
422 much larger models that use Vicuna-13B as backbones [7, 14, 38].
423 However, its performance on knowledge-based VQA tasks does not
424 outperform the prior SOTA. Our hypothesis for this is that its
425 language model (3B UL2) is too small to contain all the necessary
426 knowledge needed for correctly solving these tasks.

428 When compared to their Instruct variants, both
429 PaLI-3-VPD and PaLI-X-VPD outperform on 11/12
430 tasks. Specifically, for PaLI-X, VPD obtains a +1.6
431 improvement on GQA (which is heavily focused on compositional
432 questions, spatial relationship, and localization), +1.2
433 on TallyQA for complex questions, and +1.2 on MMBench.
434 These results suggest that **VPD is a more effective method**
435 **for creating instruction-tuning data which enables VLMs**
436 **to improve their visual reasoning ability.**

Per-task fine-tuning (specialist). The results for the specialist
437 models are shown in Table 2. **PaLI-X-VPD sets a**
438 **new SOTA on all benchmarks.** Note how the specialist
439 models tend to have higher scores than the generalist one.
440 We propose several hypotheses for this performance gain:
441 (1) As shown in Table 2, the score gap tends to be larger on
442 free-form VQA tasks. This may be due the fact that human
443 annotations on these datasets have ambiguities. For example,
444 for the question “Who is looking up?”, GQA [26] labels
445 are “man” or “woman” while OK-VQA [45] have more
446 detailed labels, for example, “worker” or “cook”. Per-task
447 fine-tuning alleviates this annotation ambiguity and lets the
448 model focus on the annotation style of that task. For multiple-
449 choice and counting tasks, the answer has less ambiguity, and
450 the score gaps are much smaller. (2) The performance gap
451 between PaLI-3 (5B) *specialist* and *generalist* is larger
452 than that of PaLI-X. We hypothesize that this shows models
453 with larger scale has more multi-task capacity; (3) Per-task
454 fine-tuning adds more in-domain training data, which gener-
455 ally improves performance. 456

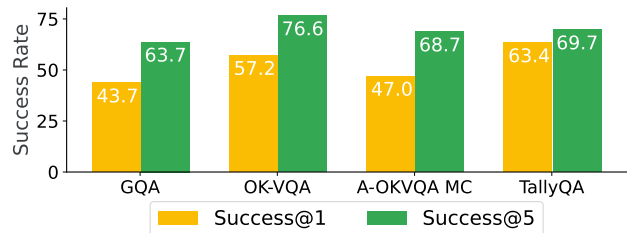


Figure 4. Success rate of top-1 program and top-5 programs on the training set during our data synthesis process.

Analysis: sampling multiple programs is key to good data
457 **generation.** Fig. 4 shows the success rate of finding at least
458

Specialist Models	GQA	OK-VQA	A-OKVQA			TallyQA		
	test-dev	val	Multi-choice val	Direct Answer test	val	Simple test	Complex test	
InstructBLIP (Vicuna-7B) [14]	-	62.1	75.7	73.4	64.0	62.1	-	-
PaLI-3 (5B) [11]	-	60.1	-	-	-	-	83.3	70.5
PaLI (17B) [9]	-	64.5	-	-	-	-	81.7	70.9
CogVLM (17B) [62]	65.2	64.7	-	-	-	-	-	-
PaLI-X (55B) [10]	-	66.1	-	-	-	-	86.0	75.6
PaLI-3-VPD (5B) <i>generalist</i>	61.3	57.5	78.5	-	56.5	-	83.1	70.9
PaLI-3-VPD (5B) <i>specialist</i>	64.7	60.3	79.7	76.5	65.5	63.6	83.3	70.8
PaLI-X-VPD (55B) <i>generalist</i>	64.9	64.6	84.5	-	62.7	-	86.2	76.6
PaLI-X-VPD (55B) <i>specialist</i>	67.3	66.8	85.2	80.4	71.1	68.2	86.2	76.4

Table 2. Comparison of per-task fine-tuning results of specialist models. PaLI-X-VPD sets a new SOTA for all the tasks.

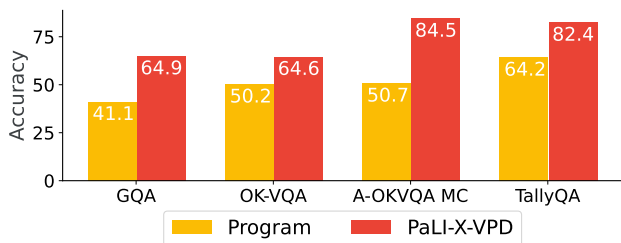


Figure 5. Accuracy of visual programs and PaLI-X-VPD on validation sets.

one program that passes the filtering stage, when the LLM generates the top-1 or top-5 programs, respectively. There is a dramatic increase in success rate from 1 program to 5: +45% on GQA and A-OKVQA, +33% on OK-VQA, and +10% on TallyQA. This design choice greatly improves our data synthesis efficiency, and, as a consequence, adding more CoT data that requires complex reasoning in our training set. We also conduct an analysis in §4.2 to compare the performance of VPD models with that of the visual programs they are distilled from.

Analysis: comparison with visual programs Figure 5 does a side-by-side comparison of the accuracy of visual programs and that of PaLI-X-VPD on GQA, OK-VQA, A-OKVQA (multiple choice), and TallyQA (simple and complex combined). We report results on the validation sets, so PaLI-X-VPD was not distilled with these exact visual programs, but with visual programs generated in a similar manner on the training set. The results indicate that PaLI-X-VPD has much higher accuracy than visual programs on all tasks. This raises an interesting question: why is the student model more accurate than its teacher? One explanation is that our pipeline allows us to leverage labeled data to improve the quality of the visual programs. When ground truth labels are available, we can choose a correct program among 5 candidates, rather than only relying on a single candidate. As supported by the results in Figure 4,

this greatly improves the accuracy of our visual programs as the teacher, thus making them more helpful for distillation.

4.3. Human evaluation on rationales

In this section, we focus on the quality of model outputs. We performed this analysis using human annotators, who are asked to evaluate both the correctness of the final answer and the quality of the rationale behind it.

Models. We compare PaLI-X-VPD with PaLI-X Instruct. Among possible baselines, we chose PaLI-X Instruct because it is trained to generate long-form answers [64], which allows us to assess if it is the quality of PaLI-X-VPD’s answers that the annotators prefer, rather than its length. Since PaLI-X Instruct is also instruction-tuned, it can be prompted to provide long answers to alleviate this confounder.

Annotation protocol. We run inference with each of the two models on a combination of 600 samples from GQA (test-dev), A-OKVQA (val), and TallyQA (Simple and Complex) and record their answers and rationales. We then ask 3 human annotators to evaluate each model answer. We use prior work from natural language processing (NLP) as inspiration, and build upon it for selecting the evaluation criteria [73]. Given an image and a query, for each model-generated rationale, we ask human annotators to score the model answers along the following criteria: (1) **correctness**—is the final answer correct? (2) **explainability**—does the model explain its rationale for reaching the final answer? (3) **factuality**—is every step in the rationale factually correct (with respect to the image and external knowledge)? (4) **consistency**—does step and final answer logically follow from the previous ones? Note that a rationale may have the wrong answer while being consistent. We also conduct a **side-by-side comparison**, and ask annotators which of the two answers—the one provided by PaLI-X-VPD or by PaLI-X Instruct—they prefer in terms of quality of the answer and explanation. More details about the annotation protocol are in §C.

Human evaluation results. The results of the evaluation are

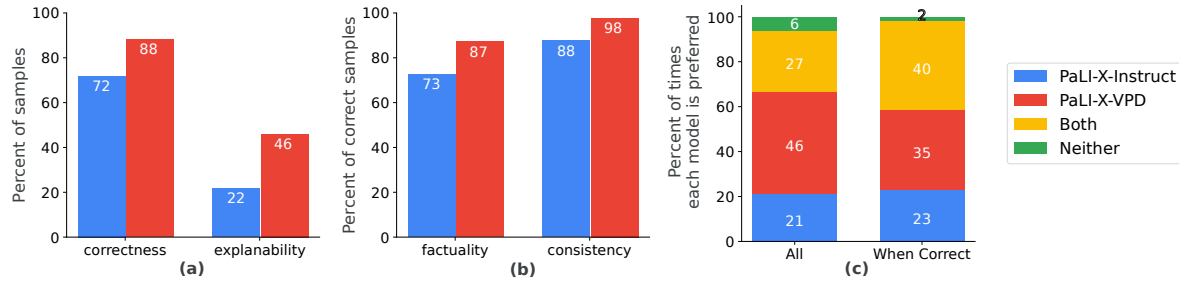


Figure 6. Human evaluation results assessing answer and rational quality for PaLI-X Instruct and PaLI-X VPD: (a) percentage of answers that are correct and have explanations; (b) rationale factuality and consistency for the answers that contain an explanation; (c) preference between the two models, when aggregating across all samples (“All”) and across those with correct answers (“When correct”).

522 first averaged among the human raters per sample, then aggregated across samples. **Our PaLI-X-VPD model far out-**
 523 **performs PaLI-X Instruct along all criteria.** Fig. 6
 524 (a) shows the correctness of the final answer provided by
 525 each model (i.e. accuracy), and the proportion of the sam-
 526 ples where the model provides a CoT rationale to explain
 527 its answer. PaLI-X-VPD’s gain in accuracy of +16.7%
 528 is even more impressive than in evaluations in §4.2 based
 529 on benchmark labels—this is because human annotators are
 530 better able to assess correctness for ambiguous questions
 531 with different possible interpretations (e.g., the model’s an-
 532 swer “*In the living room.*” to the question “*Where is the*
 533 *couch located?*” was considered correct by annotators, even
 534 when the benchmark answer was “*On the right.*”). More-
 535 over, the explainability results confirm that our model is
 536 able to explain its own answer on +24% more samples than
 537 the instruction tuned model. Additionally, Fig. 6 (b) shows
 538 **impressive rationale quality**: among the samples where
 539 an explanation is provided, PaLI-X-VPD’s rationales are
 540 **factual 87.2% of times**, and **consistent 97.8% of times**, a
 541 gain of +14.6% and +10%, respectively, when compared to
 542 PaLI-X Instruct.
 543

544 We also asked the annotators which of the two answers
 545 given by the two models they prefer. We aggregated these
 546 results in two ways: (1) across all samples, (2) across the
 547 samples where both models answered the question correctly.
 548 The results in Fig. 6 (c) show that PaLI-X-VPD is preferred
 549 on 25% more samples than PaLI-X Instruct in the “All”
 550 case, and on 12% more samples when both are correct. This
 551 suggests that even when PaLI-X-VPD makes mistakes, it
 552 still provides a better answer. Such examples are shown in
 553 §C. These results confirm that a model fine-tuned with VPD
 554 leads to more faithful and consistent answers.

555 5. Experiments on content moderation

556 Prior experiments focus on training generalist VLMs. Here,
 557 we explore the effectiveness of VPD at quickly adapting mod-
 558 els to real-world applications from a different domain than
 559 the training data. We experiment on Hateful Memes [29], a

560 content moderation dataset where the task is to classify if a
 561 meme contains hateful content. The target labels are “yes”
 562 or “no”, and models are evaluated in terms of classification
 563 accuracy and AUC-ROC. We experiment with two settings:
 564 *supervised*, in which the models are trained on the provided
 565 training set with 8, 500 labels, and *unsupervised*, in which no
 566 human labels are provided. We provide qualitative examples
 567 of the rationales our models give in Figure 7.

Model	Acc	AUC-ROC
<i>Unsupervised / Zero-Shot Methods</i>		
Flamingo (9B) [2]	57.0	-
InstructBLIP (Vicuna-13B) [14]	59.6	-
MiniGPT-V2 (7B) [6]	57.8	-
Generated Programs (ours)	69.7	70.1
PaLI-X-VPD (generalist)	61.4	66.8
PaLI-X-VPD (specialist w/ 0-shot CoT)	70.8	78.3
<i>Supervised Methods</i>		
VisualBERT [35]	69.5	75.4
Flamingo (80B) [2]	-	86.6
Previous SOTA [46]	78.8	86.7
PaLI-X-VPD (label-only FT)	77.6	88.0
PaLI-X-VPD (specialist w/ CoT)	80.8	89.2
Human [29]	84.7	82.7

Table 3. Results on Hateful Memes [29] seen test set. We improve SOTA for both supervised and unsupervised settings. Surprisingly, unsupervised PaLI-X-VPD outperforms supervised VisualBERT.

568 5.1. Supervised setting

569 We first experiment with the supervised setting. We fine-tune
 570 the models on 8,500 labels in two ways:

- 571 1. **Label-only fine-tuning**: To establish a strong baseline,
 572 we first experiment with the traditional supervised setting,
 573 where we fine-tune the model to output “yes” or “no”.
 574
- 575 2. **PaLI-X-VPD (specialist)**: We use the complete VPD
 576 pipeline to train this model. We select an execution trace
 577 that leads to the correct label, and use it to tune our VLM.

577 **Results.** Our PaLI-X-VPD (specialist) **sets a new**
 578 **SOTA for this task**, with an accuracy of 80.8% and
 579 AUC-ROC of 89.2%. It outperforms the label-only fine-
 580 tuning baseline and significantly improves the SOTA met-

581 rics, **achieving nearly human-level accuracy, and super-**
582 **human AUC-ROC.**

583 5.2. Unsupervised setting

584 How would VPD perform if we do not have any human-
585 annotated labels? We experiment with three methods:

- 586 1. Generated Program: The program generated by PaLM-
587 2 [3] for solving this task consists of following main steps:
588 (1) get image description with PaLI-X [10]; (2) use an
589 OCR tool to extract embedded texts; (3) given the image
590 description and OCR texts, ask PaLM-2 to explain if this
591 meme is hateful.
- 592 2. PaLI-X-VPD (generalist): In a zero-shot setting, we di-
593 rectly prompt our PaLI-X-VPD (generalist) with:
594 “The text is <OCR text>. Is this a hateful meme?” We com-
595 pute the probability of PaLI-X-VPD generating “yes”
596 or “no” and measure accuracy and AUC-ROC.
- 597 3. PaLI-X-VPD (specialist with zero-shot CoT): We follow
598 our VPD pipeline and convert the execution traces of the
599 program into CoTs, and then fine-tune PaLI-X-VPD
600 to output these CoTs. Since no groundtruth labels are
601 available, no filtering is done during the process.

602 **VPD significantly improves performance even when no**
603 **labels are available.** As shown in Table 3, PaLI-X-VPD
604 (generalist) outperforms all other VLMs (61.4%). In-
605 terestingly, the generated programs themselves get much
606 higher accuracy (69.7%) than on the previous datasets,
607 perhaps because PaLM-2 is better suited at analyzing a
608 meme than typical VQA datasets. As shown in Figure 7,
609 PaLI-X-VPD (generalist) is relatively insensitive
610 to hateful content, while programs are better. Moreover,
611 our PaLI-X-VPD (specialist) sets a new zero-shot
612 SOTA on Hateful Memes, and even outperforms supervised
613 VisualBERT [35]. To understand this impressive gain, we
614 manually inspect the model outputs and find that the model
615 has learnt PaLM-2’s reasoning process via VPD, as seen in
616 Figure 7. Besides, our model has better understanding of
617 the image. As seen in Figure 7 (a), the program, despite
618 giving the correct answer, has not mentioned anything about
619 “black people”, because it is not covered in the image cap-
620 tion. In comparison, our PaLI-X-VPD (specialist)
621 includes it and gives a better explanation.

622 **More qualitative analysis.** We observe that PaLI-X-VPD
623 (supervised specialist) is able to capture the nuances in the
624 meme, while unsupervised methods fail, as shown in Figure 7
625 (c). Also, there are some failure cases like (d), in which even
626 our supervised model fails while the program succeeds on
627 hateful content detection.

628 6. Limitations and Future Directions

629 We mainly identify two directions for improvement. First,
630 VPD can be further scaled-up with more diverse tasks. Sec-

ond, we find that there are some problems that our visual
program framework (ViperGPT [16]) cannot solve. Further
improving the programs will yield bigger gains for VPD. We
list the limitations below, along with future work that may
address these challenges.

**Scaling-up VPD with LLM-generated questions and bet-
ter filtering strategies.** Our current setting limits the data
source of VPD to existing VQA tasks. Future work can scale-
up and diversify the tasks using LLMs, as demonstrated in
Self-Instruct [63]. One challenge is that there lacks human-
annotated labels for LLM-generated tasks. Experiments in
§5 show that VPD works well even without labels. Neverthe-
less, it would be better if we have some filtering strategies,
as shown in §4.2. Since even the strongest VLMs like PaLI-
X [10] and GPT-4V [49] cannot give reliable answers, fact-
checking strategies for multimodal chain-of-thought data is
a promising future direction.

Agents, rather than static programs. There exist complex
visual-language tasks that cannot be easily solved with one
program. Recent work [66, 70, 71] have explored the idea
of LLM agents, where LLMs interact with an environment
and do planning interactively. We may be able to leverage
this idea in our scenario, and have an LLM update its plans
given the new information gathered from vision tools.

Adding fine-grained and dense labeling tools. We find that
programs struggles in complex scenarios when objects are
occluded, or the scene contains too many objects. Adding
dense-labeling tools like Segment Anything [31] may ad-
dress this challenge. Other recent work (e.g., LISA [33])
have proposed combining segmentation with LLMs. Future
work can make dense labeling tools available in VPD in a
similar way, which will further boost VLM performance.

7. Conclusion

We introduce VPD, a framework for distilling the reasoning
abilities of LLMs along with the capabilities of vision tools
into VLMs. VPD synthesizes training data for VLMs by
generating programs that can leverage external tools. We use
this technique to fine-tune some of the best existing VLMs
(PaLI-3 and PaLI-X) and established new SOTA results on
8 classical VQA and 2 zero-shot multimodal benchmarks.
According to human evaluations, VPD-tuned models provide
more accurate answers and better rationales. Experiments on
the Hateful Memes show how VPD can also adapt models
to new real-world domains, even when no labeled data is
available, also establishing new SOTAs. We also point out
directions to further improve VPD. We believe VPD will
grant future VLMs better multimodal reasoning abilities.



Figure 7. Example outputs of different methods on Hateful Memes [29] dev set. Here we prompt our PaLI-X-VPD models with “Is this a hateful meme? Explain the rationale to answer the question.” to get models’ long-form rationales. The unsupervised methods include zero-shot PaLI-X-VPD (generalist), our generated program, and PaLI-X-VPD (specialist with zero-shot CoTs). We also include our supervised method, i.e., PaLI-X-VPD (specialist). We mark whether their outputs match the gold answers.

678

References

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

714

715

716

717

718

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

- [1] Manoj Acharya, Kushal Kafle, and Christopher Kanan. Tal-lyqa: Answering complex counting questions. In *Proceedings of the AAAI conference on artificial intelligence*, pages 8076–8084, 2019. 1, 5, 16, 17
- [2] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katie Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob Menick, Sebastian Borgeaud, Andrew Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikolaj Binkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karen Simonyan. Flamingo: a visual language model for few-shot learning, 2022. 1, 3, 6, 8
- [3] Rohan Anil, Andrew M. Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos, Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, Eric Chu, Jonathan H. Clark, Laurent El Shafey, Yanping Huang, Kathy Meier-Hellstern, Gaurav Mishra, Erica Moreira, Mark Omernick, Kevin Robinson, Sebastian Ruder, Yi Tay, Ke-fan Xiao, Yuanzhong Xu, Yujing Zhang, Gustavo Hernandez Abrego, Junwhan Ahn, Jacob Austin, Paul Barham, Jan Botha, James Bradbury, Siddhartha Brahma, Kevin Brooks, Michele Catasta, Yong Cheng, Colin Cherry, Christopher A. Choquette-Choo, Aakanksha Chowdhery, Clément Crepy, Shachi Dave, Mostafa Dehghani, Sunipa Dev, Jacob Devlin, Mark Díaz, Nan Du, Ethan Dyer, Vlad Feinberg, Fangxiaoyu Feng, Vlad Fienber, Markus Freitag, Xavier Garcia, Sebastian Gehrmann, Lucas Gonzalez, Guy Gur-Ari, Steven Hand, Hadi Hashemi, Le Hou, Joshua Howland, Andrea Hu, Jeffrey Hui, Jeremy Hurwitz, Michael Isard, Abe Ittycheriah, Matthew Jagielski, Wenhao Jia, Kathleen Kenealy, Maxim Krikun, Sneha Kudugunta, Chang Lan, Katherine Lee, Benjamin Lee, Eric Li, Music Li, Wei Li, YaGuang Li, Jian Li, Hyeontaek Lim, Hanzhao Lin, Zhongtao Liu, Frederick Liu, Marcello Maggioni, Aroma Mahendru, Joshua Maynez, Vedant Misra, Maysam Moussalem, Zachary Nado, John Nham, Eric Ni, Andrew Nystrom, Alicia Parrish, Marie Pellat, Martin Polacek, Alex Polozov, Reiner Pope, Siyuan Qiao, Emily Reif, Bryan Richter, Parker Riley, Alex Castro Ros, Aurko Roy, Brennan Saeta, Rajkumar Samuel, Renee Shelby, Ambrose Slone, Daniel Smilkov, David R. So, Daniel Sohn, Simon Tokumine, Dasha Valter, Vijay Vasudevan, Kiran Vodrahalli, Xuezhi Wang, Pidong Wang, Zirui Wang, Tao Wang, John Wieting, Yuhuai Wu, Kelvin Xu, Yunhan Xu, Linting Xue, Pengcheng Yin, Jiahui Yu, Qiao Zhang, Steven Zheng, Ce Zheng, Weikang Zhou, Denny Zhou, Slav Petrov, and Yonghui Wu. Palm 2 technical report, 2023. 1, 3, 4, 9, 15
- [4] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond, 2023. 1, 3, 6
- [5] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information*

- processing systems*, 33:1877–1901, 2020. 3 735
- [6] Jun Chen, Deyao Zhu, Xiaoqian Shen, Xiang Li, Zechun Liu, Pengchuan Zhang, Raghuraman Krishnamoorthi, Vikas Chandra, Yunyang Xiong, and Mohamed Elhoseiny. Minigt-v2: large language model as a unified interface for vision-language multi-task learning, 2023. 3, 8 736
737
738
739
740
- [7] Keqin Chen, Zhao Zhang, Weili Zeng, Richong Zhang, Feng Zhu, and Rui Zhao. Shikra: Unleashing multimodal llm’s referential dialogue magic, 2023. 1, 2, 3, 6 741
742
743
- [8] Wenhui Chen, Xueguang Ma, Xinyi Wang, and William W. Cohen. Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks, 2022. 1, 3 744
745
746
747
- [9] Xi Chen, Xiao Wang, Soravit Changpinyo, AJ Piergiovanni, Piotr Padlewski, Daniel Salz, Sebastian Goodman, Adam Grycner, Basil Mustafa, Lucas Beyer, Alexander Kolesnikov, Joan Puigcerver, Nan Ding, Keran Rong, Hassan Akbari, Gaurav Mishra, Linting Xue, Ashish Thapliyal, James Bradbury, Weicheng Kuo, Mojtaba Seyedhosseini, Chao Jia, Burcu Karagol Ayan, Carlos Riquelme, Andreas Steiner, Anelia Angelova, Xiaohua Zhai, Neil Houlsby, and Radu Soricut. Pali: A jointly-scaled multilingual language-image model, 2022. 1, 3, 7 748
749
750
751
752
753
754
755
756
757
- [10] Xi Chen, Josip Djolonga, Piotr Padlewski, Basil Mustafa, Soravit Changpinyo, Jialin Wu, Carlos Riquelme Ruiz, Sebastian Goodman, Xiao Wang, Yi Tay, Siamak Shakeri, Mostafa Dehghani, Daniel Salz, Mario Lucic, Michael Tschannen, Arsha Nagrani, Hexiang Hu, Mandar Joshi, Bo Pang, Ceslee Montgomery, Paulina Pietrzyk, Marvin Ritter, AJ Piergiovanni, Matthias Minderer, Filip Pavetic, Austin Waters, Gang Li, Ibrahim Alabdulmohsin, Lucas Beyer, Julien Amelot, Kenton Lee, Andreas Peter Steiner, Yang Li, Daniel Keysers, Anurag Arnab, Yuanzhong Xu, Keran Rong, Alexander Kolesnikov, Mojtaba Seyedhosseini, Anelia Angelova, Xiaohua Zhai, Neil Houlsby, and Radu Soricut. Pali-x: On scaling up a multilingual vision and language model, 2023. 2, 3, 4, 5, 7, 9, 17 758
759
760
761
762
763
764
765
766
767
768
769
770
771
- [11] Xi Chen, Xiao Wang, Lucas Beyer, Alexander Kolesnikov, Jialin Wu, Paul Voigtlaender, Basil Mustafa, Sebastian Goodman, Ibrahim Alabdulmohsin, Piotr Padlewski, Daniel Salz, Xi Xiong, Daniel Vlasic, Filip Pavetic, Keran Rong, Tianli Yu, Daniel Keysers, Xiaohua Zhai, and Radu Soricut. Pali-3 vision language models: Smaller, faster, stronger, 2023. 1, 4, 5, 7, 17 772
773
774
775
776
777
778
- [12] Zhoujun Cheng, Tianbao Xie, Peng Shi, Chengzu Li, Rahul Nadkarni, Yushi Hu, Caiming Xiong, Dragomir Radev, Mari Ostendorf, Luke Zettlemoyer, et al. Binding language models in symbolic languages. *arXiv preprint arXiv:2210.02875*, 2022. 3 779
780
781
782
783
- [13] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk 784
785
786
787
788
789
790
791
792

- 793 Michalewski, Xavier Garcia, Vedant Misra, Kevin Robin-
794 son, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan,
795 Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan
796 Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, An-
797 drew M. Dai, Thanumalayan Sankaranarayana Pillai, Marie
798 Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Olek-
799 sandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang,
800 Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta,
801 Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean,
802 Slav Petrov, and Noah Fiedel. Palm: Scaling language model-
803 ing with pathways, 2022. 3
- [14] Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat
804 Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung,
805 and Steven Hoi. Instructblip: Towards general-purpose vision-
806 language models with instruction tuning, 2023. 2, 6, 7, 8
- [15] Mostafa Dehghani, Josip Djolonga, Basil Mustafa, Piotr
807 Padlewski, Jonathan Heek, Justin Gilmer, Andreas Steiner,
808 Mathilde Caron, Robert Geirhos, Ibrahim Alabdulmohsin,
809 Rodolphe Jenatton, Lucas Beyer, Michael Tschannen, Anurag
810 Arnab, Xiao Wang, Carlos Riquelme, Matthias Minderer,
811 Joan Puigcerver, Utku Evci, Manoj Kumar, Sjoerd van
812 Steenkiste, Gamaleldin F. Elsayed, Aravindh Mahendran,
813 Fisher Yu, Avital Oliver, Fantine Huot, Jasmijn Bastings,
814 Mark Patrick Collier, Alexey Gritsenko, Vighnesh Birodkar,
815 Cristina Vasconcelos, Yi Tay, Thomas Mensink, Alexander
816 Kolesnikov, Filip Pavetić, Dustin Tran, Thomas Kipf, Mario
817 Lučić, Xiaohua Zhai, Daniel Keysers, Jeremiah Harmsen,
818 and Neil Houlsby. Scaling vision transformers to 22 billion
819 parameters, 2023. 17
- [16] Surís Dídac, Sachit Menon, and Carl Vondrick. Vipergpt:
820 Visual inference via python execution for reasoning. *arXiv*
821 *preprint arXiv:2303.08128*, 2023. 1, 2, 3, 9, 20
- [17] Ruofei Du, Eric Turner, Maksym Dzitsiuk, Luca Prasso, Ivo
822 Duarte, Jason Dourgarian, Joao Afonso, Jose Pascoal, Josh
823 Gladstone, Nuno Cruces, Shahram Izadi, Adarsh Kowdle,
824 Konstantine Tsotsos, and David Kim. DepthLab: Real-Time
825 3D Interaction With Depth Maps for Mobile Augmented Re-
826 ality. In *Proceedings of the 33rd Annual ACM Symposium on*
827 *User Interface Software and Technology*. ACM, 2020. 4
- [18] Weixi Feng, Wanrong Zhu, Tsu-Jui Fu, Varun Jampani, Ar-
828 jun Reddy Akula, Xuehai He, Sugato Basu, Xin Eric Wang,
829 and William Yang Wang. Layoutgpt: Compositional visual
830 planning and generation with large language models. *ArXiv*,
831 *abs/2305.15393*, 2023. 2
- [19] Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Ba-
832 tra, and Devi Parikh. Making the V in VQA matter: Elev-
833 ating the role of image understanding in Visual Question
834 Answering. In *Conference on Computer Vision and Pattern*
835 *Recognition (CVPR)*, 2017. 5, 17
- [20] Tanmay Gupta and Aniruddha Kembhavi. Visual program-
836 ming: Compositional visual reasoning without training. *2023*
837 *IEEE/CVF Conference on Computer Vision and Pattern*
838 *Recognition (CVPR)*, 2023. 1, 2, 3
- [21] Cheng-Yu Hsieh, Chun-Liang Li, Chih-Kuan Yeh, Hootan
839 Nakhost, Yasuhisa Fujii, Alexander Ratner, Ranjay Krishna,
840 Chen-Yu Lee, and Tomas Pfister. Distilling step-by-step!
841 outperforming larger language models with less training data
842 and smaller model sizes, 2023. 2, 3, 4
- [22] Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-
843 Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen.
844 Lora: Low-rank adaptation of large language models, 2021.
845 5, 17
- [23] Yushi* Hu, Hang* Hua, Zhengyuan Yang, Weijia Shi, Noah A
846 Smith, and Jiebo Luo. Promptcap: Prompt-guided task-aware
847 image captioning. *arXiv preprint arXiv:2211.09699*, 2022. 2,
848 3, 4
- [24] Yushi Hu, Chia-Hsuan Lee, Tianbao Xie, Tao Yu, Noah A.
849 Smith, and Mari Ostendorf. In-context learning for few-shot
850 dialogue state tracking. In *Findings of the Association for*
851 *Computational Linguistics: EMNLP 2022*, pages 2627–2643,
852 Abu Dhabi, United Arab Emirates, 2022. Association for
853 Computational Linguistics. 4
- [25] Ziniu Hu, Ahmet Iscen, Chen Sun, Kai-Wei Chang, Yizhou
854 Sun, David A Ross, Cordelia Schmid, and Alireza Fathi. Avis:
855 Autonomous visual information seeking with large language
856 model agent, 2023. 1, 3
- [26] Drew A Hudson and Christopher D Manning. Gqa: A new
857 dataset for real-world visual reasoning and compositional
858 question answering. In *Proceedings of the IEEE/CVF con-*
859 *ference on computer vision and pattern recognition*, pages
860 6700–6709, 2019. 1, 5, 6, 16, 17
- [27] Norman P. Jouppi, Cliff Young, Nishant Patil, David Patter-
861 son, Gaurav Agrawal, Raminder Bajwa, Sarah Bates, Suresh
862 Bhatia, Nan Boden, Al Borchers, Rick Boyle, Pierre-luc
863 Cantin, Clifford Chao, Chris Clark, Jeremy Coriell, Mike
864 Daley, Matt Dau, Jeffrey Dean, Ben Gelb, Tara Vazir Ghaem-
865 maghami, Rajendra Gottipati, William Gulland, Robert Hag-
866 mann, C. Richard Ho, Doug Hogberg, John Hu, Robert
867 Hundt, Dan Hurt, Julian Ibarz, Aaron Jaffey, Alek Ja-
868 worski, Alexander Kaplan, Harshit Khaitan, Daniel Kille-
869 brew, Andy Koch, Naveen Kumar, Steve Lacy, James Laudon,
870 James Law, Diemthu Le, Chris Leary, Zhuyuan Liu, Kyle
871 Lucke, Alan Lundin, Gordon MacKean, Adriana Mag-
872 giore, Maire Mahony, Kieran Miller, Rahul Nagarajan, Ravi
873 Narayanaswami, Ray Ni, Kathy Nix, Thomas Norrie, Mark
874 Omernick, Narayana Penukonda, Andy Phelps, Jonathan
875 Ross, Matt Ross, Amir Salek, Emad Samadiani, Chris Sev-
876 ern, Gregory Sizikov, Matthew Snelman, Jed Souter, Dan
877 Steinberg, Andy Swing, Mercedes Tan, Gregory Thorson, Bo
878 Tian, Horia Toma, Erick Tuttle, Vijay Vasudevan, Richard
879 Walter, Walter Wang, Eric Wilcox, and Doe Hyun Yoon. In-
880 datacenter performance analysis of a tensor processing unit.
881 In *Proceedings of the 44th Annual International Symposium*
882 *on Computer Architecture*. ACM, 2017. 18
- [28] Ehsan Kamaloo, Nouha Dziri, Charles L. A. Clarke, and
883 Davood Rafiei. Evaluating open-domain question answering
884 in the era of large language models, 2023. 4, 23
- [29] Douwe Kiela, Hamed Firooz, Aravind Mohan, Vedanuj
885 Goswami, Amanpreet Singh, Pratik Ringshia, and Davide
886 Testuggine. The Hateful Memes Challenge: Detecting Hate
887 Speech in Multimodal Memes, 2020. 2, 8, 10
- [30] Diederik P. Kingma and Jimmy Ba. Adam: A method for
888 stochastic optimization. *CoRR*, *abs/1412.6980*, 2015. 18
- [31] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao,
889 Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer White-
890 891 892 893 894 895 896 897 898 899 900 901 902 903 904 905 906 907

- 908 head, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross
909 Girshick. Segment anything, 2023. 9
- 910 [32] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson,
911 Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalan-
912 tidis, Li-Jia Li, David A. Shamma, Michael S. Bernstein, and
913 Fei-Fei Li. Visual genome: Connecting language and vision
914 using crowdsourced dense image annotations, 2016. 3, 5, 17
- 915 [33] Xin Lai, Zhuotao Tian, Yukang Chen, Yanwei Li, Yuhui Yuan,
916 Shu Liu, and Jiaya Jia. Lisa: Reasoning segmentation via
917 large language model, 2023. 9
- 918 [34] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi.
919 Blip-2: Bootstrapping language-image pre-training with
920 frozen image encoders and large language models. *ArXiv*,
921 abs/2301.12597, 2023. 1, 3
- 922 [35] Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh,
923 and Kai-Wei Chang. VisualBERT: A Simple and Performant
924 Baseline for Vision and Language. *ArXiv*, abs/1908.03557,
925 2019. 8, 9
- 926 [36] Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin
927 Zhao, and Ji-Rong Wen. Evaluating object hallucination in
928 large vision-language models, 2023. 5, 17
- 929 [37] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays,
930 Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence
931 Zitnick. Microsoft coco: Common objects in context. In
932 *European conference on computer vision*, pages 740–755.
933 Springer, 2014. 3
- 934 [38] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee.
935 Improved baselines with visual instruction tuning, 2023. 2, 6
- 936 [39] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee.
937 Visual instruction tuning. *arXiv preprint arXiv:2304.08485*,
938 2023. 1, 2, 3, 4
- 939 [40] Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang
940 Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He,
941 Ziwei Liu, Kai Chen, and Dahua Lin. Mmbench: Is your
942 multi-modal model an all-around player?, 2023. 2, 5, 17
- 943 [41] Jiasen Lu, Christopher Clark, Rowan Zellers, Roozbeh Mot-
944 taghi, and Aniruddha Kembhavi. Unified-io: A unified
945 model for vision, language, and multi-modal tasks. *ArXiv*,
946 abs/2206.08916, 2022. 1, 3
- 947 [42] Pan Lu, Baolin Peng, Hao Cheng, Michel Galley, Kai-Wei
948 Chang, Ying Nian Wu, Song-Chun Zhu, and Jianfeng Gao.
949 Chameleon: Plug-and-play compositional reasoning with
950 large language models, 2023. 3
- 951 [43] Qing Lyu, Shreya Havaldar, Adam Stein, Li Zhang, Delip
952 Rao, Eric Wong, Marianna Apidianaki, and Chris Callison-
953 Burch. Faithful chain-of-thought reasoning, 2023. 1, 3
- 954 [44] Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and
955 Roozbeh Mottaghi. Ok-vqa: A visual question answering
956 benchmark requiring external knowledge. In *Proceedings*
957 *of the IEEE/cvf conference on computer vision and pattern*
958 *recognition*, pages 3195–3204, 2019. 1
- 959 [45] Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and
960 Roozbeh Mottaghi. Ok-vqa: A visual question answering
961 benchmark requiring external knowledge. *2019 IEEE/CVF*
962 *Conference on Computer Vision and Pattern Recognition*
963 *(CVPR)*, 2019. 5, 6, 17
- [46] Jingbiao Mei, Jinghong Chen, Weizhe Lin, Bill Byrne,
and Marcus Tomalin. Improving hateful memes detection
via learning hatefulness-aware embedding space through
retrieval-guided contrastive learning, 2023. 8
- [47] Matthias Minderer, Alexey Gritsenko, and Neil Houlsby. Scal-
ing open-vocabulary object detection, 2023. 4
- [48] Anand Mishra, Shashank Shekhar, Ajeet Kumar Singh, and
Anirban Chakraborty. Ocr-vqa: Visual question answering by
reading text in images. In *ICDAR*, 2019. 5, 17
- [49] OpenAI. Gpt-4 technical report, 2023. 1, 3, 9
- [50] Zhiliang Peng, Wenhui Wang, Li Dong, Yaru Hao, Shaohan
Huang, Shuming Ma, and Furu Wei. Kosmos-2: Grounding
multimodal large language models to the world, 2023. 3
- [51] Timo Schick, Jane Dwivedi-Yu, Roberto Dessi, Roberta
Raileanu, Maria Lomeli, Luke Zettlemoyer, Nicola Cancedda,
and Thomas Scialom. Toolformer: Language models can
teach themselves to use tools, 2023. 3
- [52] Dustin Schwenk, Apoorv Khandelwal, Christopher Clark,
Kenneth Marino, and Roozbeh Mottaghi. A-okvqa: A bench-
mark for visual question answering using world knowledge,
2022. 5, 16, 17
- [53] Dustin Schwenk, Apoorv Khandelwal, Christopher Clark,
Kenneth Marino, and Roozbeh Mottaghi. A-okvqa: A bench-
mark for visual question answering using world knowledge.
arXiv preprint arXiv:2206.01718, 2022. 1
- [54] Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, Weim-
ing Lu, and Yueting Zhuang. Hugginggpt: Solving ai tasks
with chatgpt and its friends in hugging face, 2023. 3
- [55] Amanpreet Singh, Vivek Natarjan, Meet Shah, Yu Jiang, Xin-
lei Chen, Devi Parikh, and Marcus Rohrbach. Towards vqa
models that can read. In *Proceedings of the IEEE Confer-
ence on Computer Vision and Pattern Recognition*, pages
8317–8326, 2019. 5, 17
- [56] Sanjay Subramanian, Medhini Narasimhan, Kushal
Khangonkar, Kevin Yang, Arsha Nagrani, Cordelia Schmid,
Andy Zeng, Trevor Darrell, and Dan Klein. Modular visual
question answering via code generation, 2023. 1
- [57] Yi Tay, Mostafa Dehghani, Vinh Q. Tran, Xavier Garcia, Ja-
son Wei, Xuezhi Wang, Hyung Won Chung, Siamak Shakeri,
Dara Bahri, Tal Schuster, Huaixiu Steven Zheng, Denny Zhou,
Neil Houlsby, and Donald Metzler. UL2: Unifying Language
Learning Paradigms, 2022. 17
- [58] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Am-
jad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya
Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas
Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cu-
curull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin
Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman
Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan
Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Is-
abel Kloumann, Artem Korenev, Punit Singh Koura, Marie-
Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich,
Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov,
Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton,
Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schel-
ten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian,
Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, 964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020

1021	Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models, 2023. 3		
1022		[71] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. <i>arXiv preprint arXiv:2210.03629</i> , 2022. 9	1079
1023			1080
1024			1081
1025			1082
1026	[59] Jianfeng Wang, Zhengyuan Yang, Xiaowei Hu, Linjie Li, Kevin Lin, Zhe Gan, Zicheng Liu, Ce Liu, and Lijuan Wang. Git: A generative image-to-text transformer for vision and language. <i>arXiv preprint arXiv:2205.14100</i> , 2022. 1, 3	[72] Qinghao Ye, Haiyang Xu, Jiabo Ye, Ming Yan, Haowei Liu, Qi Qian, Ji Zhang, Fei Huang, and Jingren Zhou. mplug-owl2: Revolutionizing multi-modal large language model with modality collaboration, 2023. 1, 3, 6	1083
1027			1084
1028			1085
1029			1086
1030	[60] Peng Wang, An Yang, Rui Men, Junyang Lin, Shuai Bai, Zhikang Li, Jianxin Ma, Chang Zhou, Jingren Zhou, and Hongxia Yang. Ofa: Unifying architectures, tasks, and modalities through a simple sequence-to-sequence learning framework. <i>CoRR</i> , abs/2202.03052, 2022. 3	[73] Xi Ye and Greg Durrett. The unreliability of explanations in few-shot prompting for textual reasoning. In <i>Advances in Neural Information Processing Systems</i> , pages 30378–30392. Curran Associates, Inc., 2022. 7	1087
1031			1088
1032			1089
1033			1090
1034			1091
1035	[61] Peifeng Wang, Zhengyang Wang, Zheng Li, Yifan Gao, Bing Yin, and Xiang Ren. Scott: Self-consistent chain-of-thought distillation, 2023. 3	[74] Xiang Yue, Xingwei Qu, Ge Zhang, Yao Fu, Wenhao Huang, Huan Sun, Yu Su, and Wenhui Chen. Mammoth: Building math generalist models through hybrid instruction tuning. <i>arXiv preprint arXiv:2309.05653</i> , 2023. 3	1092
1036			1093
1037			1094
1038	[62] Weihang Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang, Lei Zhao, Xixuan Song, Jiazheng Xu, Bin Xu, Juanzi Li, Yuxiao Dong, Ming Ding, and Jie Tang. Cogvlm: Visual expert for pretrained language models, 2023. 1, 3, 7	[75] Andy Zeng, Adrian Wong, Stefan Welker, Krzysztof Choromanski, Federico Tombari, Aavek Purohit, Michael Ryoo, Vikas Sindhwani, Johnny Lee, Vincent Vanhoucke, et al. Socratic models: Composing zero-shot multimodal reasoning with language. <i>arXiv preprint arXiv:2204.00598</i> , 2022. 3	1095
1039			1096
1040			1097
1041			1098
1042			1099
1043	[63] Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and Hannaneh Hajishirzi. Self-instruct: Aligning language models with self-generated instructions, 2022. 5, 9	[76] Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training, 2023. 17	1100
1044			1101
1045			1102
1046			1103
1047	[64] Yaqing Wang, Jialin Wu, Tanmaya Dabral, Jiageng Zhang, Geoff Brown, Chun-Ta Lu, Frederick Liu, Yi Liang, Bo Pang, Michael Bendersky, and Radu Soricut. Non-intrusive adaptation: Input-centric parameter-efficient fine-tuning for versatile multimodal modeling, 2023. 2, 5, 7	[77] Haotian Zhang, Pengchuan Zhang, Xiaowei Hu, Yen-Chun Chen, Liunian Harold Li, Xiyang Dai, Lijuan Wang, Lu Yuan, Jenq-Neng Hwang, and Jianfeng Gao. Glipv2: Unifying localization and vision-language understanding, 2022. 3	1104
1048			1105
1049			1106
1050			1107
1051			1108
1052	[65] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models, 2022. 3	[78] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: Enhancing vision-language understanding with advanced large language models, 2023. 1, 3, 6	1109
1053			1110
1054			1111
1055			1112
1056	[66] Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun Zhang, Shaokun Zhang, Jiale Liu, Ahmed Hassan Awadallah, Ryen W White, Doug Burger, and Chi Wang. Autogen: Enabling next-gen llm applications via multi-agent conversation. 2023. 9	[79] Xuekai Zhu, Biqing Qi, Kaiyan Zhang, Xingwei Long, and Bowen Zhou. Pad: Program-aided distillation specializes large models in reasoning, 2023. 3	1113
1057			
1058			
1059			
1060			
1061	[67] Yiheng Xu, Hongjin Su, Chen Xing, Boyu Mi, Qian Liu, Weijia Shi, Binyuan Hui, Fan Zhou, Yitao Liu, Tianbao Xie, Zhoujun Cheng, Siheng Zhao, Lingpeng Kong, Bailin Wang, Caiming Xiong, and Tao Yu. Lemur: Harmonizing natural language and code for language agents, 2023. 1		
1062			
1063			
1064			
1065			
1066	[68] Zhengyuan Yang, Zhe Gan, Jianfeng Wang, Xiaowei Hu, Yumao Lu, Zicheng Liu, and Lijuan Wang. An empirical study of gpt-3 for few-shot knowledge-based vqa. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , pages 3081–3089, 2022. 3		
1067			
1068			
1069			
1070			
1071	[69] Zhengyuan Yang, Linjie Li, Kevin Lin, Jianfeng Wang, Chung-Ching Lin, Zicheng Liu, and Lijuan Wang. The dawn of llms: Preliminary explorations with gpt-4v (ision). <i>arXiv preprint arXiv:2309.17421</i> , 9, 2023. 1		
1072			
1073			
1074			
1075	[70] Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Ehsan Azarnasab, Faisal Ahmed, Zicheng Liu, Ce Liu, Michael Zeng, and Lijuan Wang. Mm-react: Prompting chatgpt for multimodal reasoning and action, 2023. 3, 9		
1076			
1077			
1078			

A. Examples of our Data-Synthesis Pipeline

1114

We first present step-by-step examples of our data-synthesis pipeline. As discussed in §3, the pipeline contains four steps: program generation, program execution, program filtering, and converting execution traces into chains-of-thought. Here we show three filtered examples from three datasets we use. As shown, LLMs like PaLM-2 [3] is able to generate complex programs to solve compositional vision tasks. The synthesized CoTs are used to fine-tune our VLMs.

1115

1116

1117

1118

(a)



Is the bookshelf to the right or to the left of the chair that is to the left of the vase?

PaLM-2 Generated Program

```

1 def execute_command(image):
2     image_patch = ImagePatch(image)
3
4     # Find the chair that is to the left of the vase
5     chair_patches = image_patch.find("chair")
6     vase_patches = image_patch.find("vase")
7
8     chair_patch_on_the_left_of_the_vase = None
9
10    for vase_patch in vase_patches:
11
12        # find all the chairs on the left of the vase.
13        chair_patches_on_the_left = []
14        for chair_patch in chair_patches:
15            if chair_patch.horizontal_center < vase_patch.horizontal_center:
16                chair_patch_on_the_left = chair_patch
17                chair_patches_on_the_left.append(chair_patch_on_the_left)
18
19        # the chair on the left of the vase should be the one that is closest on the vase.
20        if len(chair_patches_on_the_left) > 0:
21            chair_patch_on_the_left_of_the_vase = sorted(chair_patches_on_the_left,
22                                                         key=lambda patch: distance(vase_patch, chair_patch))[0]
23            break
24
25    print(f"the chair at {chair_patch_on_the_left} is to the left of the vase at {vase_patch}.")
26
27    # Find the bookshelf
28    bookshelf_patch = image_patch.find("bookshelf")[0]
29
30    # Check if the bookshelf is to the right or to the left of the chair
31    if bookshelf_patch.horizontal_center > chair_patch_on_the_left.horizontal_center:
32        print(f"the bookshelf at {bookshelf_patch} is to the right of the chair at {chair_patch}.")
33        return "right"
34    elif bookshelf_patch.horizontal_center < chair_patch_on_the_left.horizontal_center:
35        print(f"the bookshelf at {bookshelf_patch} is to the left of the chair at {chair_patch}.")
36        return "left"

```

Execution Trace

```

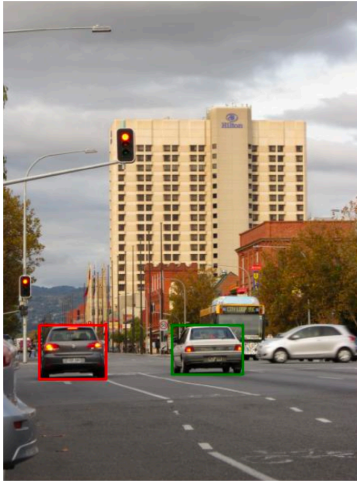
1 Calling find function. Detect chair
2 Detection result: 599 64 655 107 chair and 624 143 836 245 chair and 586 321 782 395 chair and 603 467 771 549 chair
3 Calling find function. Detect vase
4 Detection result: 761 0 889 70 vase and 676 615 756 653 vase
5 the chair at 603 467 771 549 is to the left of the vase at 676 615 756 653.
6 Calling find function. Detect bookshelf
7 Detection result: 505 244 714 359 bookshelf
8 the bookshelf at 505 244 714 359 is to the left of the chair at 603 467 771 549.
9 Program output: left

```

Conversion to CoT

The vase is at **676 615 756 653**. To its left, the chair is at **603 467 771 549**. The bookshelf is at **505 244 714 359**. Therefore, the bookshelf is to the left of the chair.

(b)



How many cars have the brake lights on?

PaLM-2 Generated Program

```

1 def execute_command(image):
2     image_patch = ImagePatch(image)
3
4     # Find all cars. Notice that find() only support single word queries.
5     car_patches = image_patch.find("car")
6
7     # Count the number of cars that have the brake lights on.
8     num_cars_with_brake_lights_on = 0
9     for car_patch in car_patches:
10        if "yes" in car_patch.visual_question_answering("Are the brake lights on?"):
11            print(f"the car at {car_patch} has the brake lights on.")
12            num_cars_with_brake_lights_on += 1
13        else:
14            print(f"the car at {car_patch} does not have the brake lights on.")
15
16    return formatting_answer(str(num_cars_with_brake_lights_on))

```

Execution Trace

```

1 Calling find function. Detect car
2 Detection result: 669 103 779 286 car and 669 468 769 664 car and 668 705 747 991 car
3 Calling visual_question_answering function.
4 Question: Are the brake lights on?
5 Answer: yes
6 the car at 669 103 779 286 has the brake lights on.
7 Calling visual_question_answering function.
8 Question: Are the brake lights on?
9 Answer: yes
10 the car at 669 468 769 664 has the brake lights on.
11 Calling visual_question_answering function.
12 Question: Are the brake lights on?
13 Answer: no
14 the car at 668 705 747 991 does not have the brake lights on.
15 Program output: 2

```

Conversion to COT

The cars at **669 103 779 286** and **669 468 769 664** have the brake lights on. Thus, there are 2 cars with the brake lights on.

(c)



What is usually found in the same room as the word on the sign spelled backwards?

PaLM-2 Generated Program

```

1 def execute_command(image):
2     image_patch = ImagePatch(image)
3
4     # Find the word on the sign
5     word_on_the_sign = image_patch.visual_question_answering("What is the word on the sign?")
6
7     # Spell the word on the sign backwards
8     word_on_the_sign_backwards = word_on_the_sign[::-1]
9     print(f"The word on the sign backward is {word_on_the_sign_backwards}.")
10
11    # Find what is usually found in the same room as the word on the sign spelled backwards
12    object_in_the_same_room = language_question_answering(f"What is usually found in the same room as {word_on_the_sign_backwards}?")
13
14    return formatting_answer(object_in_the_same_room)

```

Execution Trace

```

1 Image Patch 165 0 832 1000 Calling visual_question_answering function.
2 Question: What is the word on the sign?
3 Answer: stop
4 The word on the sign backward is pots.
5 Calling language_question_answering function.
6 Question: What is usually found in the same room as pots?
7 Answer: pans
8 Program output: pans

```

Conversion to COT

The word on the sign is "stop". "Stop" spelled backwards is "pots". Pans are usually found in the same room as pots.

Figure 8. Examples of our data-synthesis pipeline. (a) is from GQA [26]; (b) is from TallyQA [1]; (c) is from A-OKVQA [52].

B. Experimental Details

1119

B.1. Details of Model Architectures

1120

For both PaLI models, the architectures are similar: images are encoded into visual tokens individually via a visual encoder. Then, the vision tokens along with the textual inputs are passed through an encoder-decoder UL2 Transformer [57]. The PaLI models were then pre-trained with image-text pairs to perform multimodal tasks. Specifically, PaLI-3 [11] uses a pre-trained 2B SigLIP [76] as visual encoder, and a 3B UL2. The image resolution is 812×812 . PaLI-X [10] uses a pre-trained VIT-22B [15] as visual encoder, and a 32B UL2. The image resolution is 756×756 . Please refer to the PaLI-3 [11] and PaLI-X [10] papers about more architecture details.

1121

1122

1123

1124

1125

1126

B.2. Details of Data

1127

The details of the data mixture of academic task-oriented VQA datasets used in VPD training are shown in Table 4. We only use a subset of each dataset’s training set. # labels refers to the total number of examples (containing image, query, and answer) we use. # CoTs refers to the number of examples that we have synthesized CoTs using our programs. In total, there are 89.6K CoTs used during training.

1128

1129

1130

1131

Dataset	Description	# labels	# CoTs
VQAv2 [19]	General	100.0K	
OCR-VQA [48]	OCR	50.0K	
GQA [26]	Compositional	86.0K	38.0K
OK-VQA [45]	Knowledge	9.0K	6.7K
A-OKVQA [52]	Knowledge	17.1K	11.2K
TallyQA [1]	Counting	48.4K	33.7K
Total		310.5K	89.6K

Table 4. Data mixture of academic task-oriented VQA datasets used in VPD training.

Details of each evaluation benchmark we use are in Table 5. For free-form question answering, we run inference with the prompt “Answer with a single word or phrase.”, using greedy decoding without any constraint on the model’s output space. For multiple-choice questions, we run inference with the prompt “Answer with the option letter from the given choices directly.” and generate the option letter.

1132

1133

1134

Dataset	Description	# split	# Metrics
VQAv2 [19]	General VQA. General questions about entities, colors, materials, etc.	test-dev	VQA Score
GQA [26]	Compositional VQA. Built on the scene-graphs in Visual Genome [32]. More compositional questions and spatial relation questions.	test-dev	EM
OK-VQA [45]	Knowledge-based VQA. Questions that need external knowledge to be answered.	val	VQA Score
A-OKVQA [52]	An advanced version of OK-VQA that is more challenging. – Multiple Choice (MC): choose 1 of the 4 options. – Direct Answer (DA): compare with 10 free-form human answers	val, test val, test	EM VQA Score
TallyQA [1]	Counting questions. – Simple: synthesized simple counting questions – Complex: human-written complex counting questions	test-simple test-complex	EM EM
TextVQA [55]	VQA on images that contain text	val	VQA Score
POPE [36]	Benchmark on VLM hallucination. Binary questions of whether an object exists in the image.	dev	EM
MMBench [40]	Comprehensive benchmark on VLMs with multiple-choice questions. Covering 20 ability dimensions across 3 levels (e.g., coarse perception, fine-grained perception, attribute reasoning, relation reasoning, logic reasoning, etc.)	dev	EM

Table 5. Summary of evaluation benchmarks.

1135

B.3. Training Details

1136

We use LoRA [22] to fine-tune both PaLI-3 [11] and PaLI-X [10]. For generalist training, we add LoRA weights on each linear layer in the attention blocks and multilayer perceptron (MLP) blocks for both the encoder and decoder in the UL2

1137

1138

1139 transformer. For both models, we use $rank = 8$. We use a cosine learning rate schedule, with warm-up ratio 1% and peak
1140 learning rate $1e - 4$. For all models and all settings, we use a batch size of 128 and fine-tune the pre-trained model for 8,000
1141 steps. In terms of training time, we train PaLI-X-VPD with 128 TPU-v3 [27] and it takes about 2 days to finish training. For
1142 PaLI-3-VPD, we use 32 TPU-v4 and training takes about 20 hours. We still observe a steady loss drop when we terminate
1143 training, which indicates that more computation may lead to even better performance. For per-task fine-tuning, to avoid
1144 overfitting, we reduce the number of training parameters. For both models, we only add LoRA weights to encoder layers. We
1145 use LoRA $rank = 4$ for PaLI-X-VPD and $rank = 8$ for PaLI-3-VPD. The peak learning rate is $1e - 4$ and we use cosine
1146 learning schedule, with warmup ratio 1%. For all per-task fine-tuning experiments, we use a batch size of 64. We train for 1
1147 epoch on GQA, and 3 epochs on all other datasets. We use the AdamW [30] optimizer with $\beta_1 = 0.9$ and $\beta_2 = 0.98$, and
1148 bfloat16 for all experiments.

1149 C. Human Evaluation

1150 We asked our human annotators to first evaluate each model’s answer, using the criteria described in §4.3. After rating each
1151 model answer separately, we also asked them to choose a preferred answer among the two. However, we observed that there
1152 are cases where one or both models have similar answers, or both answers are incorrect so it would be difficult to choose a
1153 favorite, so we allowed annotators to also choose “Both” or “Neither”, giving the annotators the following instruction: “Please
1154 try to choose “Answer 1 is better” or “Answer 2 is better” whenever possible. We also give you the option to choose “Both are
1155 equally good.” or “Both are too bad to make a choice.” for the cases when it is hard to make a choice either because both
1156 answers are correct and similar, or because both answers are wrong so it makes no sense to choose a favorite.”

1157 We show some examples from our human evaluation in Table 6. The table contains the images and corresponding
1158 text queries (column 2), the answers provided by the two models we compared—PaLI-X Instruct (column 3) and
1159 PaLI-X-VPD (column 4)—along with the corresponding annotations given by the human annotators. The human annotations
1160 are aggregated across 3 raters per sample. Finally, column 5 shows which of the two answers was preferred by the human
1161 raters. When a model’s answer includes a bounding box, we annotate it on the image for convenience.

- 1162 • Example #1 shows a common situation where PaLI-X-VPD succeeds where PaLI-X Instruct fails. By being trained
1163 with programs that include calls to an object detection tool, PaLI-X-VPD has learned to produce answers that localize the
1164 object in question in the image, which prods the model to correctly perform tasks such as counting.
- 1165 • Example #2 shows a type of question where neither model produces an explanation, where one is arguably not necessary.
1166 However, in spite the lack of explanation, PaLI-X-VPD’s answer is more accurate.
- 1167 • Example #3 shows an example where from Tally QA where PaLI-X-VPD is able to correctly localize and count the
1168 number of tennis balls in the image, as asked by the question, when PaLI-X Instruct counts incorrectly and also does
1169 not provide a rationale.
- 1170 • Example #4 shows a situation where PaLI-X-VPD produces an answer that is deemed both factual and consistent by the
1171 human raters, and yet the answer is incorrect. Indeed we can notice that everything that PaLI-X-VPD is accurate with
1172 respect to the context image, and yet it does not address the question as expected by the human raters.
- 1173 • Example #5 shows a case where both models are incorrect, but the annotators prefer PaLI-X-VPD’s answer, despite the
1174 fact that both models essentially give the same incorrect final count of 0 yellow busses. However, PaLI-X-VPD provides
1175 a consistent rationale, from which we can tell the reason for its mistake: it correctly identifies 2 out of the 3 buses in the
1176 image, missing the partially occluded yellow bus in the back. It seems the annotators prefer this kind of mistake over one
1177 without an explanation.






Index	Image & Query	PaLI-Instruct Answer & Human Eval	PaLI-VPD Answer & Human Eval	Annotator Preference
1.	 <p>On which side is the knife?</p>	<p>The knife is on the left side of the table.</p> <ul style="list-style-type: none"> • Correct: ✗ • Has Explanation: ✗ • Explanation Factuality: N/A • Explanation Consistency: N/A 	<p>The knife is at 55 693 412 884, which is on the right side of the image.</p> <ul style="list-style-type: none"> • Correct: ✓ • Has Explanation: ✓ • Explanation Factuality: ✓ • Explanation Consistency: ✓ 	PaLI-VPD
2.	 <p>What are these fruits called?</p>	<p>The fruits in the image are apples and broccoli.</p> <ul style="list-style-type: none"> • Correct: ✗ • Has Explanation: ✗ • Explanation Factuality: N/A • Explanation Consistency: N/A 	<p>The fruits are called apples.</p> <ul style="list-style-type: none"> • Correct: ✓ • Has Explanation: ✗ • Explanation Factuality: N/A • Explanation Consistency: N/A 	PaLI-VPD
3.	 <p>How many tennis balls are in the photo?</p>	<p>There are two tennis balls in the photo.</p> <ul style="list-style-type: none"> • Correct: ✗ • Has Explanation: ✗ • Explanation Factuality: N/A • Explanation Consistency: N/A 	<p>There is a tennis ball at 826 665 869 721. Thus, there is 1 tennis ball.</p> <ul style="list-style-type: none"> • Correct: ✓ • Has Explanation: ✓ • Explanation Factuality: ✓ • Explanation Consistency: ✓ 	PaLI-VPD
4.	 <p>Why is one player wearing gloves?</p>	<p>The player is wearing gloves because he is a batter. Batters wear gloves to protect their hands from the ball.</p> <ul style="list-style-type: none"> • Correct: ✓ • Has Explanation: ✗ • Explanation Factuality: ✓ • Explanation Consistency: ✓ 	<p>The player is wearing gloves to catch the ball.</p> <ul style="list-style-type: none"> • Correct: ✗ • Has Explanation: ✓ • Explanation Factuality: ✓ • Explanation Consistency: ✓ 	PaLI-Instruct
5.	 <p>How many yellow buses are in the picture?</p>	<p>There are no yellow buses in the picture.</p> <ul style="list-style-type: none"> • Correct: ✗ • Has Explanation: ✗ • Explanation Factuality: N/A • Explanation Consistency: N/A 	<p>There are 2 buses in the picture. Neither of them is yellow. Thus, there are 0 yellow buses in the picture.</p> <ul style="list-style-type: none"> • Correct: ✗ • Has Explanation: ✓ • Explanation Factuality: ✗ • Explanation Consistency: ✓ 	PaLI-VPD

Table 6. Human Annotation Examples. We show the answers provided by the two models, PaLI-X Instruct and PaLI-X-VPD, along with the aggregated evaluation scores by human raters along the criteria introduced in §4.3.

1178 **D. Prompts**

1179 In this section we present the prompts used in our data synthesis pipeline. Refer to §A for step-by-step examples of the
1180 programs, execution traces, and the converted CoTs.

1181 **D.1. Prompt for Code Generation**

1182 For each image and query, we put the query and a model-generated image caption in the prompt. An LLM takes this prompt
1183 and generate the program to answer the query. We modify the original ViperGPT [16] prompt to adapt to the vision tools we
1184 use in this paper.

```

1185 1 class ImagePatch:
1186 2     # A Python class containing a crop of an image centered around a particular object, as well as relevant information.
1187 3     # Attributes
1188 4     # -----
1189 5     # cropped_image : array_like
1190 6     #     An array-like of the cropped image taken from the original image.
1191 7     # left, lower, right, upper : int
1192 8     #     An int describing the position of the (left/lower/right/upper) border of the crop's bounding box in the original image.
1193 9     # Methods
1194 10    # -----
1195 11    # find(object_name: str)->List[ImagePatch]
1196 12    #     Returns a list of new ImagePatch objects containing crops of the image centered around any objects found in the
1197 13    #     image matching the object_name.
1198 14    # visual_question_answering(question: str=None)->str
1199 15    #     Returns the answer to a basic question asked about the image. If no question is provided, returns the answer to "What is this?".
1200 16    # image_caption()->str
1201 17    #     Returns a short description of the image crop.
1202 18    # expand_patch_with_surrounding()->ImagePatch
1203 19    #     Returns a new ImagePatch object that contains the current ImagePatch and its surroundings.
1204 20    # overlaps(patch: ImagePatch)->Bool
1205 21    #     Returns True if the current ImagePatch overlaps with another patch and False otherwise
1206 22    # compute_depth()->float
1207 23    #     Returns the median depth of the image patch. The bigger the depth, the further the patch is from the camera.
1208 24
1209 25    def __init__(self, image, left: int = None, lower: int = None, right: int = None, upper: int = None):
1210 26    #     Initializes an ImagePatch object by cropping the image at the given coordinates and stores the coordinates as
1211 27    #     attributes. If no coordinates are provided, the image is left unmodified, and the coordinates are set to the
1212 28    #     dimensions of the image.
1213 29    # Parameters
1214 30    # -----
1215 31    # image: PIL.Image
1216 32    #     An array-like of the original image.
1217 33    # left, lower, right, upper : int
1218 34    #     An int describing the position of the (left/lower/right/upper) border of the crop's bounding box in the original image.
1219 35    #     The coordinates (y1,x1,y2,x2) are with respect to the upper left corner the original image.
1220 36    #     To be closer with human perception, left, lower, right, upper are with respect to the lower left corner of the squared image.
1221 37    #     Use left, lower, right, upper for downstream tasks.
1222 38
1223 39    self.original_image = image
1224 40    self.size_x, self.size_y = image.size
1225 41
1226 42    if left is None and right is None and upper is None and lower is None:
1227 43        self.x1 = 0
1228 44        self.y1 = 0
1229 45        self.x2 = 999
1230 46        self.y2 = 999
1231 47    else:
1232 48        self.x1 = left
1233 49        self.y1 = 999 - upper
1234 50        self.x2 = right
1235 51        self.y2 = 999 - lower
1236 52
1237 53    self.cropped_image = image.crop((int(self.x1/1000*self.sz), int(self.y1/1000*self.sz),
1238 54                                     int(self.x2/1000*self.sz), int(self.y2/1000*self.sz)))
1239 55
1240 56    self.width = self.x2 - self.x1
1241 57    self.height = self.y2 - self.y1
1242 58
1243 59    # all coordinates use the upper left corner as the origin (0,0).
1244 60    # However, human perception uses the lower left corner as the origin.
1245 61    # So, need to revert upper/lower for language model
1246 62    self.left = self.x1
1247 63    self.right = self.x2
1248 64    self.upper = 999 - self.y1
1249 65    self.lower = 999 - self.y2
1250 66
1251 67    self.horizontal_center = (self.left + self.right) / 2
1252 68    self.vertical_center = (self.lower + self.upper) / 2
1253 69
1254 70    self.patch_description_string = f"({self.y1} {self.x1} {self.y2} {self.x2})"
1255 71
1256 72    def __str__(self):
1257 73    return self.patch_description_string
1258 74
1259 75    def compute_depth(self):
1260 76    # compute the depth map on the full image. Returns a np.array with size 192*192
1261 77    # Parameters
1262 78    # -----
1263 79    # Returns
1264 80    # -----
1265 81    # float
1266 82    #     the median depth of the image crop
1267 83
1268 84    # Examples
1269 85    # -----
1270 86    # >>> return the image patch of the bar furthest away
1271 87    # >>> def execute_command(image)->ImagePatch:
1272 88    # >>> image_patch = ImagePatch(image)

```

```
89 # >>> bar_patches = image_patch.find("bar")
90 # >>> bar_patches . sort(key=lambda bar: bar.compute_depth())
91 # >>> return bar_patches[-1]
92
93 return depth(self .cropped_image)
94
95 def find(self, object_name: str):
96 # Returns a list of ImagePatch objects matching object_name contained in the crop if any are found.
97 # The object_name should be as simple as example, including only nouns
98 # Otherwise, returns an empty list .
99 # Note that the returned patches are not ordered
100 # Parameters
101 # -----
102 # object_name : str
103 #     the name of the object to be found
104
105 # Returns
106 # -----
107 # List[ImagePatch]
108 #     a list of ImagePatch objects matching object_name contained in the crop
109
110 # Examples
111 # -----
112 # >>> # find all the kids in the images
113 # >>> def execute_command(image) -> List[ImagePatch]:
114 # >>>     image_patch = ImagePatch(image)
115 # >>>     kid_patches = image_patch.find("kid")
116 # >>>     return kid_patches
117
118 print(f"Calling find function . Detect {object_name}.")
119 det_patches = detect(self .cropped_image, object_name)
120 print(f"Detection result : { ' and ' . join([str(d) + ' ' + object_name for d in det_patches ])}")
121
122 return det_patches
123
124
125 def expand_patch_with_surrounding(self):
126 # Expand the image patch to include the surroundings . Now done by keeping the center of the patch
127 # and returns a patch with double width and height
128
129 # Examples
130 # -----
131 # >>> # How many kids are not sitting under an umbrella?
132 # >>> def execute_command(image):
133 # >>>     image_patch = ImagePatch(image)
134 # >>>     kid_patches = image_patch.find("kid")
135
136 # >>> # Find the kids that are under the umbrella.
137 # >>>     kids_not_under_umbrella = []
138
139 # >>> for kid_patch in kid_patches:
140 # >>>     kid_with_surrounding = kid_patch.expand_patch_with_surrounding()
141 # >>>     if "yes" in kid_with_surrounding . visual_question_answering("Is the kid under the umbrella?"):
142 # >>>         print(f"the kid at {kid_patch} is sitting under an umbrella.")
143 # >>>     else:
144 # >>>         print(f"the kid at {kid_patch} is not sitting under an umbrella.")
145 # >>>         kids_not_under_umbrella.append(kid_patch)
146
147 # >>> # Count the number of kids under the umbrella.
148 # >>>     num_kids_not_under_umbrella = len(kids_not_under_umbrella)
149
150 # >>>     return formatting_answer(str(num_kids_not_under_umbrella))
151
152 new_left = max(self .left - self .width / 2, 0)
153 new_right = min(self .right + self .width / 2, 999)
154 new_lower = max(self .lower - self .height / 2, 0)
155 new_upper = min(self .upper + self .height / 2, 999)
156
157 return ImagePatch(self .original_image, new_left, new_lower, new_right, new_upper)
158
159
160 def visual_question_answering(self, question: str = None) -> str:
161 # Returns the answer to a basic question asked about the image.
162 # The questions are about basic perception, and are not meant to be used for complex reasoning
163 # or external knowledge.
164
165 # Parameters
166 # -----
167 # question : str
168 #     A string describing the question to be asked.
169
170 # Examples
171 # -----
172
173 # >>> # What is the name of the player in this picture?
174 # >>> def execute_command(image) -> str:
175 # >>>     image_patch = ImagePatch(image)
176 # >>>     return formatting_answer(image_patch .visual_question_answering("What is the name of the player?"))
177
178 # >>> # What color is the foo?
179 # >>> def execute_command(image) -> str:
180 # >>>     image_patch = ImagePatch(image)
181 # >>>     return formatting_answer(foo_patch .visual_question_answering("What color is the foo?"))
182
183 # >>> # What country serves this kind of food the most?
184 # >>> def execute_command(image) -> str:
185 # >>>     image_patch = ImagePatch(image)
186 # >>>     food_name = image_patch .visual_question_answering("What kind of food is served?")
187 # >>>     country = language_question_answering(f"What country serves {food_name} most?", long_answer=False)
188 # >>>     return formatting_answer(country)
189
190 # >>> # Is the second bar from the left quaxy?
191 # >>> def execute_command(image) -> str:
192 # >>>     image_patch = ImagePatch(image)
193 # >>>     bar_patches = image_patch.find("bar")
194 # >>>     bar_patches . sort(key=lambda x: x .horizontal_center )
```

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194

1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403
1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457
1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500

```
195 # >>> bar_patch = bar_patches[1]
196 # >>> return formatting_answer(bar_patch.visual_question_answering("Is the bar quuxy?"))
197
198 answer = vqa(self.cropped_image, question)
199
200 print(f"Calling visual_question_answering function.")
201 print(f"Question: {question}")
202 print(f"Answer: {answer}")
203
204 return answer
205
206 def image_caption(self) -> str:
207     # Returns a short description of the image.
208     return image_caption(self.cropped_image)
209
210 def overlaps(self, patch) -> bool:
211     # check if another image patch overlaps with this image patch
212     # if patch overlaps with current patch, return True. Otherwise return False
213
214     if patch.right < self.left or self.right < patch.left:
215         return False
216     if patch.lower > self.upper or self.lower > patch.upper:
217         return False
218     return True
219
220 def language_question_answering(question: str, long_answer: bool = False) -> str:
221     # Answers a text question using a language model like PaLM and GPT-3. The input question is always a formatted string with a variable in it.
222     # Default is short-form answers, can be made long-form responses with the long_answer flag.
223
224     # Parameters
225     # -----
226     # question: str
227     #     the text question to ask. Language model cannot understand the image. Must not contain any reference to 'the image' or 'the photo', etc.
228     # long_answer: bool
229     #     whether to return a short answer or a long answer. Short answers are one or at most two words, very concise.
230     #     Long answers are longer, and may be paragraphs and explanations. Default is False.
231
232     # Examples
233     # -----
234     # >>> # What is the city this building is in?
235     # >>> def execute_command(image) -> str:
236     # >>>     image_patch = ImagePatch(image)
237     # >>>     building_name = image_patch.visual_question_answering("What is the name of the building?")
238     # >>>     return formatting_answer(language_question_answering(f"What city is {building_name} in?", long_answer=False))
239
240     # >>> # Who invented this object?
241     # >>> def execute_command(image) -> str:
242     # >>>     image_patch = ImagePatch(image)
243     # >>>     object_name = image_patch.visual_question_answering("What is this object?")
244     # >>>     return formatting_answer(language_question_answering(f"Who invented {object_name}?", long_answer=False))
245
246     # >>> # Explain the history behind this object.
247     # >>> def execute_command(image) -> str:
248     # >>>     image_patch = ImagePatch(image)
249     # >>>     object_name = image_patch.visual_question_answering("What is the object?")
250     # >>>     return formatting_answer(language_question_answering(f"What is the history behind {object_name}?", long_answer=True))
251
252     print(f"Calling language_question_answering")
253     print(f"Question: {question}")
254
255     answer = language_model_qa(question, long_answer).lower().strip()
256     print(f"Answer: {answer}")
257     return answer
258
259 def distance(patch_a: Union[ImagePatch, float], patch_b: Union[ImagePatch, float]) -> float:
260     # Returns the distance between the edges of two ImagePatches, or between two floats.
261     # If the patches overlap, it returns a negative distance corresponding to the negative intersection over union.
262
263     # Parameters
264     # -----
265     # patch_a : ImagePatch
266     # patch_b : ImagePatch
267
268     # Examples
269     # -----
270     # # Return the qux that is closest to the foo
271     # >>> def execute_command(image):
272     # >>>     image_patch = ImagePatch(image)
273     # >>>     qux_patches = image_patch.find('qux')
274     # >>>     foo_patches = image_patch.find('foo')
275     # >>>     foo_patch = foo_patches[0]
276     # >>>     qux_patches.sort(key=lambda x: distance(x, foo_patch))
277     # >>>     return qux_patches[0]
278
279     return dist(patch_a, patch_b)
280
281
282 def formatting_answer(answer) -> str:
283     # Formatting the answer into a string that follows the task's requirement
284     # For example, it changes bool value to "yes" or "no", and clean up long answer into short ones.
285     # This function should be used at the end of each program
286
287     final_answer = ""
288     if isinstance(answer, str):
289         final_answer = answer.strip()
290
291     elif isinstance(answer, bool):
292         final_answer = "yes" if answer else "no"
293
294     elif isinstance(answer, list):
295         final_answer = ", ".join([str(x) for x in answer])
296
297     elif isinstance(answer, ImagePatch):
298         final_answer = answer.image_caption()
299
300
```

```

301 else :
302     final_answer = str(answer)
303
304     print(f"Program output: {final_answer}")
305     return final_answer
306
307
308 Given an image and a query, write the function execute_command using Python and the ImagePatch class (above), and the other functions above that could be executed to provide an answer to the query.
309 For reference, a model generated image description is also provided, so that the function can be customized for the given image. The image description is model-generated and may not be reliable, so do not trust it.
310
311 Consider the following guidelines :
312 - Use base Python (comparison, sorting) for basic logical operations, left / right / up/down, math, etc.
313 - Use the language_question_answering function to access external information and answer informational questions NOT concerning the image.
314 - The program should print out the intermediate traces as it runs. So add print function in the program if needed.
315
316 For usual cases, follow the guidelines below:
317 - For simple visual queries, directly call visual_question_answering to get the answer.
318 - For queries that need world knowledge, commonsense knowledge, and language reasoning, use visual_question_answering, language_question_answering, and sometimes image_caption to get the answer.
319 - For queries that require counting and spatical relations, in addition to the above functions, use find function to help getting the answer.
320 - For queries involve "behind" and "front", consider using compute_depth function.
321
322
323 Some examples:
324 Image description : a woman is walking several dogs
325 Query: how many dogs are to left of the person?
326 Function:
327 def execute_command(image):
328     image_patch = ImagePatch(image)
329     person_patch = image_patch.find("person")[0]
330     dog_patches = image_patch.find("dog")
331
332     # Count the number of dogs whose leftmost x-coordinate is less than the person.
333     num_dogs_left = 0
334     for dog_patch in dog_patches:
335         if dog_patch.left < person_patch.horizontal_center :
336             print(f"dog at {dog_patch} is on the left of human.")
337             num_dogs_left += 1
338
339     return formatting_answer(num_dogs_left)
340
341 # [other in-context examples]
342
343 Image description : INSERT_IMAGE_CAPTION_HERE
344 Query: INSERT_QUERY_HERE
345 Function:

```

1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511
1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1529

D.2. Prompt for Result verification

1530

After running each program, we get an output. As discussed in §3.1, we adopt the method of [28] and use an LLM to determine if the program output matches human answers. The LLM takes the program output and reference answers as input. The prompt is as follows:

1531

1532

1533

```

1 Given a visual question, several human annotator answers, and a candidate answer, determine if the candidate is correct.
2 The candidate is considered correct if it is allowed to have formatting differences compared with the human answers.
3 If the candidate is correct, return the gold answer it matches. Otherwise, return None.
4
5 Question: INSERT_QUESTION_HERE
6 Answers: INSERT_ANSWERS_HERE
7 Candidate: INSERT_CANDIDATE_HERE
8 Is the candidate correct?

```

1534

1535

1536

1537

1538

1539

1540

1541

D.3. Prompt for CoT conversion

1542

Finally, once a program is filtered, we convert its execution trace into chain-of-thought using an LLM. The LLM takes the query, program, execution trace, program output as input, and summarizes the execution trace into a chain-of-thought rationale. The prompt we use as as follows:

1543

1544

1545

```

1 Given an image and a question, I wrote the function execute_command using Python and the ImagePatch class (above), and the other functions above that could be executed to provide an answer to the query.
2 As shown in the code, the code will print execution traces.
3 I need you to rewrite the execution trace into a natural language rationale that leads to the answer.
4
5 Consider the following guidelines :
6 - Use the bounding box information in the rationale.
7 - Referencing the execution trace, write a reasoning chain that leads to the most common human answer. Notice that the output should be the same as the human answer, not necessarily the program output.
8
9
10 Some examples:
11 Question: How many wheels does the plane have?
12 Program:
13 def execute_command(image):
14     image_patch = ImagePatch(image)
15
16     # Find the plane in the image
17     plane_patch = image_patch.find("plane")[0]
18
19     # Count the number of wheels on the plane
20     num_wheels = 0
21     for wheel in plane_patch.find("wheel"):
22         num_wheels += 1
23
24     return formatting_answer(str(num_wheels))
25 Execution trace :
26 Calling find function. Detect plane
27 Detected plane at 153 25 647 972
28 Calling find function. Detect wheel

```

1546

1547

1548

1549

1550

1551

1552

1553

1554

1555

1556

1557

1558

1559

1560

1561

1562

1574
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586

```
29 Detected wheel at 603 471 649 515
30 Detected wheel at 621 85 646 113
31 Detected wheel at 615 383 645 428
32 Program output: 3
33 Rationale: The plane at 153 25 647 972 has wheels at 603 471 649 515, 621 85 646 113, and 615 383 645 428.Thus, it has 3 wheels.
34
35 [Other demonstration examples]
36
37 Question: INSERT_QUESTION_HERE
38 INSERT_PROGRAM_HERE
39 Execution trace :
40 INSERT_EXECUTION_TRACE_HERE
41 Rationale:
```