# CLOVA: A <u>C</u>losed-<u>LO</u>op <u>V</u>isual <u>A</u>ssistant with Tool Usage and Update

Zhi Gao[1,2], Yuntao Du[2], Xintong Zhang[2,3], Xiaojian Ma[2],
Wenjuan Han[3], Song-Chun Zhu[1,2,4], Qing Li[2 ✉]

[1]School of Intelligence Science and Technology, Peking University [2]State Key Laboratory of General Artificial Intelligence, BIGAI
[3]Beijing Jiaotong University [4]Department of Automation, Tsinghua University
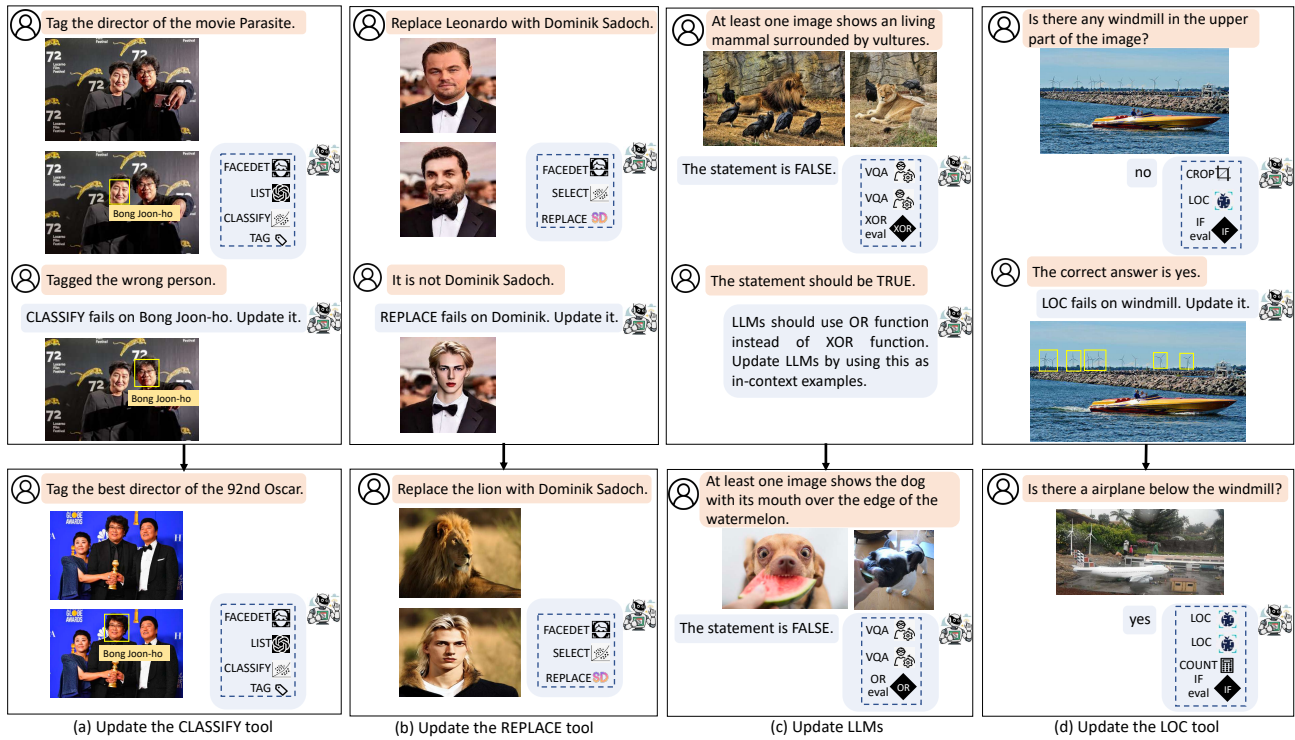
https://clova-tool.github.io

Figure 1. CLOVA is a general visual assistant that updates both LLMs and visual tools via inference, reflection, and learning in a closed-loop framework. During inference, CLOVA uses LLMs to integrate visual tools to accomplish given tasks. In reflection, CLOVA identifies tools that require updating based on human feedback. Finally, in learning, CLOVA collects data and updates the tools accordingly.

## Abstract

*Utilizing large language models (LLMs) to compose off-the-shelf visual tools represents a promising avenue of research for developing robust visual assistants capable of addressing diverse visual tasks. However, these methods often overlook the potential for **continual learning**, typically by freezing the utilized tools, thus limiting their adaptation to environments requiring new knowledge. To tackle this challenge, we propose CLOVA, a <u>C</u>losed-<u>LO</u>op <u>V</u>isual <u>A</u>ssistant, which operates within a framework encompassing inference, reflection, and learning phases. During the inference phase, LLMs generate programs and execute corresponding tools to complete assigned tasks. In the reflection phase, a multimodal global-local reflection scheme analyzes human feedback to determine which tools require updating. Lastly, the learning phase employs three flexible approaches to automatically gather training data and introduces a novel prompt tuning scheme to update the tools, allowing CLOVA to efficiently acquire new knowledge. Experimental findings demonstrate that CLOVA surpasses existing tool-usage methods by 5% in visual question answering and multiple-image reasoning, by 10% in knowledge tagging, and by 20% in image editing. These results underscore the significance of the continual learning capability in general visual assistants.*

---

✉ Corresponding author: Qing Li (dylan.liqing@gmail.com).

# 1. Introduction

The advancement of large language models (LLMs) [48, 69] alongside various visual tools (e.g., neural networks and OpenCV functions) [6, 29, 46, 57, 58] offers feasible avenues for constructing general visual assistants. When confronted with a task accompanied by language instructions, a common approach involves harnessing LLMs to generate programs, which are then executed using readily available visual tools to solve the task as dictated by the generated program [8, 11, 34, 42, 62, 67]. For instance, when posed with the query "*What is the person to the left of the umbrella doing?*", a viable solution entails LLMs sequentially performing the following steps: (1) utilizing a detection tool to locate the umbrella, (2) cropping the image region to the left of the umbrella, (3) using the detection tool to locate the person, and (4) finally querying a Visual Question Answering (VQA) tool with the question "*What is the person doing?*". Being highly compositional, such tool-usage methods demonstrate impressive performance and appealing explainability in tackling complex reasoning tasks, including VQA [12, 31, 32], mathematical reasoning [14, 15], and image editing [40, 75, 80, 81].

However, the potential for continual learning has been largely overlooked in existing tool-usage methods. Most of them simply freeze the used tools, which limits their applicability to environments where new knowledge is required, as depicted in Fig. 1. For instance, a user might instruct a visual assistant to label the face of the movie director Bong Joon-ho in a photograph. However, if the face recognition tool employed by the assistant fails to recognize Bong Joon-ho, it may provide an incorrect response. In such scenarios, it is expected that the assistant can learn this missing information about Bong Joon-ho and generalize it to other photographs. Thus, it is imperative to endow visual assistants with the learning capability, enabling them to swiftly acquire new knowledge from failures.

In this paper, we propose CLOVA, a <u>C</u>losed-<u>LO</u>op <u>V</u>isual <u>A</u>ssistant that updates used tools via closed-loop learning [33, 35] to better adapt to new environments, as illustrated in Figure 1. CLOVA consists of three phases: inference, reflection, and learning. During inference, CLOVA employs Large Language Models (LLMs) to generate programs and execute corresponding tools to accomplish the task at hand. Subsequently, in the reflection phase, CLOVA utilizes human feedback to provide critiques, identifying tools that require updates. Finally, in the learning phase, CLOVA autonomously collects data and updates tools accordingly. Thus, CLOVA facilitates autonomous tool updating, thereby continually enhancing their ability to adapt to diverse environments.

To establish such a closed-loop learning framework, we must address three key challenges. Firstly, identifying tools that require updates is difficult due to the multi-step nature of generated programs and the diversity of errors within them. Secondly, automatically collecting training data is necessary as the knowledge to be learned is unpredictable. Thirdly, efficiently updating tools presents another obstacle, considering their scale and the quality of the collected data. Visual tools typically involve large neural networks, making them inefficient to update, and naive fine-tuning could result in unacceptable catastrophic forgetting [44]. Moreover, the presence of noise within the collected data further complicates the training process.

We propose several techniques to tackle these challenges. First, we introduce a multimodal global-local reflection scheme, which resorts to LLMs to identify tools that need to be updated from both global and local aspects. For the second challenge, three data collection manners are employed, including inferring answers by LLMs, searching on the Internet, and searching from open-vocabulary datasets. Lastly, we develop a training-validation prompt tuning scheme for the tools, which includes instance-wise prompt tuning and a subsequent prompt validation stage, where learned prompts that fail to predict the validation data will be discarded. The learning phase also updates LLMs by storing correct examples and incorrect examples with critiques as in-context examples, which will be used in future inference. As a result, CLOVA efficiently updates tools in challenging environments with noisy data, while avoiding catastrophic forgetting.

We apply CLOVA to compositional VQA and multiple-image reasoning tasks, using the GQA [20] and NLVRv2 [66] datasets. Additionally, we manually collect data for image editing and factual knowledge tagging tasks. CLOVA outperforms existing tool-usage methods by 5% in compositional VQA and multiple-image reasoning tasks, by 10% in knowledge tagging tasks, and by 20% in image editing tasks, showing the significance of the learning capability for general visual assistants.

In summary, our contributions are three-fold:
- We build CLOVA, a visual assistant that updates its tools within a closed-loop learning framework for better adaptation to new environments.
- We propose a multimodal global-local reflection scheme, capable of identifying tools in need of updates.
- We employ three flexible manners to automatically collect training data and introduce a novel training-validation prompt tuning scheme to update tools efficiently while avoiding catastrophic forgetting.

# 2. Related Work

## 2.1. General Visual Assistant

Benefiting from the advancements of LLMs [48, 69] and visual tools [26, 46, 57, 58], visual assistants have achieved great progresses. Some methods concatenate

| Method | Visual Tool | Reflection | Update LLMs | Update VTs |
|---|---|---|---|---|
| ART [50] | ✗ | ✗ | Prompt | - |
| TRICE [55] | ✗ | Global | Instruction + RL | - |
| ToolkenGPT [13] | ✗ | - | ✗ | - |
| Toolformer [60] | ✗ | - | Fine-tune | - |
| VISPROG [11] | ✓ | ✗ | ✗ | ✗ |
| Visual ChatGPT [75] | ✓ | ✗ | ✗ | ✗ |
| HuggingGPT [62] | ✓ | ✗ | ✗ | ✗ |
| ViperGPT [67] | ✓ | ✗ | ✗ | ✗ |
| GPT4TOOLs [80] | ✓ | ✗ | Instruction | ✗ |
| OpenAGI [9] | ✓ | ✗ | RL | ✗ |
| AssistGPT [8] | ✓ | Global | Prompt | ✗ |
| **CLOVA (Ours)** | ✓ | Global+Local | Prompt | Prompt |

Table 1. Comparisons with representative tool-usage methods, where VTs means visual tools.

and train LLMs with visual tools in an end-to-end manner, where representative work includes LLaVA [38], Otter [27], MMICL [83], Kosmos-2 [53], and Flamingo [1], *etc*. In addition, some work extends the idea of tool usage for AI assistants from natural language processing [4, 13, 50, 55, 56, 60] to computer vision. By providing in-context examples, VISPROG [11] and ViperGPT [67] generate programs to use visual tools. Following this idea, some work improves performance by collecting instruction-following data [39, 52, 80], adding more tools [36, 62], and designing more dedicated tool-usage procedures [16, 40–42, 75, 76, 81]. The most related work to CLOVA is AssistGPT [8] and OpenAGI [9]. The two methods update LLMs after development through in-context learning and reinforcement learning, respectively. Different from them, CLOVA can update both LLMs and visual tools via its reflection and learning phases. This allows CLOVA to better adapt to new environments. In addition, the closed-loop framework enables us to set a separate training stage for tool-usage methods, going beyond zero-shot or few-shot visual assistants. Comparisons between CLOVA and some representative tool-usage methods are shown in Tab. 1.

### 2.2. Reflection of LLMs

Reflection has become a remedy in case LLMs cannot generate good responses in a single attempt [5, 49, 51, 79]. Reflection methods send outputs back to LLMs to obtain critiques and further improve the outputs. These critiques take the form of scores or natural language [45, 47, 64, 87]. To generate better critiques, some methods employ instruction tuning [59, 78] or reinforcement learning [3, 55]. Recently, Huang *et al.* [18] revealed that LLMs struggle to provide accurate critiques for complex tasks. One way to address this issue is incorporating external information such as human-desired results into LLMs [2, 77, 79]. Unlike existing methods that rely solely on feedback in the language modality, our method generates reflection using all multimodal intermediate results. In addition, our method incorporates both global and local aspects for reflection, instead of only the global aspect. These bring more effective critiques for compositional tasks.
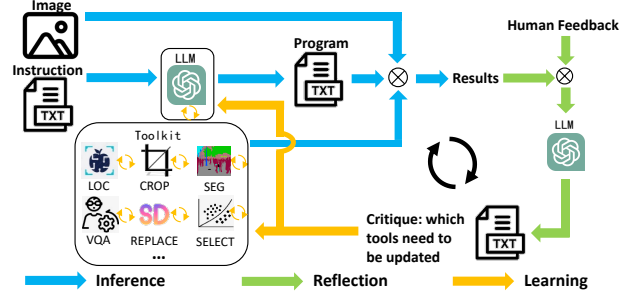


Figure 2. Framework of CLOVA.

### 2.3. Prompt-based Learning

Prompt-based learning is an efficient technique to update neural networks, achieving impressive performance in both NLP [22, 63] and computer vision [19, 61, 88, 89]. Prompt engineering and prompt tuning are two kinds of commonly used methods. Prompt engineering develops interpretable tokens (*e.g.*, texts and image regions) to guide model prediction, which are usually obtained by manually designing [57, 68], retrieval [82], and model generation [17, 54]. Prompt tuning learns vectors as prompts via gradient-based optimization. VPT [21] and CoOp [86] learn prompts for vision encoders and text encoders, respectively. To handle diverse data, CoCoOp [85] learns to generate prompts for unknown classes, ProDA [43] builds a Gaussian distribution for prompts, and MaPLe [24] learns both text and visual prompts. In addition, some methods employ prompts for continual learning, which learn prompts for different classes and produce adaptive prompts during inference. Representative methods include PIVOT [71], DualPrompt [73], and L2P [74]. Different from existing methods, our training-validation prompt tuning scheme discards harmful prompts, leading to more stable learning processes when the quality of training data is subpar.

## 3. Method

### 3.1. Overview

As shown in Fig. 2, CLOVA has three phases: inference, reflection, and learning. In the inference phase, CLOVA uses LLMs to generate programs and executes corresponding tools to solve the task. The reflection phase introduces a multimodal global-local reflection scheme that uses LLMs to generate critiques, identifying which tool needs to be updated. During learning, we employ three manners to collect training data and use a training-validation prompt tuning scheme to update the tools.

### 3.2. Inference

Our inference phase is based on VISPROG [11], while the difference is that CLOVA first uses LLMs to generate plans and then generates programs based on the plans, instead of

| Tool Type | Tool Name | Tool Description | Data Collection |
|---|---|---|---|
| Tools to be updated | LOC | Use the OWL-ViT model [46] for object localization | Open-vocabulary dataset |
| | VQA | Use the BLIP model [29] for VQA | LLM inference |
| | SEG | Use the maskformer model [6] for panoptic segmentation | Open-vocabulary dataset |
| | SELECT | Use the CLIP model [57] to select the most relevant object, given a text description | Internet |
| | CLASSIFY | Use the CLIP model [57] to classify given images | Internet |
| | REPLACE | Use the stable diffusion inpainting model [58] to replace one object with another desirable object | Internet |
| Tools not to be updated | FACEDET | Use the DSFD model [28] for face detection | N/A |
| | LIST | Use the text-davinci-002 model of OpenAI for knowledge retrieval | N/A |
| | EVAL | Use the Python function eval() to process string expressions for answers | N/A |
| | RESULT | Use the Python function dict() to output the intermediate and final results | N/A |
| | COUNT | Use Python function len() to count the number of input bounding boxes or masks | N/A |
| | CROP | Use Python function PIL.crop() to crop images | N/A |
| | COLORPOP | Use Python function PIL.convert() to keep desirable objects in color and other regions gray | N/A |
| | BGBLUR | Use Python function PIL.GaussianBlur() to blur the background | N/A |
| | EMOJI | Use emojis in the Python packages AngLy(pypi) to hide someone's face | N/A |

Table 2. Used tools in CLOVA, categorized based on whether the tool is updated in our method. Details of tool updates are in Sec. 3.4
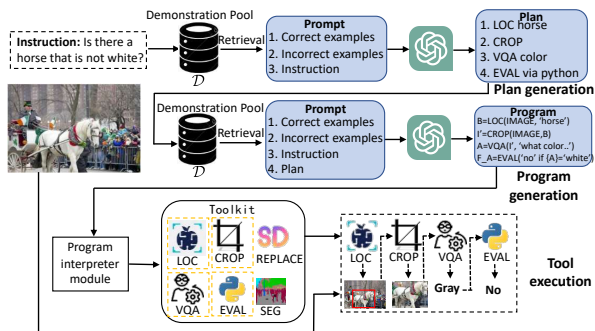


Figure 3. Illustration of the inference phase in CLOVA.

directly generating programs. Plans can be seen as intermediate reasoning chains that benefit the inference and reflection phases. Given a task, CLOVA selects in-context examples from a demonstration pool $\mathcal{D}$ (the construction of $\mathcal{D}$ will be detailed in Sec. 3.4.2), including correct examples and incorrect examples with error critiques. These examples are used to create prompts that are then sent to LLMs for plan and program generation. Finally, the program is parsed to execute visual tools (see Fig. 3).

**Plan generation.** The demonstration pool $\mathcal{D}$ is composed by $\mathcal{D} = \{\mathcal{D}_{p,s}, \mathcal{D}_{p,f}, \mathcal{D}_{c,s}, \mathcal{D}_{c,f}\}$, where $\mathcal{D}_{p,s}$ and $\mathcal{D}_{p,f}$ contain correct and incorrect examples for plan generation respectively, and $\mathcal{D}_{c,s}$ and $\mathcal{D}_{c,f}$ contain correct and incorrect examples for program generation respectively. Given a task, we use the BERT model [23] to extract features of the given instruction and examples stored in $\mathcal{D}_{p,s}$ and $\mathcal{D}_{p,f}$. Then, we combine similar examples in $\mathcal{D}_{p,s}$ and $\mathcal{D}_{p,f}$ with the instruction to create a prompt. Finally, we send the prompt to LLMs to generate the plan in a one-go manner.

**Program generation.** Similar to plan generation, we use LLMs to generate programs in a one-go manner. We select correct and incorrect examples of programs from $\mathcal{D}_{c,s}$ and $\mathcal{D}_{c,f}$. We combine these examples with the plan as a prompt and then send the prompt to LLMs for program generation.

**Tool execution.** We utilize the interpreter module in [11] to parse the program, extracting tool names, inputs, and outputs of each step. CLOVR activates tools from a toolkit
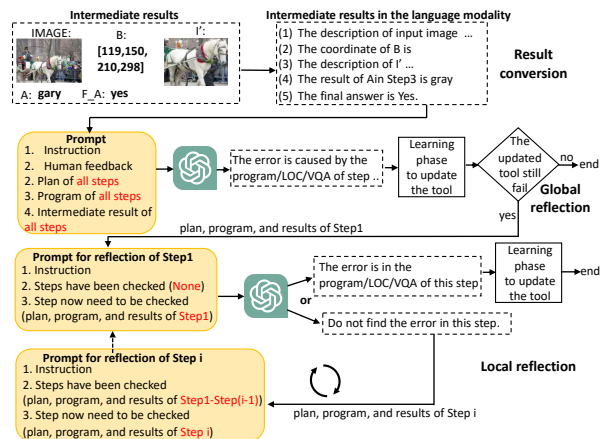


Figure 4. Illustration of the reflection phase in CLOVA.

$\mathcal{T}$ that contains 15 tools, including neural networks and Python functions, as shown in Tab. 2.

## 3.3. Reflection

In the inference phase, if a task is not solved correctly, the multimodal global-local reflection scheme uses LLMs to generate critiques, identifying which tool needs to be updated, as shown in Fig. 4.

**Result conversion.** Since LLMs often struggle to identify errors by themselves [18, 65, 70], we provide the human feedback, our wrong results, and intermediate results of each step for LLMs to better identify the error source. This requires us to convert visual results into textual form. For this purpose, we use the BLIP [29] model to convert images into languages by captioning.

**Global reflection.** CLOVA first uses global reflection to generate critiques in a one-go manner. The prompts are composed of task inputs, feedback on the task (*e.g.*, desirable results in VQA tasks, or human comments in image editing tasks), generated plans and programs, and intermediate results at each step. We send the prompts to LLMs to generate critiques that are used to update tools in the learning phase.

**Local reflection.** If CLOVA still fails after the tools are updated via global reflection and the learning phase–meaning the actual tools that lead to the faulty response are still to be found, we resort to local reflection to analyze each step of the program. Prompts are composed of the task inputs, feedback on the task, the steps that have been checked, and the current step that needs to be checked. Each step includes plans, programs, and intermediate results. We send the prompts to LLMs to infer whether this step has errors and the reasons. Local reflection continues until an error location and reasons are identified for a step.

## 3.4. Learning

### 3.4.1 Updating tools with prompt tuning

After identifying tools that need to be updated from the reflection phase, CLOVA then moves to the learning phase to collect training data and goes through training-validation prompt tuning to update the tools, as shown in Fig. 5.

**Data collection.** Since the tools that need to be updated can be rather different (a full list can be found in Tab. 2), we explore three manners to collect data online. (1) We use LLMs to generate training data for the VQA tool. If reflection concludes that the VQA tool makes errors, we combine the desirable response of the whole task and intermediate results of each step to prompt LLMs into inferring the correct output of the VQA tool. The question and the inferred output are then used to update the VQA tool. (2) We gather training data from open-vocabulary visual object grounding datasets (*e.g.*, LVIS [10]) for the LOC and SEG tools. For example, if the reflection phase indicates that LOC does not work well for the visual concept "*horse*", CLOVA will select images and bounding boxes of horses from LVIS to update LOC. (3) We collect data by searching on the Internet for the SELECT, CLASSIFY, and REPLACE tools. For instance, If CLASSIFY is marked as unable to recognize "*horse*" during the reflection phase, CLOVA will search the Internet for images of horses to update CLASSIFY.

**Prompt tuning and validation.** Given the collected data, we invoke training-validation prompt tuning to update tools. Note that, instead of learning a single prompt for all collected data, we choose to learn a prompt for each training instance collected. Each learned prompt will then be validated by running the tool with it on validation data (held out from collected data except for the VQA tool, where VQA will be validated on the original visual question it failed on) and seeing if the tool can produce desirable responses (*e.g.* correctly localizing a horse for the LOC tool). As a result, we discard prompts that do not lead to the desirable responses, possibly due to the faulty training instances they were trained on, alleviating the issue of noisy collected data. Finally, we build a prompt pool $\mathcal{P}$ for each tool. Take the LOC tool as an example. After training and validation, CLOVA stores the visual concept (*e.g.*, "horse" in the
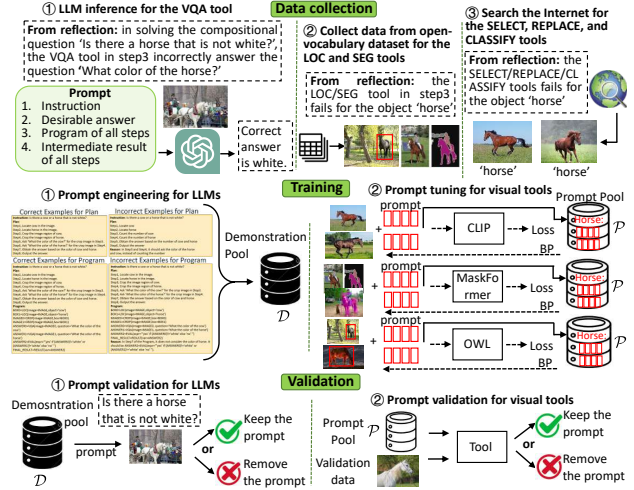


Figure 5. Illustration of the learning phase in CLOVA.

task of localizing horses) with its learned prompts and visual features of all collected instances in $\mathcal{P}$, formulated as $\mathcal{P} = \left\{ v_j : \left[ [f_{j1}, \cdots, f_{jn}], [p_{j1}, \cdots, p_{jn}] \right] \right\}_{j=1}^{m}$, where $v_j$ is the name of the $j$-the concept (*e.g.*, "horse"), $m$ concepts are stored totally. For the concept $v_j$, $f_{ji}$ and $p_{ji}$ are the feature and learned prompt using the $i$-th instance, respectively, and $n$ instances are learned for $v_j$.

In summary, a tool is formulated as $T_{\theta,\mathcal{P}}$, where $\theta$ is the parameter of neural networks. The forward process of a sample $x$ is $T_{\theta,\mathcal{P}}(x) = \theta([x, \mathcal{P}])$, where we concatenate $x$ with a retrieved prompt from $\mathcal{P}$ as the input for the tool. We update the tool $T_{\theta,\mathcal{P}}$ by tuning $\mathcal{P}$ while $\theta$ being fixed,

$$\min_{\mathcal{P}} \mathbb{E}_{(x,y)} \mathcal{L}\big(T_{\theta,\mathcal{P}}(x), y\big),$$

where $(x, y)$ is collected data, $\mathcal{L}$ is the loss function of $T_{\theta,\mathcal{P}}$.

**Prompt ensemble.** During inference, we use prompt ensemble to retrieve and utilize prompts from the learned prompt pool $\mathcal{P}$. Specifically, given a generated program, we first identify the visual concept for each involved tool in the program. For example, given an image editing task, "*Replace the dog with a cat*" where the SEG, SELECT, and REPLACE tools will be used, "*dog*" is the visual concept for the SEG and SELECT tools, and "*cat*" is the visual concept for the REPLACE tool. Then, in each step of tool usage with an input image $x$, if the visual concept is not in $\mathcal{P}$ of the corresponding tool, the prompt $p'$ for $x$ will be set as a zero-vector, *i.e.* the concept has not been learned for the tool so we just use the original tool (using zero-vector as a prompt); if the visual concept can be found in $\mathcal{P}$, *i.e.* the tool was updated with the visual concept before, we aggregate the prompts corresponding to this concept based on the similarity between features stored with these prompts and the feature extracted from the current input $x$: $p'$ is computed by $p' = \frac{\sum_{i=1}^{n} w_i \cdot p_{ji}}{\sum_{i=1}^{n} w_i}$, where we compute the cosine similarity between feature $f_x$ of $x$ and features $f_{ji}$ in $\mathcal{P}$ as the weight $w_i$.

### 3.4.2 Updating LLMs with demonstrations

Besides the visual tools, CLOVA can also update its LLMs. As we mentioned above, CLOVA utilizes a demonstration pool $\mathcal{D}$ to provide relevant examples for the LLMs. After working on new data, the plan, program, and reflection will be stored in $\mathcal{D}$ as correct or incorrect examples, based on whether the data is correctly proceeded. We also have a validation process that uses the original instruction as validation data to evaluate stored in-context examples. As the size of $\mathcal{D}$ grows, LLMs use more examples and therefore strengthen reasoning skills.

## 4. Experiments

### 4.1. Setting

Following VISPROG [11], we evaluate our method on four tasks: compositional VQA, multiple-image reasoning, language-guided image editing, and factual knowledge tagging, which requires visual perception, compositional reasoning, and image generation and manipulation abilities.

To comprehensively evaluate the learning capability of CLOVA, we set a separate training stage before deployment, which iteratively learns new knowledge via inference, reflection, and updating phases. In the test stage (*i.e.*, after development), we do not update LLMs and visual tools, and evaluate the performance only via the inference phase.

In the compositional VQA task, the GQA dataset [20] is used. We randomly select 500 samples from its train split as the training data, and 500 samples from its test-dev split as the test data. We report the top-1 accuracy. In the multiple-image reasoning task, we use the NLVRv2 dataset [66] that provides two images and a statement. We need to judge whether the statement is true or false. Similarly, we randomly select 500 samples from its train split as our training data, and 500 samples from its dev split as our test data.

Similar to VISPROG, we manually collect data for the language-guided image editing and factual knowledge tagging tasks. To better evaluate the learning capability, we collect fine-grained visual concepts that visual tools may not have learned, such as "*Replace the lion in the image with pine grosbeak*", where pine grosbeak is a fine-grained bird of Passeriformes. In the image editing task, we collect 129 images with 193 instructions, where 27 images with 78 instructions are used for training, and the rest are test data. We manually check whether edited images are semantically correct. The factual knowledge tagging task needs to identify persons or objects with bounding boxes in images. We collect 86 images with 177 instructions for this task, where 10 images with 88 instructions are used for training and the rest are used as the test data. We report the F1 score for this task.

In plan and program generation, prompts contain 4 correct examples and 4 incorrect examples. The demonstration

| | Method | GQA | NLVRv2 | Editing | Tagging |
|---|---|---|---|---|---|
| E2E | Otter [27] | 48.2 | 48.2 | - | - |
| | MMICL [83] | 64.4 | 62.2 | - | - |
| Tool | GPT4TOOLs [75] | 41.2 | 45.4 | 17.8 | |
| | Visual ChatGPT [75] | 43.2 | 51.6 | 21.7 | - |
| | InternGPT [40] | 44.8 | 39.4 | - | - |
| | HuggingGPT [62] | 46.0 | 44.0 | - | - |
| | ViperGPT [67] | 47.2 | - | - | - |
| | VISPROG [11] | 49.8 | 60.8 | 40.2 | 0.393 |
| | CLOVA (Ours) | **54.6** | **65.6** | **65.4** | **0.502** |

Table 3. Comparisons in the four tasks. We report accuracies on GQA, NLVRv2, and image editing tasks, and F1 score on the knowledge tagging task. 'E2E' means end-to-end methods.

pool $\mathcal{D}$ is initialized having about 20 in-context examples.

### 4.2. Main Results

We compare CLOVA with tool-usage methods: VISPROG [11], GPT4TOOLs [80], Visual ChatGPT [75], InternGPT [40], HuggingGPT [62], and ViperGPT [67]. We use their official codes, where all methods use the GPT-3.5 model. In addition, we also compare CLOVA with two advanced end-to-end models: Otter [27] and MMICL [83], which do well in GQA and NLVRv2 datasets. Results on the four tasks are shown in Tab. 3. We observe that CLOVA achieves the best performance among tool-usage methods. CLOVA has at least $4.8\%$ improvements on the GQA dataset and $4.9\%$ improvements on the NLVRv2 dataset. The reason is that CLOVA learns how to generate programs for the two tasks and update the VQA and LOC tools for better image perception. CLOVA performs competitively and even outperforms Otter and MMICL.

On image editing and knowledge tagging tasks, we do not compare our method with InterGPT, HuggingGPT, and ViperGPT, since they either need object masks or cannot accurately locate objects. In addition, most methods cannot finish the tagging task. Thus, we compare our method with VISPROG and Visual ChatGPT. As we collect find-grained data, it is challenging for off-the-shelf classification, segmentation, and image generation tools. Since GPT4TOOLs and Visual ChatGPT cannot use OpenCV functions and do not have the learning capability, they get bad performance on the image editing task. VISPROG can use OpenCV functions, but it cannot learn new knowledge. Its main fault is the inability to recognize or generate fine-grained concepts. Compared with them, the learning capability of CLOVA brings more than $20\%$ and $10\%$ improvements to image editing and knowledge tagging tasks, respectively.

### 4.3. Qualitative Results

In Fig. 6, we visualize four cases to illustrate the reflection and learning capability in CLOVA. It identifies tools that need to be updated, no matter LLMs or visual tools. Prompt engineering guides LLMs to generate correct programs for similar instructions. Visual tools learn new concepts via our

Figure 6. Case study of CLOVA on four example tasks.

data collection and prompt turning schemes. In Fig. 7, we visualize an example of global and local reflection. When the instruction is complex, the global reflection does not accurately identify which step has the error. Using global reflection as in-context examples still cannot generate correct programs. In contrast, local reflection successfully identifies the error step, and using local reflection generates correct programs, showing the effectiveness of local reflection.

## 4.4. Ablation Studies

We conduct ablation studies on the reflection and learning phases, using the GQA and NLVRv2 datasets. For reflection, we evaluate only using global reflection, only using local reflection, not using multimodal intermediate results, and not generating plans. We separately evaluate learning schemes for LLMs and visual tools. We evaluate only storing correct examples or incorrect examples for updating

Figure 7. Visualization of the global reflection and local reflection in a VQA task.

| | Method | GQA | NLVRv2 |
|---|---|---|---|
| | w/o local reflection | 52.0 | 65.2 |
| | w/o global reflection | 53.6 | 64.2 |
| Reflection | w/o intermediate results | 48.8 | 61.2 |
| | w/o plan | 50.0 | 62.6 |
| | Ours | **54.6** | **65.6** |
| | w/o incorrect cases | 46.1 | 61.4 |
| Prompt Engineering | w/o correct cases | 48.2 | 63.2 |
| for LLMs | w/o validation | 44.2 | 61.0 |
| | Ours | **54.6** | **65.6** |
| Prompt Tuning | w/o validation | 42.8 | 62.8 |
| for visual tools | Ours | **54.6** | **65.6** |

Table 4. Ablation on the GQA and NLVRv2 dataset.

| Dataset | Method | LLaMA2-7B | GPT-3.5 | GPT-4 |
|---|---|---|---|---|
| | Baseline | 39.2 | 46.4 | 52.6 |
| GQA | + Update LLMs | 56.8 | 51.6 | 56.6 |
| | + Update visual tools | 60.2 | 54.6 | 60.4 |
| | Baseline | 50.0 | 60.2 | 64.8 |
| NLVRv2 | + Update LLMs | 59.2 | 63.6 | 68.8 |
| | + Update visual tools | 63.8 | 65.6 | 69.2 |

Table 5. Different LLMs on the GQA and NLVRv2 datasets.

| Method | GQA | NLVRv2 | Editing | Tagging |
|---|---|---|---|---|
| LLama2-7B | 39.2 | 50.0 | 31.2 | 0.308 |
| LLama2-7B + Ours | **60.2** | **63.8** | **47.6** | **0.357** |
| Mistral-7B | 20.4 | 34.6 | 29.0 | 0.205 |
| Mistral-7B + Ours | **31.4** | **42.2** | 46.5 | **0.303** |

Table 6. Results on two open-source LLMs.

LLMs. We also evaluate removing the validation process in prompt engineering and prompt tuning processes. Results are shown in Tab. 4. We find that these components are necessary for CLOVA to achieve better performance.

We evaluate CLOVA using different LLMs: LLaMA2-7B, GPT-3.5-turbo, and GPT-4. Results on GQA and NLVRv2 are shown in Tab. 5. We find that CLOVA leads to improvements in both strong LLMs (GPT-4) and weaker LLMs (GPT-3.5 and LLaMA2-7). We observe that CLOVA even achieves higher improvements on open-source LLMs (*i.e.*, LLaMA2-7B), 21% on the GQA dataset and 13.8% on the NLVRv2 dataset, bringing the significance of studying the learning capability of visual assistants. We further conduct experiments to evaluate CLOVA on two open-source LLMs: LLaMA2-7B and Mistral-7B. Results are shown in Tab. 6. We observe that CLOVA achieves significant improvements again.

## 5. Conclusion and Future Work

In this paper, we have presented CLOVA, a general visual assistant that can adapt to new environments via inference, reflection, and learning in a closed-loop framework. In the inference phase, using both correct and incorrect examples for prompts benefits to generate better plans and programs. Our reflection scheme is capable of identifying tools that need to be updated. Through three data collection manners and the validation-learning prompt tuning scheme in the learning phase, CLOVA can efficiently improve its tools. Experimental results on four tasks and different LLMs show the effectiveness of CLOVA as a general visual assistant with learning abilities.

In the current method, we assume there is no selection or loop structure in programs, and assume there is at most one tool that needs updates in a task. The two assumptions cannot always hold in the real world. We could add selection and loop in-context examples for program generation and iterate the reflection and learning phases to update multiple tools. Besides, we could make some deployment designs to save the response time in our loop, including a foreground process and a background process. The former performs inference and gathers human feedback, while the latter does reflection, updates tools, and periodically synchronizes the weights of tools.

The framework of CLOVA can be easily generalized to new tools. (1) We will try multimodal LLMs (*e.g.*, LLaVA-1.5 [37]) to replace LLMs. In this case, we will evaluate the effectiveness of visual information for plan and program generation, reflection, and answer inference. (2) We can add more up-to-date tools (*e.g.*, BLIP2 [30] as a VQA tool and SAM [26] as an SEG tool), by just designing and programming their forward and prompt tuning processes.

# CLOVA: A <u>C</u>losed-<u>LO</u>op <u>V</u>isual <u>A</u>ssistant with Tool Usage and Update

## Supplementary Material

## 6. Framework of CLOVA

The pseudo-code of CLOVA is summarized in Algorithm 1.

---

**Algorithm 1** CLOVA

---

**Input:** LLMs, visual tools, instruction data $\mathcal{T} = \{T_1, T_2, \cdots, T_t\}$, demonstration pool $\mathcal{D}$, prompt pool $\mathcal{P} = \emptyset$.
**Output:** Updated $\mathcal{D}$, updated $\mathcal{P}$
1: **for** $i = 1, 2, \ldots, t$ **do**
2:    Perform inference for $T_i$ by generating plan and program.
3:    **if** $T_i$ is correctly solved **then**
4:       Save the plan and program to $\mathcal{D}$.
5:    **else**
6:       Convert intermediate results into language.
7:       Perform global reflection.
8:       Perform training-validation prompt tuning, store in-context examples into $\mathcal{D}$ and prompts in $\mathcal{P}$.
9:       **if** Updated tools solve $T_i$ incorrectly **then**
10:          Perform local reflection.
11:          Perform training-validation prompt tuning, store in-context examples into $\mathcal{D}$ and prompts in $\mathcal{P}$.
12:       **end if**
13:    **end if**
14: **end for**

---

## 7. Comparisons with related methods

In Tab. 7, we present a more detailed comparison table with more related methods. In the table, 'Global' means the global reflection, 'Local' means the local reflection, 'Instruction' means instruction-following tuning, 'RL' means reinforcement learning, and 'Prompt' means using prompt examples as in-context learning. We observe that many tool-usage methods do not have a reflection capability, and few methods use global reflection to improve the plans or programs generated by LLMs. Different from them, CLOVA uses both global reflection and local reflection to identify tools that need to be updated, capable of handling complex instructions. Moreover, as we all know, our method is the first work to update visual tools, through which the visual assistant can better adapt to new environments.

## 8. Prompt Examples

### 8.1. In-context examples

In the inference phase of our method, CLOVA generates plans and programs based on in-context examples that in-clude correct examples and incorrect examples with criteria. Here we show some correct examples and incorrect examples for the compositional VQA, multi-image reasoning, image editing, and factual knowledge tagging tasks, as shown in Figs. 8 to 11.

### 8.2. Prompts in inference

In the inference phase, we use LLMs to generate plans and programs. We show examples of prompts for plan generation and program generation in Figs. 12 and 13, respectively.

### 8.3. Prompts in reflection

In the reflection phase, we use LLMs for global reflection and local reflection. We show two examples of prompts for global reflection and local reflection in Figs. 14 and 15, respectively.

### 8.4. Prompts in learning

In the learning phase, we use LLMs to infer answers for the VQA tool, and then tune prompts of the VQA tool using the question and inferred answer. One example of prompts sent to LLMs for answer inferring is shown in Fig. 16, respectively.

## 9. Details of Tool update

### 9.1. Update VQA tool

#### 9.1.1 Model

We use the BLIP [29] model for the VQA tool. One BLIP model is composed of three components: an image encoder, an image-grounded question encoder, and an answer coder, which are used to extract image features, extract question features, and generate answers, respectively. In our prompt tuning scheme, the visual features are extracted from the image encoder followed by average pooling, whose dimension is 768. We concatenate learnable prompts with the inputs to the answer decoder, guiding the answer decoder to generate desirable answers. The architecture of BLIP is shown in Fig. 17.

#### 9.1.2 Training

The size of learnable prompts is $\mathbb{R}^{16 \times 768}$, that is we learn 16 vectors as prompts, and the dimension of each one is 768. Since it is non-trivial to define visual concepts for the VQA tool, we do not store concepts in the prompt pool of

(a) Correct examples for plan generation

**Instruction**: Is the vehicle in the top of the image?
**Plan**:
Step1, Locate the upper region of the given image, and obtain bounding boxes of the upper region.
Step2, Crop the upper region from the given image, based on bounding boxes of the upper region. The bounding boxes are obtained in Step1.
Step3, Locate vehicle in the upper region of the given image, and obtain bounding boxes of vehicle. The upper region is cropped in Step2.
Step4, Count the number of vehicle, based on bounding boxes of vehicle. The bounding boxes are obtained in Step3.
Step5, Obtain whether the answer is 'yes' or 'no', based on the number of vehicles. The number is obtained in Step4. If the number is greater than zero, the answer is 'yes'; On the contrary, the answer is 'no'.
Step6, Visualize results.

(b) Incorrect examples for plan generation

**Instruction**: On which side of the picture are the yellow chairs, the right or the left?
**Plan**:
Step1, Locate the yellow chairs in the given image and obtain bounding boxes of the yellow chairs.
Step2, Crop the region where the yellow chairs are located from the given image, based on the bounding boxes of the yellow chairs. The bounding boxes are obtained in Step1.
Step3, Obtain whether the yellow chairs are on the right or left side of the picture, based on the position of the cropped region. If the region is on the right side of the picture, the answer is 'right'; On the contrary, the answer is 'left'.
Step4, Visualize the results.
**Reason**: The Step1-Step3 in the plan are correct. Using the position of the cropped region is not supported. It should first crop the left region of the image, and then detect where there are yellow chairs. If not, the answer is right.

(c) Correct examples for program generation

**Instruction** : Is the vehicle in the top of the image?
**Plan**:
Step1, Locate the upper region of the given image, and obtain bounding boxes of the upper region.
Step2, Crop the upper region from the given image, based on bounding boxes of the upper region. The bounding boxes are obtained in Step1.
Step3, Locate vehicle in the upper region of the given image, and obtain bounding boxes of vehicle. The upper region is cropped in Step2.
Step4, Count the number of vehicle, based on bounding boxes of vehicle. The bounding boxes are obtained in Step3.
Step5, Obtain whether the answer is 'yes' or 'no', based on the number of vehicles. The number is obtained in Step4. If the number is greater than zero, the answer is 'yes'; On the contrary, the answer is 'no'.
Step6, Visualize results.
**Program**:
BOX0=LOC(image=IMAGE,object='TOP')
IMAGE0=CROP(image=IMAGE,box=BOX0)
BOX1=LOC(image=IMAGE0,object='vehicle')
ANSWER0=COUNT(box=BOX1)
ANSWER1=EVAL(expr="'yes' if {ANSWER0} > 0 else 'no'")
FINAL_RESULT=RESULT(var=ANSWER1)

(d) Incorrect examples for program generation

**Instruction** : Is the lamp different in color than the shirt?
**Plan**:
Step1, Locate the lamp in the given image, and obtain bounding boxes of lamp.
Step2, Crop the region of the lamp from the given image, based on bounding boxes of lamp. The bounding boxes are obtained in Step1.
Step3, Asking the image region of lamp, 'What color is the lamp?'. The image region of lamp is cropped in Step2.
Step4, Locate the shirt in the given image, and obtain bounding boxes of shirt.
Step5, Crop the region of the shirt from the given image, based on bounding boxes of shirt. The bounding boxes are obtained in Step4.
Step6, Asking the image region of shirt, 'What color is the shirt?'. The image region of lamp is cropped in Step5.
Step7, Obtain whether the answer is 'yes' or 'no', based on the color of lamp and the color of shirt. The color of lamp and shirt is obtained in Step3 and Step6, respectively. If their color are the same, the answer is 'yes'; On the contrary, the answer is 'no'.
Step8, Visualize results.
**Program**:
BOX0=LOC(image=IMAGE,object='lamp')
IMAGE0=CROP(image=IMAGE,box=BOX0)
BOX1=LOC(image=IMAGE0,object='shirt')
ANSWER0=COUNT(box=BOX1)
ANSWER1=EVAL(expr="'yes' if {ANSWER0} > 0 else 'no'")
FINAL_RESULT=RESULT(var=ANSWER1)
**Reason:** The plan is correct, and can address the task. But the Step3-Step6 of the program does not match the plan. The plan locates, crops, and asks color of the lamp and shirt, but the program counts the number of lamp.

Figure 8. Demonstration examples for compositional VQA tasks in $\mathcal{D}_{p,s}$, $\mathcal{D}_{p,f}$, $\mathcal{D}_{c,s}$, and $\mathcal{D}_{c,f}$.

| Method | Visual Tool | Reflection | Update LLMs | Update Tools |
|---|---|---|---|---|
| ART [50] | ✗ | ✗ | Prompt | - |
| LATM [4] | ✗ | Global | Prompt | - |
| TRICE [55] | ✗ | Global | Instruction + RL | - |
| ToolkenGPT [13] | ✗ | - | ✗ | - |
| Toolformer [60] | ✗ | - | Fine-tune | - |
| VISPROG [11] | ✓ | ✗ | ✗ | ✗ |
| Visual ChatGPT [75] | ✓ | ✗ | ✗ | ✗ |
| InternGPT [40] | ✓ | ✗ | ✗ | ✗ |
| HuggingGPT [62] | ✓ | ✗ | ✗ | ✗ |
| ViperGPT [67] | ✓ | ✗ | ✗ | ✗ |
| ToT [16] | ✓ | ✗ | ✗ | ✗ |
| Chameleon [42] | ✓ | ✗ | ✗ | ✗ |
| ControlLLM [41] | ✓ | ✗ | ✗ | ✗ |
| MM-REACT [81] | ✓ | ✗ | ✗ | ✗ |
| VideoAgent [7] | ✓ | ✗ | ✗ | ✗ |
| Llava-plus [39] | ✓ | ✗ | Instruction | ✗ |
| Gorilla [52] | ✓ | ✗ | Instruction | ✗ |
| GPT4TOOLs [80] | ✓ | ✗ | Instruction | ✗ |
| MLLM-Tool [72] | ✓ | ✗ | Instruction | ✗ |
| VIoTGPT [84] | ✓ | ✗ | Instruction | ✗ |
| OpenAGI [9] | ✓ | ✗ | Reinforcement Learning | ✗ |
| AssistGPT [8] | ✓ | Global | Prompt | ✗ |
| **CLOVA (Ours)** | ✓ | Global+Local | Prompt | Prompt |

Table 7. Comparisons with representative tool-usage methods.

the VQA tool, and all learned prompts and stored together. We use questions and inferred answers of incorrect cases (detailed in Section 3.4) to update the VQA tool. We also store correct cases with zero vectors as prompts. We use the language modeling loss [29] to train learnable prompts, where the Adam optimizer is used and the learning rate is $1e-3$. We train the prompts 100 steps for each instance.

### 9.1.3 Inference

The prompt ensemble process of the VQA tool has two steps. (1) We roughly select out 20 prompts from the prompt pool as candidates, by computing the similarity between the given query instance and stored instances in the prompt pool. (2) We use prompt ensemble (detailed in Section 3.4) to aggregate the 20 prompts for a query instance. In other words, we do not aggregate all stored prompts for a query instance, but 20 similar instances.

### 9.2. Update LOC tool

#### 9.2.1 Model

In CLOVA, we use the OWL-ViT model [25] for object localization as the LOC tool. The architecture of OWL-ViT is shown in Fig. 18. OWL-ViT model uses a standard vision

transformer as the image encoder and a similar transformer architecture as the text encoder. It removes the token pooling and final projection layer, and instead linearly projects each output token representation to obtain per-object image embeddings for classification. Besides, box coordinates are obtained by passing token representations through a small MLP. The text embeddings, which are called queries, are obtained by passing category names or other textual object descriptions through the text encoder. At inference time, given a set of candidate class names and an image, the model predicts a bounding box and a probability with which each query, and filters out the bounding box with the prediction confidence less than 0.1.

#### 9.2.2 Training

In this study, upon identifying the need to update the LOC tool for learning a specific concept, we employ instance-wise prompts to update the OWL-ViT model. To achieve this, CLOVA first collect training data from open-vocabulary datasets. We adopt LVIS dataset [10] for OWL-ViT model. Taking the concept "glove" as an example, CLOVA randomly selects 200 samples from the LVIS dataset, whose class labels contain a glove as the training data. During training, visual features are extracted from

**(a) Correct examples for plan generation**

**Instruction**: An image shows one bare hand with the thumb on the right holding up a belly-first, head-up crab, with water in the background.

**Plan**:

Step1, Asking the left image, 'Does the image shows one bare hand with the thumb on the right holding a crab?'.

Step2, Asking the right image, 'Does the image shows one bare hand with the thumb on the right holding a crab?'.

Step3, Asking the left image, 'Is the crab belly-first and head-ups?'.

Step4, Asking the right image, 'Is the crab belly-first and head-ups?'.

Step5, Asking the left image, 'Is there water in the background?'.

Step6, Asking the right image, 'Is there water in the background?'.

Step7, Obtain whether the answer is 'yes' or 'no' by executing Python expression, based on the answers obtained in Step1, Step3, and Step5. If all answers obtained in Step1, Step3, and Step5 are 'yes', the answer in this step is 'yes'; On the contrary, the answer is 'no'.

Step8, Obtain whether the answer is 'yes' or 'no' by executing Python expression, based on the answers obtained in Step2, Step4, and Step6. If all answers obtained in Step2, Step4, and Step6 are 'yes', the answer in this step is 'yes'; On the contrary, the answer is 'no'.

Step9, Obtain whether the answer is 'yes' or 'no' by executing Python expression, based on the answers obtained in Step7 and Step8. If one of the two answers is equal to 'yes' and the rest one answer is equal to 'no', the answer in this step is 'yes'; On the contrary, the answer is 'no'.

Step10, Visualize results.

**(b) Incorrect examples for plan generation**

**Instruction**: There are two pairs of hands wearing gloves.

**Plan**:

Step1, Asking the left image, 'How many pairs of hands are in the image?'.

Step2, Asking the right image, 'How many pairs of hands are in the image?'.

Step3, Asking the left image, 'Are the hands wearing gloves?'.

Step4, Asking the right image, 'Are the hands wearing gloves?'.

Step5, Obtain whether the answer is 'yes' or 'no' by executing Python expression, based on the number obtained in Step1 and answers obtained in Step3. If the number is equal to '2' and the answer is equal to 'yes', the answer in this step is 'yes'; On the contrary, the answer is 'no'.

Step6, Obtain whether the answer is 'yes' or 'no' by executing Python expression, based on the number obtained in Step2 and answers obtained in Step4. If the number is equal to '2' and the answer is equal to 'yes', the answer in this step is 'yes'; On the contrary, the answer is 'no'.

Step7, Obtain whether the answer is 'yes' or 'no' by executing Python expression, based on the answers obtained in Step5 and Step6. If one of the two answers is equal to 'yes' and the rest one answer is equal to 'no', the answer in this step is 'yes'; On the contrary, the answer is 'no'.

Step8, Visualize results.

**Reason**: In Step5 of the plan, the plan judge whether there are two pairs of hands in the left image. It is wrong, the plan should judge whether the sum of pairs of hands in the two images are equal to two, and then judge whether the two pairs of hands are wearing gloves.

**(c) Correct examples for program generation**

**Instruction**: An image shows one bare hand with the thumb on the right holding up a belly-first, head-up crab, with water in the background.

**Plan**:

Step1, Asking the left image, 'Does the image shows one bare hand with the thumb on the right holding a crab?'.

Step2, Asking the right image, 'Does the image shows one bare hand with the thumb on the right holding a crab?'.

Step3, Asking the left image, 'Is the crab belly-first and head-ups?'.

Step4, Asking the right image, 'Is the crab belly-first and head-ups?'.

Step5, Asking the left image, 'Is there water in the background?'.

Step6, Asking the right image, 'Is there water in the background?'.

Step7, Obtain whether the answer is 'yes' or 'no' by executing Python expression, based on the answers obtained in Step1, Step3, and Step5. If all answers obtained in Step1, Step3, and Step5 are 'yes', the answer in this step is 'yes'; On the contrary, the answer is 'no'.

Step8, Obtain whether the answer is 'yes' or 'no' by executing Python expression, based on the answers obtained in Step2, Step4, and Step6. If all answers obtained in Step2, Step4, and Step6 are 'yes', the answer in this step is 'yes'; On the contrary, the answer is 'no'.

Step9, Obtain whether the answer is 'yes' or 'no' by executing Python expression, based on the answers obtained in Step7 and Step8. If one of the two answers is equal to 'yes' and the rest one answer is equal to 'no', the answer in this step is 'yes'; On the contrary, the answer is 'no'.

Step10, Visualize results.

**Program**:

```
ANSWER0=VQA(image=LEFT,question='Does the image shows one bare hand with the thumb on the right holding a crab?')
ANSWER1=VQA(image=RIGHT,question='Does the image shows one bare hand with the thumb on the right holding a crab?')
ANSWER2=VQA(image=LEFT,question='Is the crab belly-first and head-ups?')
ANSWER3=VQA(image=RIGHT,question='Is the crab belly-first and head-ups?')
ANSWER4=VQA(image=LEFT,question='Is there water in the background?')
ANSWER5=VQA(image=RIGHT,question='Is there water in the background?')
ANSWER6=EVAL(expr='{ANSWER0} and {ANSWER2} and {ANSWER4}')
ANSWER7=EVAL(expr='{ANSWER1} and {ANSWER3} and {ANSWER5}')
ANSWER8=EVAL(expr='{ANSWER6} xor {ANSWER7}')
FINAL_ANSWER=RESULT(var=ANSWER8)
```

**(d) Incorrect examples for program generation**

**Instruction**: A mitten is being worn in one image and the mittens are not being worn in the other image.

**Plan**:

Step1, Asking the left image, 'Is a mitten being worn in the image?'.

Step2, Asking the right image, 'Is a mitten being worn in the image?'.

Step3, Obtain whether the answer is 'yes' or 'no' by executing Python expression, based on the answers obtained in Step1 and Step2. If one of the two answers is equal to 'yes' and the rest one answer is equal to 'no', the answer in this step is 'yes'; On the contrary, the answer is 'no'.

Step4, Visualize results.

**Program**:

```
ANSWER0=VQA(image=LEFT,question='Is a mitten being worn in the image?')
ANSWER1=VQA(image=RIGHT,question='Is a mitten being worn in the image?')
ANSWER2=EVAL(expr='{ANSWER0} and {ANSWER1}')
FINAL_ANSWER=RESULT(var=ANSWER2)
```

**Reason**: The plan is correct, and can address the given question. But the Step3 of the program does not match the plan. The program should use 'xor' instead of 'and'.

Figure 9. Demonstration examples for multi-image reasoning tasks in $\mathcal{D}_{p,s}$, $\mathcal{D}_{p,f}$, $\mathcal{D}_{c,s}$, and $\mathcal{D}_{c,f}$.

(a) Correct examples for plan generation

**Instruction**: Hide the face of Nicole Kidman with face_with_tongue.
**Plan**:
Step1, Detect face regions from the given image, and obtain bounding boxes of face regions.
Step2, Select the face region of Nicole Kidman, based on the bounding boxes obtained in Step1.
Step3, Add the emoji face_with_tongue to the face region of Nicole Kidman in the given image, where the face region of Nicole Kidman is obtained in Step2.
Step4, Visualize results.

(b) Incorrect examples for plan generation

**Instruction**: Hide the face of Nicole Kidman with face_with_tongue.
**Plan**:
Step1, Detect face regions from the given image, and obtain bounding boxes of face regions.
Step2, Add the emoji face_with_tongue to the face regions the given image, where the face region is obtained in Step1.
Step4, Visualize results.
**Reason**: The plan is incorrect. In Step2, the program should use the SELECT model to classify face regions of Nicole Kidman, instead of adding emoji to all faces.

(c) Correct examples for program generation

**Instruction**: Hide the face of Nicole Kidman with face_with_tongue.
**Plan**:
Step1, Detect face regions from the given image, and obtain bounding boxes of face regions.
Step2, Select the face region of Nicole Kidman, based on the bounding boxes obtained in Step1.
Step3, Add the emoji face_with_tongue to the face region of Nicole Kidman in the given image, where the face region of Nicole Kidman is obtained in Step2.
Step4, Visualize results.
Program:
OBJ0=FACEDET(image=IMAGE)
OBJ1=SELECT(image=IMAGE,object=OBJ0,query='Nicole Kidman',category=None)
IMAGE0=EMOJI(image=IMAGE,object=OBJ1,emoji='face_with_tongue')
FINAL_RESULT=RESULT(var=IMAGE0)

(d) Incorrect examples for program generation

**Instruction**: Hide the face of Nicole Kidman with face_with_tongue.
**Plan**:
Step1, Detect face regions from the given image, and obtain bounding boxes of face regions.
Step2, Select the face region of Nicole Kidman, based on the bounding boxes obtained in Step1.
Step3, Add the emoji face_with_tongue to the face region of Nicole Kidman in the given image, where the face region of Nicole Kidman is obtained in Step2.
Step4, Visualize results.
Program:
OBJ0=FACEDET(image=IMAGE)
OBJ1=SELECT(image=IMAGE,object=OBJ0,query='Nicole Kidman',category=None)
IMAGE0=REPLACE(image=IMAGE,object=OBJ1,emoji='face_with_tongue')
FINAL_RESULT=RESULT(var=IMAGE0)
**Reason**: The plan is correct, and can address the given question. But the Step3 of the program does not match the plan. The program should use EMOJI functions instead of REPLACE.

Figure 10. Demonstration examples for image editing tasks in $\mathcal{D}_{p,s}$, $\mathcal{D}_{p,f}$, $\mathcal{D}_{c,s}$, and $\mathcal{D}_{c,f}$.

the backbone, and CLOVA concatenates learnable prompts with the inputs to the vision transformer decoder. The model is trained with original losses introduced by OWL-ViT where only the prompts are learned. The losses include classification loss and bounding box regression loss. The former uses focal sigmoid cross-entropy [90] while the latter uses $L_1$ loss. For classification loss, we regard the learned concept and one randomly selected class name in the image as positive labels and randomly select 13 class names as negative labels per image to avoid overfitting. The number of learned prompts is 100, and the prompt is randomly initialized. We set the maximum training step as 100. If within 100 steps, the sample could be detected cor-rectly by the model, CLOVA would save this sample and the prompt, otherwise, CLOVA will remove this sample and the corresponding prompt. Only if all the positive labels are correctly classified and the average IOU between the prediction boxes and ground-truth boxes is larger than $0.6$, we recognize this sample is correctly detected. We use Adam as the optimizer and set the learning rate as $5e - 4$. Eventually, we save the feature and the learned prompt of each correctly detected instance for this concept.

### 9.2.3 Inference

The prompt ensemble process of the LOC tool also has two steps. (1) Given a visual concept and a query image, we se-

**Instruction**: Tag the wild animals that lives on the land
**Plan**:
Step1, Locate wild animals from the given image, and obtain bounding boxes of wild animals.
Step2, List 'wild animals that lives on the land' by asking GPT.
Step3: Classify wild animals that lives on the land, based on the bounding boxes of wild animals obtained in Step1 and the wild animal list obtain in Step2.
Step4: Tag the bounding boxes and labels to the classified wild animals.
Step5: Visualize results.

(b) Incorrect examples for plan generation

**Instruction**: Tag two utensils used for drinking, which are usually in the kitchen
**Plan**:
Step1, Locate utensil from the given image, and obtain bounding boxes of utensil.
Step2, List two 'utensils used for drinking' by asking GPT.
Step3: Classify utensil, based on the bounding boxes of utensil obtained in Step1 and the utensil list obtain in Step2.
Step4: Tag the bounding boxes and labels to the classified utensil.
Step5: Visualize results.
**Reason**: In Step2, it should list 'utensils used for drinking, which are usually in the kitchen', instead of only 'utensils used for drinking'.

(c) Correct examples for program generation

**Instruction**: Tag the wild animals that lives on the land
**Plan**:
Step1, Locate wild animals from the given image, and obtain bounding boxes of wild animals.
Step2, List 'wild animals that lives on the land' by asking GPT.
Step3: Classify wild animals that lives on the land, based on the bounding boxes of wild animal obtained in Step1 and the wild animals list obtain in Step2.
Step4: Tag the bounding boxes and labels to the classified wild animals.
Step5: Visualize results.
**Program**:
OBJ0=LOC(image=IMAGE,object='wild animal')
LIST0=LIST(query='wild animals that lives on the land',max=20)
OBJ1=CLASSIFY(image=IMAGE,object=OBJ0,categories=LIST0)
IMAGE0=TAG(image=IMAGE,object=OBJ1)
FINAL_RESULT=RESULT(var=IMAGE0)

(d) Incorrect examples for program generation

**Instruction**: Tag the wild animals that lives on the land
**Plan**:
Step1, Locate wild animals from the given image, and obtain bounding boxes of wild animals.
Step2, List 'wild animals that lives on the land' by asking GPT.
Step3: Classify wild animals that lives on the land, based on the bounding boxes of wild animal obtained in Step1 and the wild animals list obtain in Step2.
Step4: Tag the bounding boxes and labels to the classified wild animals.
Step5: Visualize results.
**Program**:
OBJ0=LOC(image=IMAGE,object='wild animal')
LIST0=LIST(query='wild animals that lives on the land',max=20)
OBJ1=LOC(image=IMAGE,object=land)
IMAGE0=TAG(image=IMAGE,object=OBJ1)
FINAL_RESULT=RESULT(var=IMAGE0)
**Reason:** The plan is correct, but the program does not match the plan. In Step3 in the program, it is wrong to use the LOC model. It should use the CLASSIFY for OBJ0 obtained in Step1.

Figure 11. Demonstration examples for factual knowledge taggings task in $\mathcal{D}_{p,s}$, $\mathcal{D}_{p,f}$, $\mathcal{D}_{c,s}$, and $\mathcal{D}_{c,f}$.

lect all the prompts having the same visual concept from the prompt pool as candidates. We then filter out the prompts by making predictions with each candidate prompt, if the prediction confidence is larger than 0.1, we will use the prompt for the query image, otherwise, we will remove this candidate prompt. (2) We use prompt ensemble (detailed in Section 3.4) to aggregate all the selected prompts for a query instance and contact the prompt with the input to produce a prediction.

## 9.3. Update SEG tool

### 9.3.1 Model

We use the Maskformer [6] model for the SEG tool. One Maskformer model is composed of three components: a backbone, a pixel decoder, and a transformer decoder. The backbone extracts features of images. Then, the pixel decoder gradually upsamples image features to extract per-pixel embeddings. Finally, the transformer decoder uses

You are a **planner**. Given a question, you need to generate the plan.

Some correct cases are as follows.
**Instruction :** What color is the curtain that is to the right of the mirror?
**Plan:**
Step1, Locate mirror in the given image, and obtain bounding boxes of mirror.
Step2, Crop the region on the right side of the mirror from the given image, based on the bounding boxes of mirror. The bounding boxes are obtained in Step1.
Step3, Asking the image region on the right side of the mirror, 'What color is the curtain?'. The image region is cropped in Step2.
Step4, Visualize results.

Some incorrect cases and their reasons are as follows.
**Instruction** : Who is in the blue water?
**Plan**:
Step1, Locate blue regions in the given image, and obtain bounding boxes of the blue regions.
Step2, Crop the blue regions from the given image, based on the bounding boxes of the blue regions. The bounding boxes are obtained in Step1.
Step3, Locate people in the blue regions of the given image, and obtain bounding boxes of people. The blue regions are cropped in Step2.
Step4, Count the number of people, based on bounding boxes of people. The bounding boxes are obtained in Step3.
Step5, Obtain the names of the people in the blue water, based on the number of people. The number is obtained in Step4. If there is only one person, provide their name; if there are multiple people, provide a list of their names.
Step6, Visualize results.
**Reason**:
The Step4-Step6 in the plan is correct. After obtaining the bounding boxes of people in Step3, the plan should crop the image region of the person, and ask who is the person, instead of counting the number.

Now, you need to generate plan for the following query.
**Instruction** : Is there any snow or grass in this scene?
**Plan:**

Figure 12. Prompts for plan generation.

---

You are a **programmer**. You need to generate the program based on Instruction and Plan.

Some **correct cases** are as follows.
**Instruction:** What color is the curtain that is to the right of the mirror?
**Plan**:
Step1, Locate mirror in the given image, and obtain bounding boxes of mirror.
Step2, Crop the region on the right side of the mirror from the given image, based on the bounding boxes of mirror. The bounding boxes are obtained in Step1.
Step3, Asking the image region on the right side of the mirror, 'What color is the curtain?'. The image region is cropped in Step2.
Step4, Visualize results.
**Program**:
BOX0=LOC(image=IMAGE,object='mirror')
IMAGE0=CROP_RIGHTOF(image=IMAGE,box=BOX0)
ANSWER0=VQA(image=IMAGE0,question='What color is the curtain?')
FINAL_RESULT=RESULT(var=ANSWER0)...

Some **incorrect cases** and their reasons are as follows.
**Instruction** : What is the red piece of clothing in this photograph?
**Plan**:
Step 1: Locate the red piece of clothing in the given photograph and obtain the bounding box of the clothing.
Step 2: Crop the region of the red piece of clothing from the given photograph based on the bounding box obtained in step 1.
Step 3: Asking the image region of the red piece of clothing, "What is the red piece of clothing?"
Step 4: Visualize results.
**Program**:
BOX0=LOC(image= IMAGE,object='red piece of clothing')
IMAGE0=CROP(image= IMAGE,box=BOX0)
ANSWER0=VQA(image=IMAGE1,question='What is the red piece of clothing?')
FINAL_RESULT=RESULT(var=ANSWER0)
**Reason:** The Step3 in the program have bug, because the variable IMAGE1 is not defined. It should be IMAGE0.

Now, you need to generate program for the following query.
**Instruction**: Is the baseball man to the right or to the left of the woman?
**Plan**:
Step1, Locate the baseball man in the given image, and obtain bounding boxes of the baseball man.
Step2, Locate the woman in the given image, and obtain bounding boxes of the woman.
Step3, Obtain the relative position of the baseball man and the woman, based on the bounding boxes of the baseball man and the woman. If the baseball man is to the right of the woman, the answer is 'right'; On the contrary, the answer is 'left'.
Step4, Visualize results.
**Program**:

Figure 13. Prompts for program generation

You are a **debugger**. You need to check which model cause of the wrong answer. Errors may exist in the plan, program, or functions called by the program.

----------------------------------------------------------------------------------

**Instruction**: Is the lamp different in color than the shirt?
**Description of the Input Image**: a photography of a couple of people in a restaurant
**Human Feedback:** The correct answer is yes
**Our Wrong Answer:** no
Following are the decomposed plan, used program, and obtained result in each step.
**Plan:**
Step1, Locate the lamp in the given image, and obtain bounding boxes of lamp.
Step2, Crop the region of the lamp from the given image, based on bounding boxes of lamp. The bounding boxes are obtained in Step1.
Step3, Asking the image region of lamp, 'What color is the lamp?'. The image region of lamp is cropped in Step2.
Step4, Locate the shirt in the given image, and obtain bounding boxes of shirt.
Step5, Crop the region of the shirt from the given image, based on bounding boxes of shirt. The bounding boxes are obtained in Step4.
Step6, Asking the image region of shirt, 'What color is the shirt?'. The image region of lamp is cropped in Step5.
Step7, Obtain whether the answer is 'yes' or 'no', based on the color of lamp and the color of shirt. The color of lamp and shirt is obtained in Step3 and Step6, respectively. If their color are the same, the answer is 'yes'; On the contrary, the answer is 'no'.
Step8, Visualize results.
**Program and obtained result in each step:**
Step1 Program: BOX0=LOC(image=IMAGE,object='lamp')
Result of The coordinate of BOX0: [[45, 78, 245, 345]]
Step2 Program: IMAGE0=CROP(image=IMAGE,box=BOX0)
Result of The description of IMAGE0: a photography of a couple of people on a snowboard in the snow
Step3 Program: BOX1=LOC(image=IMAGE0,object='shirt')
Result of BOX1 is empty
Step4 Program: ANSWER0=COUNT(box=BOX1)
Result of ANSWER0: 0
Step5 Program: ANSWER1=EVAL(expr="'yes' if {ANSWER0} > 0 else 'no'")
Result of ANSWER1: no
Step6 Program: FINAL_RESULT=RESULT(var=ANSWER1)
Result of FINAL_RESULT: no
**Error Location**: program
**Reason**: The plan are correct, and can address the given question. But the Step3-Step6 of the program does not match the plan. The plan locates, crops, and asks color of the lamp and shirt, but the program counts the number of lamp.

----------------------------------------------------------------------------------

The failed case needs to be debugged is as follows.
**Instruction**: Is the wall behind a boy?
**Description of the Input Image:** a photography of a man holding a wii remote in his hand
**Human Feedback:** The correct answer is no
**Our Wrong Answer:** yes
Following are the decomposed plan, used program, and obtained result in each step.
**Plan:**
Step1, Locate the boy in the given image, and obtain bounding boxes of the boy.
Step2, Crop the image region behind the boy from the given image, based on bounding boxes of the boy. The bounding boxes are obtained in Step1.
Step3, Locate the wall in the region behind the boy, and obtain bounding boxes of the wall. The region behind the boy is cropped in Step2.
Step4, Count the number of walls, based on bounding boxes of walls. The bounding boxes are obtained in Step3.
Step5, Obtain whether the answer is 'yes' or 'no', based on the number of walls. The number is obtained in Step4. If the number is greater than zero, the answer is 'yes'; On the contrary, the answer is 'no'.
Step6, Visualize results.
**Program and obtained result in each step:**
Step1 Program: BOX0=LOC(image=IMAGE,object='boy')
The coordinate of BOX0: [[100, 43, 414, 374]]
Step2 Program: IMAGE0=CROP_BEHIND(image=IMAGE,box=BOX0)
The description of IMAGE0: a photography of a man holding a wii remote in his hand
Step3 Program: BOX1=LOC(image=IMAGE0,object='wall')
The coordinate of BOX1: [[279, 1, 498, 207], [30, 1, 498, 207]]
Step4 Program: ANSWER0=COUNT(box=BOX1)
Result of ANSWER0: 2
Step5 Program: ANSWER1=EVAL(expr="'yes' if {ANSWER0} > 0 else 'no'")
Result of ANSWER1: yes
Step6 Program: FINAL_RESULT=RESULT(var=ANSWER1)
Result of FINAL_RESULT: yes
**Error Location:**
**Reason:**

Figure 14. Prompts for global reflection.

image features with our learnable prompts to generate per-mask embeddings that are combined with pre-pixel embedding for mask prediction. In our prompt tuning scheme, the visual features are extracted from the backbone followed by average pooling, whose dimension is 256. We concatenate learnable prompts with the inputs to the transformer decoder. The architecture of Maskformer is shown in Fig. 19.

You are a **debugger**. You need to check which model cause of the wrong answer. After given one step, you need to determine if this step is correct. If it is incorrect, you need to provide the error location, and explain the reason.

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

The failed case is as follows.
**Instruction**: On which side of the photo is the stuffed bear?
**Description of the Input Image**: a photography of a display case with teddy bears
**Desirable Answer**: right
**Our Wrong Answer**: left
It has totally 4 steps. Step1 have been checked:
Step1
**Plan**:  Locate stuffed bear in the given image
**Program**: BOX0=LOC(image=IMAGE, object='stuffed bear')
**Result of The coordinate of BOX0**: [[136, 58, 184, 116], [191, 87, 227, 135]]]

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

Now you need to check Step2.
**Plan**: Crop the region of stuffed bear from the given image.
**Program**: IMAGE0=CROP(image=IMAGE,box=BOX0)
**The description of IMAGE0**: a photography of a group of stuffed animals sitting next to each other.

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

Is this Step2 correct? Errors may exist in the plan, program, or functions called by the program If they are correct, directly output "yes" and do not output any other content. If incorrect, firstly output "no", then provide the error location and reason.

Figure 15. Prompts for local reflection.

### 9.3.2 Training

The size of learnable prompts is $\mathbb{R}^{100 \times 256}$, that is we learn 100 vectors as prompts, and the dimension of each one is 256. We remove the classification loss of Maskformer, and only use the mask loss [6] to train learnable prompts, where the Adam optimizer is used and the learning rate is $1e - 1$. Given a visual concept that needs to be learned, we randomly select 50 samples having this visual concept from the LVIS dataset [10]. We train the prompts 100 steps for each instance. After training, we store the concept name, and image features of instances, and learned prompts of image features in the the prompt pool.

### 9.3.3 Inference

The prompt ensemble process of the SEG tool also has two steps. (1) Given a visual concept, we roughly select out 10 prompts having the same visual concept from the prompt pool as candidates, by computing the similarity between the given query instance and stored instances in the prompt pool. (2) We use prompt ensemble (detailed in Section 3.4) to aggregate the 10 prompts for a query instance. In other words, we do not aggregate all stored prompts for a query instance, but 10 similar instances.

### 9.4. Update SELECT and CLASSIFY tools

### 9.4.1 Model

We utilize CLIP [57] as the SELECT and CLASSIFY tools. The model consists of a Vision Transformer (ViT) as the image encoder and a Transformer-based text encoder. The image encoder and text encoder encode images and text descriptions into high-dimensional feature vectors respectively. It learns to align the representations of related content and separate unrelated content in the embedding space by utilizing a contrastive loss function. This loss function encourages CLIP to maximize the similarity between corresponding image-text pairs while minimizing it for non-corresponding pairs. In the prompt tuning scheme, learnable tokens are introduced into the image encoder and fine-tuned. These prompts are stored in the prompt pool for later use. During inference, the prompt is selected from the prompt pool and replaced with the image encoder to improve the precision of query data. The architecture of CLIP is shown in Fig. 20.

### 9.4.2 Training

We update the CLIP through the deep prompt tuning [21]. We utilize the Adam optimizer with a learning rate set to $1e - 2$. To obtain training data, CLOVA uses the concept to be learned to automatically retrieve 7 images from Google Images as positive data through web scraping. The first image is used for validation, while the remaining ones are used for training. Furthermore, CLOVA derives additional concepts related to but different from the target concept through GPT and subsequently collected data associated with these concepts as negative samples. CLOVA then uses both the positive data and negative data for training. During the training phase, we introduce 100 learnable prompts for each of the first three layers of the vision transformer to facilitate prompt tuning. We conduct prompt tuning for 100 steps per instance, followed by a validation step. If the accurate prediction is achieved on the validation data, CLOVA systematically stores features of crawled images and learned prompts in the prompt pool, otherwise, these learned prompts are discarded.

You are an **inference maker**. Your goal is that, given a failed case including a question, its desirable answer, our wrong answer, our plan, our programs, and the analysis about the wrong step, you need to infer the desirable intermediate results of the wrong step.

---------------------------------------------------------------------------------------------------------

Following are some inferring examples.
**Instruction**: It there three bottles on the table?
**The description of Input image**: a photography of a restaurant with a table set
**Human Feedback**: The correct answer is yes
**Our Wrong Answer**: no
Following are the plan, used program, and obtained result in each step.
**Plan**:
Step1, Locate table in the given image, and obtain bounding boxes of table.
Step2, Crop the region of table from the given image, based on bounding boxes of table. The bounding boxes are obtained in Step1.
Step3, Asking the image region of table, 'How many bottles?'. The image region of the table is obtained in Step2.
Step4, Obtain the answer, based on the intermediate answers obtained in Step3.
Step5, Visualize results.
**Program and obtained result in each step**:
Step1
Program: BOX0=LOC(image=IMAGE,object='table')
The coordinate of BOX0: [[40, 240, 581, 479]
Step2
Program: IMAGE0=CROP(image=IMAGE,box=BOX0)
The description of IMAGE0: a photography of a table set with a white table cloth and red plates
Step3
Program: ANSWER0=VQA(image=IMAGE0,question= 'How many bottles?')
Results of ANSWER0: two
Step4
Program: ANSWER1=EVAL(expr="'yes' if {ANSWER0} == three' else 'no'")
Result of ANSWER1: no
Step5:
Program: FINAL_RESULT=RESULT(var=ANSWER1)
Result of FINAL_RESULT: no
**Error Location**: functions called by programs
**Reason**: In the Step3 of the program, the used function 'VQA' failed to count the number of bottles, as the obtained result of ANSWER0 is 'two' instead of 'three'.
**Correct answer of the wrong step**: three

---------------------------------------------------------------------------------------------------------

Now, you will be given the failed case that needs to infer, as follows. Based on the expected answer, the intermediate results we got at each step of the program, and the analysis of the wrong step. You need to inference the desirable intermediate results of the wrong VQA step.
**Question**: What color is the dog, brown or red?
**Description of the Input Image**: a photography of a group of motorcycles parked in a field
**Desirable Answer**: red
**Our Wrong Answer**: brown
**Following are the decomposed plan, used program, and obtained result in each step**.
**Plan**:
Step1, Locate the dog in the given image and obtain bounding boxes of the dog.
Step2, Crop the region where the dog is located from the given image, based on the bounding boxes of the dog. The bounding boxes are obtained in Step1.
Step3, Asking the image region of dog, 'What color is the dog?'. The image region of the table is obtained in Step2.
Step4, Visualize results.
**Program and obtained result in each step:**
Step1
Program: BOX0=LOC(image=IMAGE,object='dog')
The coordinate of BOX0: [[51, 65, 83, 96]]
Step2
Program: IMAGE0=CROP(image=IMAGE,box=BOX0)
The description of IMAGE0: a photography of a dog is sniffing a toy in the grass
Step3
Program: ANSWER0=VQA(image=IMAGE0,question='What color is the dog?')
Result of ANSWER0: brown
Step4:
Program: FINAL_RESULT=RESULT(var=ANSWER0)
Result of FINAL_RESULT: brown
**Error Location:**
functions called by programs
**Reason:**
In Step3, the function 'VQA' failed to correctly infer the color of the dog based on the cropped region. The obtained result of ANSWER0 is 'brown', which is not one of the expected colors 'red'
**Correct answer of the wrong step:**

Figure 16. Prompts of inferring answers for the VQA tool.

### 9.4.3 Inference

The prompt ensemble process of the SELECT and CLAS-SIFY tools also has two steps. (1) Given query data, we se-lect data with its feature and prompt from the prompt pool with the same concept as the query data. Subsequently, we compute the similarity between the query data and the se-lected data. The prompts corresponding to data with high

Figure 17. The architecture of BLIP.



Figure 18. The architecture of OWL-ViT.

similarities are used for the query data. (2) We use prompt ensemble (detailed in Section 3.4) to aggregate selected prompts for the query data.

## 9.5. Update REPLACE tool

We use Stable Diffusion [58] as the REPLACE tool. In Stable Diffusion, the architecture comprises four key com-ponents: Sampler, Variational Autoencoder (VAE), UNet, and CLIPEmbedder. The Sampler and UNet focus on the actual image generation, the VAE provides a deep under-standing of image content, and the CLIPEmbedder ensures the relevance and accuracy of the generated images in rela-tion to the text inputs. We added learnable prompts to the text encoder of the CLIPEmbedder component. During the

Figure 19. The architecture of Maskformer.



Figure 20. The architecture of CLIP.

prompt tuning phase, we tune a prompt for training images and store it. The dimensionality of the prompt is 768. The architecture of Stable Diffusion is shown in Fig. 21.

### 9.5.1 Training

In order to facilitate the learning of a specific concept for the Stable Diffusion model, CLOVA employs web scraping techniques to retrieve 7 images representing this concept from Google Images. We train prompts in the text

Figure 21. The architecture of Stable Diffusion.

encoder of the Stable Diffusion model using downloaded images. During the training process, the Latent Diffusion Model(LDM) loss [58] is minimized. Subsequently, the 768-dimensional prompt obtained from the training stage is stored in the prompt pool. In terms of experimental setup, we employ the Adam optimizer. Through our experimentation, we find that setting the learning rate to $5e-3$ achieves the best learning results.

### 9.5.2 Inference

Similar to other visual tools, inference for the Stable Diffusion model also contains two steps. Given query data, CLOVA first locates prompts corresponding to the concept in the prompt pool and then loads the prompt into the text encoder of the Stable Diffusion model. Based on the query and mask, we perform editing on the input image.

## 10. More Experimental Results

### 10.1. Training-validation prompt tuning for the VQA tool

We further evaluate the proposed validation-learning prompt tuning scheme for the VQA tool, where experiments are conducted on the compositional VQA task using the GQA dataset. We use the BLIP model for the GQA dataset and compare our prompt tuning scheme with direct



Figure 22. Accuracy curves on the GQA dataset

tuning parameters. We report accuracies with using different numbers of training data. Results are shown in Fig. 22. Our method has higher performance throughout the entire training process, no matter whether the number of training data is small or large, showing the effectiveness of the proposed validation-learning prompt tuning scheme for the VQA tool.

### 10.2. Evaluation on the online setting

CLOVA can be applied to a more practical online learning setting. In this case, CLOVA is evaluated in a dynamic data stream. If it makes a correct prediction on a task, only the inference phase is activated, and the task is tagged as a correct prediction; if it makes an incorrect prediction on a task,

| Dataset | Method | LLama2-7B | GPT-3.5-turbo | GPT-4 |
|---------|--------|-----------|---------------|-------|
| GQA | Baseline | 39.2 | 46.4 | 52.6 |
| | + Update LLMs | 44.8 | 51.0 | 55.4 |
| | + Update visual tools | 50.2 | 53.0 | 57.8 |
| NLVRv2 | Baseline | 50.0 | 60.2 | 64.8 |
| | + Update LLMs | 57.4 | 61.0 | 66.4 |
| | + Update visual tools | 61.6 | 62.6 | 67.4 |

Table 8. Different LLMs on the online learning setting using the GQA and NLVRv2 datasets.

this task is tagged as a wrong prediction, and the reflection and learning phases are activated to update tools. After the data stream, we calculate the accuracy based on tagged predictions of all cases. We conduct experiments on the compositional VQA and multi-image reasoning tasks, where the GQA and NLVRv2 datasets are used. Results are shown in Tab. 8. Similar to the offline setting in Section 4.2, updating LLMs and visual tools both leads to improvements.

## 10.3. More case studies

We provide more cases to show the reflection and learning phases of CLOVA. The reflection and learning phases for LLMs are shown in Fig. 23. The reflection and learning phases for the SELECT tool are shown in Fig. 24. The reflection and learning phases for the LOC tool are shown in Fig. 25. The reflection and learning phases for the REPLACE tool are shown in Fig. 26. The reflection and learning phases for the CLASSIFY tool are shown in Fig. 27. The reflection and learning phases for the SEG tool are shown in Fig. 28.

## References

[1] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. In *NeurIPS*, pages 23716–23736, 2022.

[2] Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. Self-rag: Learning to retrieve, generate, and critique through self-reflection. In *ICLR*, 2024.

[3] Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. Constitutional ai: Harmlessness from ai feedback. *arXiv preprint arXiv:2212.08073*, 2022.

[4] Tianle Cai, Xuezhi Wang, Tengyu Ma, Xinyun Chen, and Denny Zhou. Large language models as tool makers. In *ICLR*, 2024.

[5] Xinyun Chen, Maxwell Lin, Nathanael Schärli, and Denny Zhou. Teaching large language models to self-debug. In *ICLR*, 2024.

[6] Bowen Cheng, Alex Schwing, and Alexander Kirillov. Per-pixel classification is not all you need for semantic segmentation. In *NeurIPS*, pages 17864–17875, 2021.

[7] Yue Fan, Xiaojian Ma, Rujie Wu, Yuntao Du, Jiaqi Li, Zhi Gao, and Qing Li. Videoagent: A memory-augmented multimodal agent for video understanding. *arXiv preprint arXiv:2403.11481*, 2024.

[8] Difei Gao, Lei Ji, Luowei Zhou, Kevin Qinghong Lin, Joya Chen, Zihan Fan, and Mike Zheng Shou. Assistgpt: A general multi-modal assistant that can plan, execute, inspect, and learn. *arXiv preprint arXiv:2306.08640*, 2023.

[9] Yingqiang Ge, Wenyue Hua, Jianchao Ji, Juntao Tan, Shuyuan Xu, and Yongfeng Zhang. Openagi: When llm meets domain experts. In *NeurIPS*, pages 5539–5568, 2023.

[10] Agrim Gupta, Piotr Dollar, and Ross Girshick. Lvis: A dataset for large vocabulary instance segmentation. In *CVPR*, pages 5356–5364, 2019.

[11] Tanmay Gupta and Aniruddha Kembhavi. Visual programming: Compositional visual reasoning without training. In *CVPR*, pages 14953–14962, 2023.

[12] Danna Gurari, Qing Li, Abigale J Stangl, Anhong Guo, Chi Lin, Kristen Grauman, Jiebo Luo, and Jeffrey P Bigham. Vizwiz grand challenge: Answering visual questions from blind people. In *CVPR*, pages 3608–3617, 2018.

[13] Shibo Hao, Tianyang Liu, Zhen Wang, and Zhiting Hu. Toolkengpt: Augmenting frozen language models with massive tools via tool embeddings. In *NeurIPS*, pages 45870–45894, 2023.

[14] Yining Hong, Qing Li, Daniel Ciao, Siyuan Huang, and Song-Chun Zhu. Learning by fixing: Solving math word problems with weak supervision. In *Proceedings of the AAAI conference on artificial intelligence*, pages 4959–4967, 2021.

[15] Yining Hong, Qing Li, Ran Gong, Daniel Ciao, Siyuan Huang, and Song-Chun Zhu. Smart: A situation model for algebra story problems via attributed grammar. In *Proceedings of the AAAI conference on artificial intelligence*, pages 13009–13017, 2021.

[16] Pengbo Hu, Ji Qi, Xingyu Li, Hong Li, Xinqi Wang, Bing Quan, Ruiyu Wang, and Yi Zhou. Tree-of-mixed-thought: Combining fast and slow thinking for multi-hop visual reasoning. *arXiv preprint arXiv:2308.09658*, 2023.

[17] Xiaowei Hu, Zhe Gan, Jianfeng Wang, Zhengyuan Yang, Zicheng Liu, Yumao Lu, and Lijuan Wang. Scaling up vision-language pre-training for image captioning. In *CVPR*, pages 17980–17989, 2022.

[18] Jie Huang, Xinyun Chen, Swaroop Mishra, Huaixiu Steven Zheng, Adams Wei Yu, Xinying Song, and Denny Zhou. Large language models cannot self-correct reasoning yet. In *ICLR*, 2024.

[19] Siteng Huang, Biao Gong, Yulin Pan, Jianwen Jiang, Yiliang Lv, Yuyuan Li, and Donglin Wang. Vop: Text-video co-operative prompt tuning for cross-modal retrieval. In *CVPR*, pages 6565–6574, 2023.

[20] Drew A. Hudson and Christopher D. Manning. GQA: A new dataset for real-world visual reasoning and compositional question answering. In *CVPR*, pages 6700–6709, 2019.

[21] Menglin Jia, Luming Tang, Bor-Chun Chen, Claire Cardie, Serge Belongie, Bharath Hariharan, and Ser-Nam Lim. Visual prompt tuning. In *ECCV*, pages 709–727, 2022.

Figure 23. Case studies of updating LLMs.



Figure 24. Case studies of updating the SELECT tool.
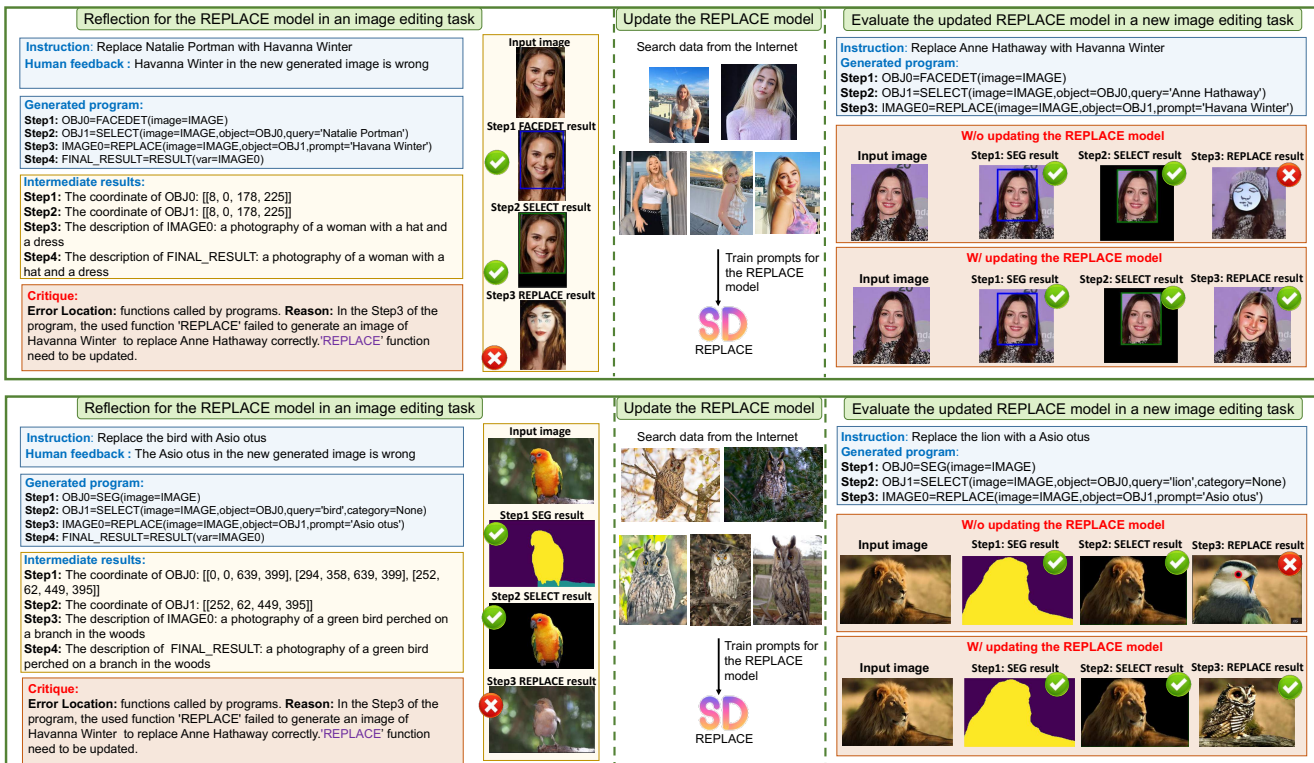
Figure 25. Case studies of updating the LOC tool.



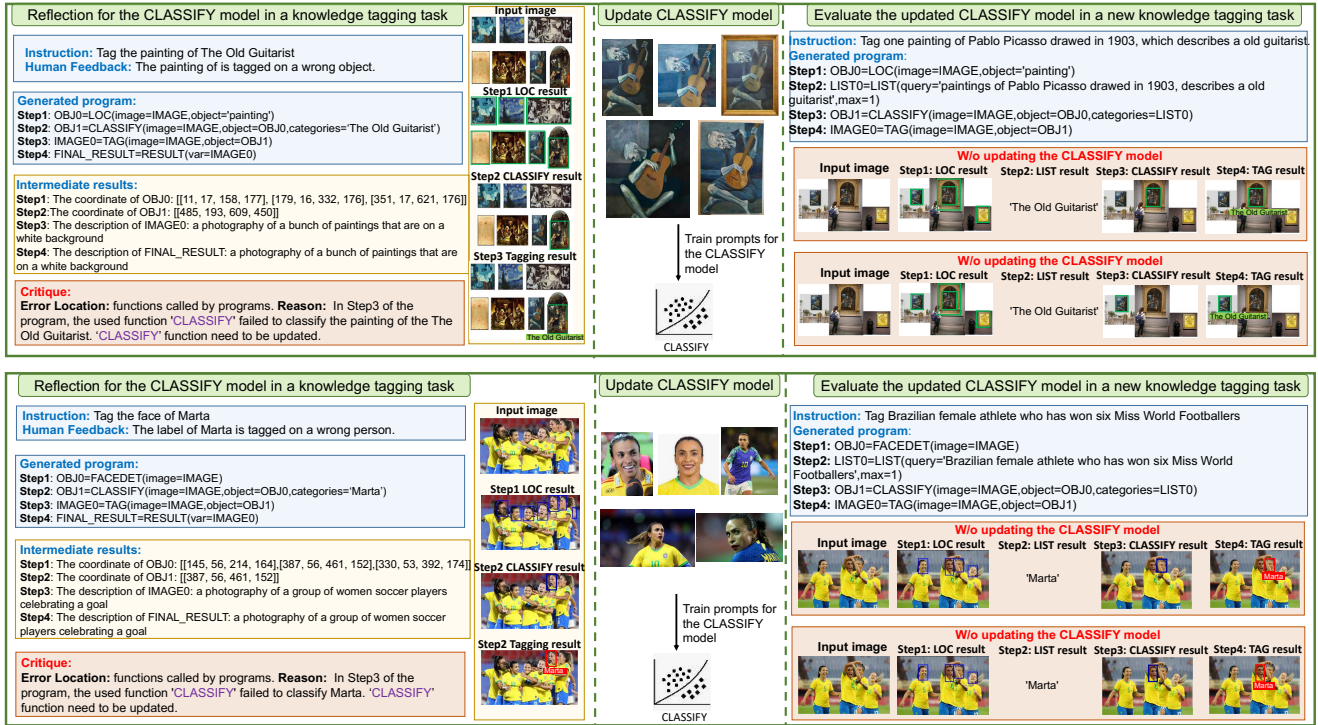Figure 26. Case studies of updating the REPLACE tool.

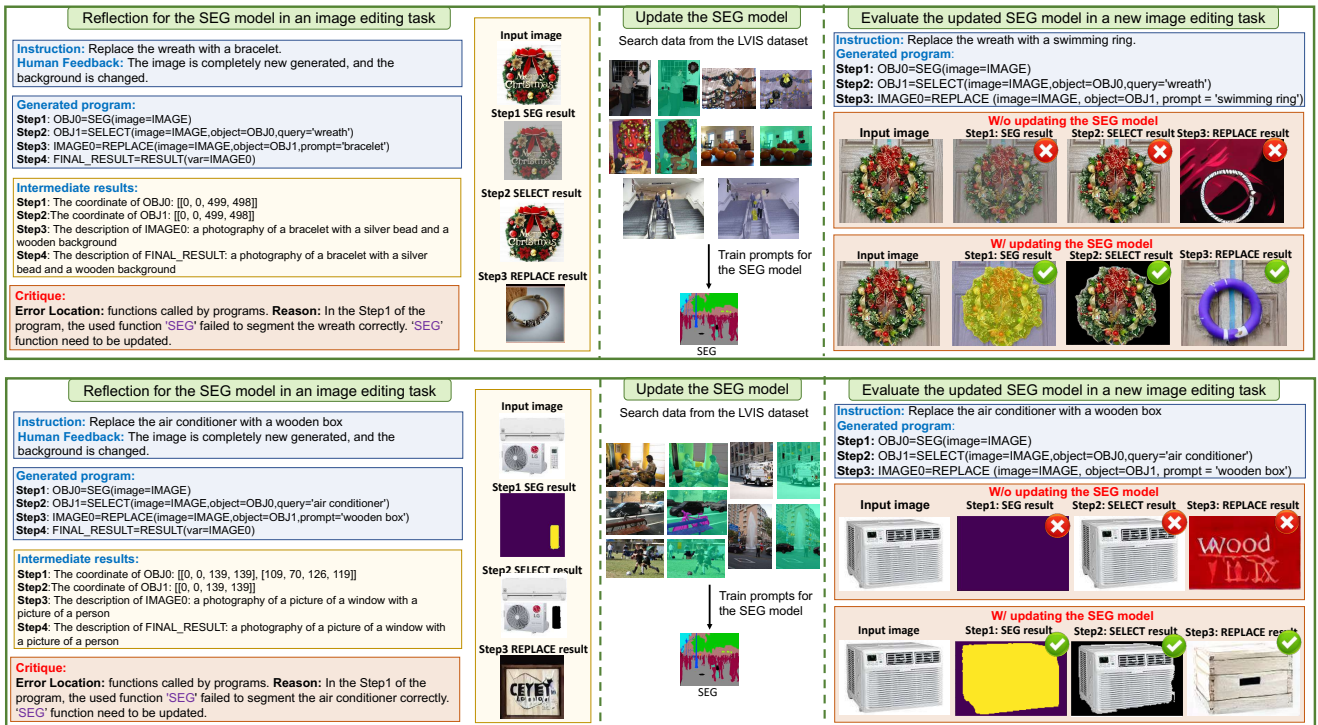Figure 27. Case studies of updating the CLASSIFY tool.



Figure 28. Case studies of updating the SEG tool.

[22] Zhengbao Jiang, Frank F Xu, Jun Araki, and Graham Neubig. How can we know what language models know? *TACL*, 8:423–438, 2020.

[23] Jacob Devlin Ming-Wei Chang Kenton and Lee Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *NAACL-HLT*, pages 4171–4186, 2019.

[24] Muhammad Uzair Khattak, Hanoona Rasheed, Muhammad Maaz, Salman Khan, and Fahad Shahbaz Khan. Maple: Multi-modal prompt learning. In *CVPR*, pages 19113–19122, 2023.

[25] Geunwoo Kim, Pierre Baldi, and Stephen McAleer. Language models can solve computer tasks. In *NeurIPS*, pages 39648–39677, 2023.

[26] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollar, and Ross Girshick. Segment anything. In *ICCV*, pages 4015–4026, 2023.

[27] Bo Li, Yuanhan Zhang, Liangyu Chen, Jinghao Wang, Jingkang Yang, and Ziwei Liu. Otter: A multi-modal model with in-context instruction tuning. *arXiv preprint arXiv:2305.03726*, 2023.

[28] Jian Li, Yabiao Wang, Changan Wang, Ying Tai, Jianjun Qian, Jian Yang, Chengjie Wang, Jilin Li, and Feiyue Huang. Dsfd: dual shot face detector. In *CVPR*, pages 5060–5069, 2019.

[29] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *ICML*, pages 12888–12900, 2022.

[30] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *ICML*, pages 19730–19742, 2023.

[31] Qing Li, Jianlong Fu, Dongfei Yu, Tao Mei, and Jiebo Luo. Tell-and-answer: Towards explainable visual question answering using attributes and captions. *EMNLP*, 2018.

[32] Qing Li, Qingyi Tao, Shafiq Joty, Jianfei Cai, and Jiebo Luo. Vqa-e: Explaining, elaborating, and enhancing your answers for visual questions. In *ECCV*, pages 552–567, 2018.

[33] Qing Li, Siyuan Huang, Yining Hong, Yixin Chen, Ying Nian Wu, and Song-Chun Zhu. Closed loop neural-symbolic learning via integrating neural perception, grammar parsing, and symbolic reasoning. In *ICML*, pages 5884–5894. PMLR, 2020.

[34] Qing Li, Siyuan Huang, Yining Hong, and Song-Chun Zhu. A competence-aware curriculum for visual concepts learning via question answering. In *European Conference on Computer Vision*, pages 141–157. Springer, 2020.

[35] Qing Li, Yixin Zhu, Yitao Liang, Ying Nian Wu, Song-Chun Zhu, and Siyuan Huang. Neural-symbolic recursive machine for systematic generalization. *ICLR*, 2024.

[36] Yaobo Liang, Chenfei Wu, Ting Song, Wenshan Wu, Yan Xia, Yu Liu, Yang Ou, Shuai Lu, Lei Ji, Shaoguang Mao, et al. Taskmatrix. ai: Completing tasks by connecting foundation models with millions of apis. *arXiv preprint arXiv:2303.16434*, 2023.

[37] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. In *NeurIPS Workshop*, 2023.

[38] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In *NeurIPS*, pages 34892–34916, 2023.

[39] Shilong Liu, Hao Cheng, Haotian Liu, Hao Zhang, Feng Li, Tianhe Ren, Xueyan Zou, Jianwei Yang, Hang Su, Jun Zhu, Lei Zhang, Jianfeng Gao, and Chunyuan Li. Llava-plus: Learning to use tools for creating multimodal agents. *2311.05437,arXiv*, 2023.

[40] Zhaoyang Liu, Yinan He, Wenhai Wang, Weiyun Wang, Yi Wang, Shoufa Chen, Qinglong Zhang, Yang Yang, Qingyun Li, Jiashuo Yu, et al. Interngpt: Solving vision-centric tasks by interacting with chatbots beyond language. *arXiv preprint arXiv:2305.05662*, 2023.

[41] Zhaoyang Liu, Zeqiang Lai, Gao Zhangwei, Erfei Cui, Zhiheng Li, Xizhou Zhu, Lewei Lu, Qifeng Chen, Yu Qiao, Jifeng Dai, and Wang Wenhai. Controlllm: Augment language models with tools by searching on graphs. *arXiv preprint arXiv:2305.10601*, 2023.

[42] Pan Lu, Baolin Peng, Hao Cheng, Michel Galley, Kai-Wei Chang, Ying Nian Wu, Song-Chun Zhu, and Jianfeng Gao. Chameleon: Plug-and-play compositional reasoning with large language models. In *NeurIPS*, pages 43447–43478, 2023.

[43] Yuning Lu, Jianzhuang Liu, Yonggang Zhang, Yajing Liu, and Xinmei Tian. Prompt distribution learning. In *CVPR*, pages 5206–5215, 2022.

[44] Yun Luo, Zhen Yang, Fandong Meng, Yafu Li, Jie Zhou, and Yuechen Zhang. An empirical study of catastrophic forgetting in large language models during continual fine-tuning. *ArXiv*, abs/2308.08747, 2023.

[45] Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, et al. Self-refine: Iterative refinement with self-feedback. In *NeurIPS*, pages 46534–46594, 2023.

[46] Matthias Minderer, Alexey Gritsenko, Austin Stone, Maxim Neumann, Dirk Weissenborn, Alexey Dosovitskiy, Aravindh Mahendran, Anurag Arnab, Mostafa Dehghani, Zhuoran Shen, et al. Simple open-vocabulary object detection. In *ECCV*, pages 728–755, 2022.

[47] Varun Nair, Elliot Schumacher, Geoffrey Tso, and Anitha Kannan. Dera: enhancing large language model completions with dialog-enabled resolving agents. *arXiv preprint arXiv:2303.17071*, 2023.

[48] OpenAI. GPT-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.

[49] Liangming Pan, Michael Saxon, Wenda Xu, Deepak Nathani, Xinyi Wang, and William Yang Wang. Automatically correcting large language models: Surveying the landscape of diverse self-correction strategies. *arXiv preprint arXiv:2308.03188*, 2023.

[50] Bhargavi Paranjape, Scott M. Lundberg, Sameer Singh, Hannaneh Hajishirzi, Luke Zettlemoyer, and Marco Tulio Ribeiro. Art: Automatic multi-step reasoning and tool-use for large language models. *ArXiv*, abs/2303.09014, 2023.

[51] Joon Sung Park, Joseph C O'Brien, Carrie J Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. Generative agents: Interactive simulacra of human behavior. In *UIST*, pages 1–22, 2023.

[52] Shishir G Patil, Tianjun Zhang, Xin Wang, and Joseph E Gonzalez. Gorilla: Large language model connected with massive apis. *arXiv preprint arXiv:2305.15334*, 2023.

[53] 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. *arXiv preprint arXiv:2306.14824*, 2023.

[54] Sarah Pratt, Ian Covert, Rosanne Liu, and Ali Farhadi. What does a platypus look like? generating customized prompts for zero-shot image classification. In *ICCV*, pages 15691–15701, 2023.

[55] Shuofei Qiao, Honghao Gui, Huajun Chen, and Ningyu Zhang. Making language models better tool learners with execution feedback. *ArXiv*, abs/2305.13068, 2023.

[56] Yujia Qin, Shengding Hu, Yankai Lin, Weize Chen, Ning Ding, Ganqu Cui, Zheni Zeng, Yufei Huang, Chaojun Xiao, Chi Han, Yi Ren Fung, Yusheng Su, Huadong Wang, Cheng Qian, Runchu Tian, Kunlun Zhu, Shihao Liang, Xingyu Shen, Bokai Xu, Zhen Zhang, Yining Ye, Bowen Li, Ziwei Tang, Jing Yi, Yuzhang Zhu, Zhenning Dai, Lan Yan, Xin Cong, Yaxi Lu, Weilin Zhao, Yuxiang Huang, Junxi Yan, Xu Han, Xian Sun, Dahai Li, Jason Phang, Cheng Yang, Tongshuang Wu, Heng Ji, Zhiyuan Liu, and Maosong Sun. Tool learning with foundation models. *arXiv preprint arXiv:2304.08354*, 2023.

[57] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *ICML*, pages 8748–8763, 2021.

[58] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, pages 10684–10695, 2022.

[59] William Saunders, Catherine Yeh, Jeff Wu, Steven Bills, Long Ouyang, Jonathan Ward, and Jan Leike. Self-critiquing models for assisting human evaluators. *arXiv preprint arXiv:2206.05802*, 2022.

[60] Timo Schick, Jane Dwivedi-Yu, Roberto Dessi, Roberta Raileanu, Maria Lomeli, Eric Hambro, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. Toolformer: Language models can teach themselves to use tools. In *NeurIPS*, pages 68539–68551, 2023.

[61] Zhenwei Shao, Zhou Yu, Meng Wang, and Jun Yu. Prompting large language models with answer heuristics for knowledge-based visual question answering. In *CVPR*, pages 14974–14983, 2023.

[62] Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, Weiming Lu, and Yueting Zhuang. Hugginggpt: Solving ai tasks with chatgpt and its friends in huggingface. In *NeurIPS*, pages 38154–38180, 2023.

[63] Taylor Shin, Yasaman Razeghi, Robert L Logan IV, Eric Wallace, and Sameer Singh. Autoprompt: Eliciting knowl-

[64] Noah Shinn, Federico Cassano, Beck Labash, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: Language agents with verbal reinforcement learning. In *NeurIPS*, pages 8634–8652, 2023.

[65] Kaya Stechly, Matthew Marquez, and Subbarao Kambhampati. Gpt-4 doesn't know it's wrong: An analysis of iterative prompting for reasoning problems. In *NeurIPS Workshop*, 2023.

[66] Alane Suhr, Stephanie Zhou, Ally Zhang, Iris Zhang, Huajun Bai, and Yoav Artzi. A corpus for reasoning about natural language grounded in photographs. In *ACL*, pages 6418–6428, 2019.

[67] Dídac Surís, Sachit Menon, and Carl Vondrick. Vipergpt: Visual inference via python execution for reasoning. In *ICCV*, pages 11888–11898, 2023.

[68] Jiajin Tang, Ge Zheng, Jingyi Yu, and Sibei Yang. Cotdet: Affordance knowledge prompting for task driven object detection. In *ICCV*, pages 3068–3078, 2023.

[69] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.

[70] Karthik Valmeekam, Matthew Marquez, and Subbarao Kambhampati. Can large language models really improve by self-critiquing their own plans? In *NeurIPS Workshop*, 2023.

[71] Andrés Villa, Juan León Alcázar, Motasem Alfarra, Kumail Alhamoud, Julio Hurtado, Fabian Caba Heilbron, Alvaro Soto, and Bernard Ghanem. Pivot: Prompting for video continual learning. In *CVPR*, pages 24214–24223, 2023.

[72] Chenyu Wang, Weixin Luo, Qianyu Chen, Haonan Mai, Jindi Guo, Sixun Dong, Xiaohua (Michael) Xuan, Zhengxin Li, Lin Ma, and Shenghua Gao. Mllm-tool: A multimodal large language model for tool agent learning. *arXiv preprint arXiv:2401.10727*, 2024.

[73] Zifeng Wang, Zizhao Zhang, Sayna Ebrahimi, Ruoxi Sun, Han Zhang, Chen-Yu Lee, Xiaoqi Ren, Guolong Su, Vincent Perot, Jennifer Dy, et al. Dualprompt: Complementary prompting for rehearsal-free continual learning. In *ECCV*, pages 631–648, 2022.

[74] Zifeng Wang, Zizhao Zhang, Chen-Yu Lee, Han Zhang, Ruoxi Sun, Xiaoqi Ren, Guolong Su, Vincent Perot, Jennifer Dy, and Tomas Pfister. Learning to prompt for continual learning. In *CVPR*, pages 139–149, 2022.

[75] Chenfei Wu, Shengming Yin, Weizhen Qi, Xiaodong Wang, Zecheng Tang, and Nan Duan. Visual chatgpt: Talking, drawing and editing with visual foundation models. *arXiv preprint arXiv:2303.04671*, 2023.

[76] Jinjin Xu, Liwu Xu, Yuzhe Yang, Xiang Li, Yanchun Xie, Yi-Jie Huang, and Yaqian Li. u-llava: Unifying multi-modal tasks via large language model. *2311.05348,arXiv*, 2023.

[77] Shicheng Xu, Liang Pang, Huawei Shen, Xueqi Cheng, and Tat-seng Chua. Search-in-the-chain: Towards the accurate, credible and traceable content generation

for complex knowledge-intensive tasks. *arXiv preprint arXiv:2304.14732*, 2023.

[78] Hao Yan, Saurabh Srivastava, Yintao Tai, Sida I Wang, Wen-tau Yih, and Ziyu Yao. Learning to simulate natural language feedback for interactive semantic parsing. In *ACL*, pages 3149–3170, 2023.

[79] Kevin Yang, Yuandong Tian, Nanyun Peng, and Dan Klein. Re3: Generating longer stories with recursive reprompting and revision. In *EMNLP*, pages 4393–4479, 2022.

[80] Rui Yang, Lin Song, Yanwei Li, Sijie Zhao, Yixiao Ge, Xiu Li, and Ying Shan. Gpt4tools: Teaching large language model to use tools via self-instruction. In *NeurIPS*, pages 71995–72007, 2023.

[81] 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. *arXiv preprint arXiv:2303.11381*, 2023.

[82] Yuanhan Zhang, Kaiyang Zhou, and Ziwei Liu. What makes good examples for visual in-context learning? In *NeurIPS*, pages 17773–17794, 2023.

[83] Haozhe Zhao, Zefan Cai, Shuzheng Si, Xiaojian Ma, Kaikai An, Liang Chen, Zixuan Liu, Sheng Wang, Wenjuan Han, and Baobao Chang. Mmicl: Empowering vision-language model with multi-modal in-context learning. In *ICLR*, 2024.

[84] Yaoyao Zhong, Mengshi Qi, Rui Wang, Yuhan Qiu, Yang Zhang, and Huadong Ma. Viotgpt: Learning to schedule vision tools towards intelligent video internet of things. *arXiv preprint arXiv:2312.00401*, 2023.

[85] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Conditional prompt learning for vision-language models. In *CVPR*, pages 16816–16825, 2022.

[86] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Learning to prompt for vision-language models. *IJCV*, 130(9):2337–2348, 2022.

[87] Pei Zhou, Aman Madaan, Srividya Pranavi Potharaju, Aditya Gupta, Kevin R. McKee, Ari Holtzman, Jay Pujara, Xiang Ren, Swaroop Mishra, Aida Nematzadeh, Shyam Upadhyay, and Manaal Faruqui. How far are large language models from agents with theory-of-mind? *arXiv*, abs/2310.03051, 2023.

[88] Jiawen Zhu, Simiao Lai, Xin Chen, Dong Wang, and Huchuan Lu. Visual prompt multi-modal tracking. In *CVPR*, pages 9516–9526, 2023.

[89] Muzhi Zhu, Hengtao Li, Hao Chen, Chengxiang Fan, Weian Mao, Chenchen Jing, Yifan Liu, and Chunhua Shen. Segprompt: Boosting open-world segmentation via category-level prompt learning. In *ICCV*, pages 999–1008, 2023.

[90] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr: Deformable transformers for end-to-end object detection. *ICLR*, 2021.