# CLOVA: A Closed-Loop Visual Assistant with Tool Usage and Update

Zhi Gao, Yuntao Du, Xintong Zhang, Xiaojian Ma, Wenjuan Han, Song-Chun Zhu, Qing Li

Project

CVPR SEATTLE, WA JUNE 17-21, 2024

## Introduction

**Motivation**: Can VLM improve from user's feedback?



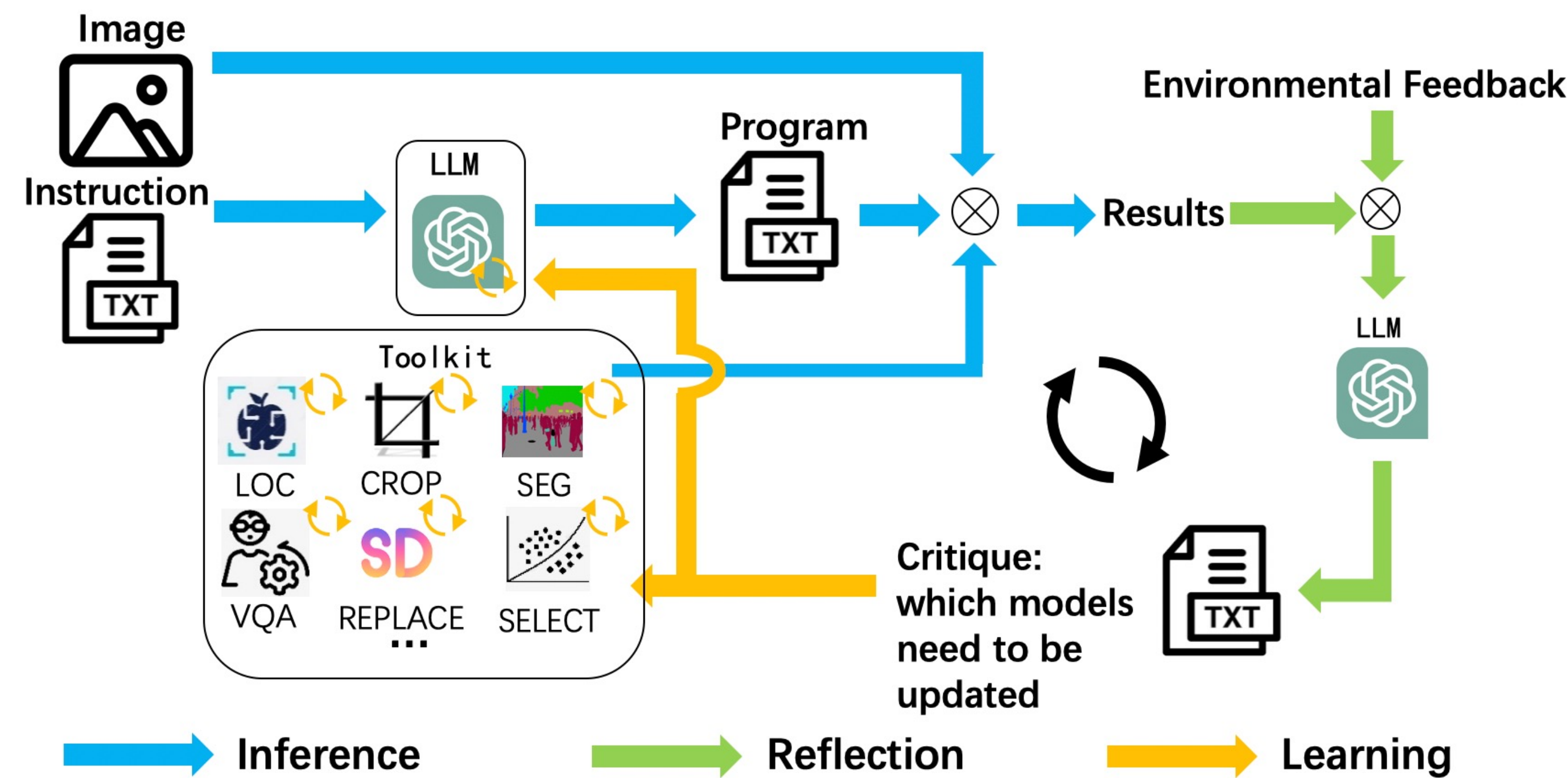(a) Update the CLASSIFY module  (b) Update the REPLACE model  (c) Update LLMs  (d) Update the LOC model
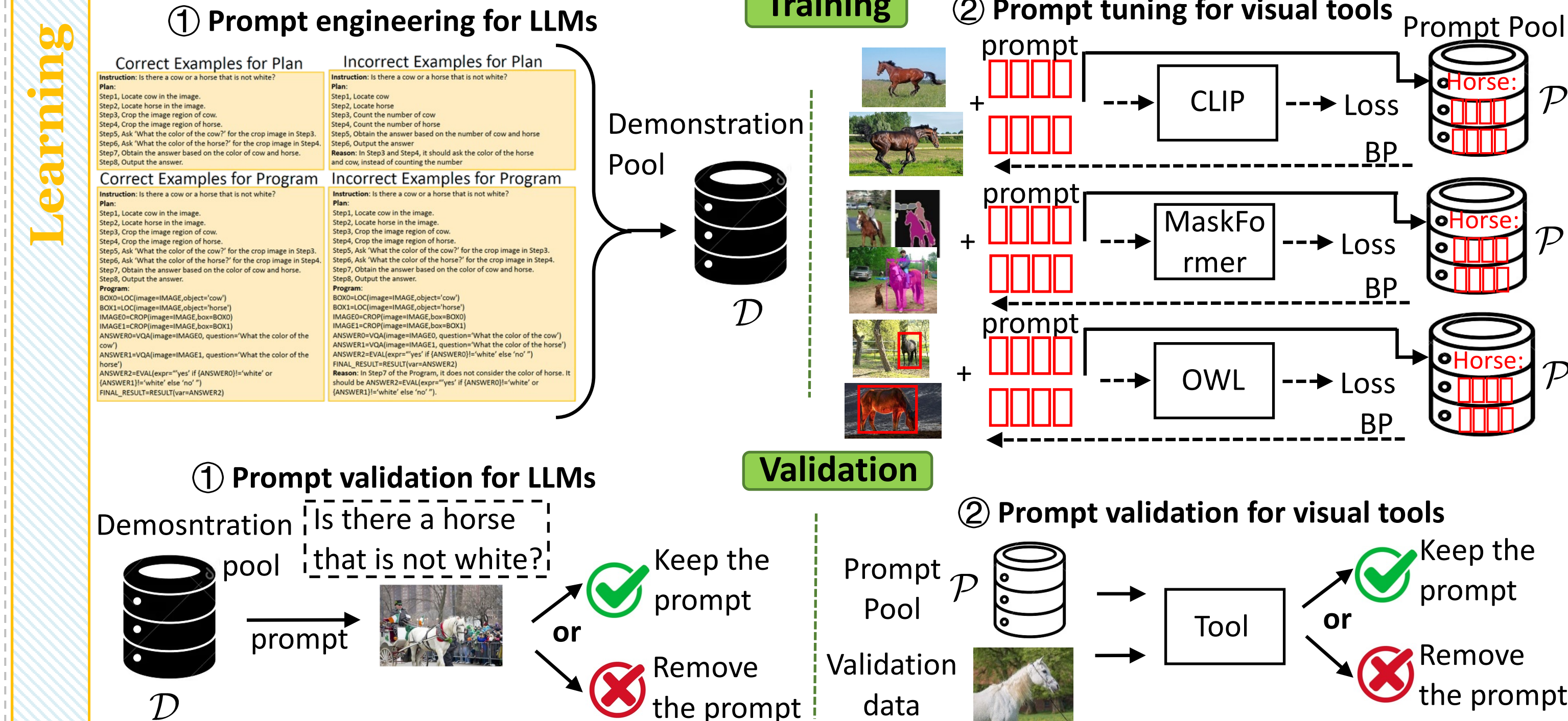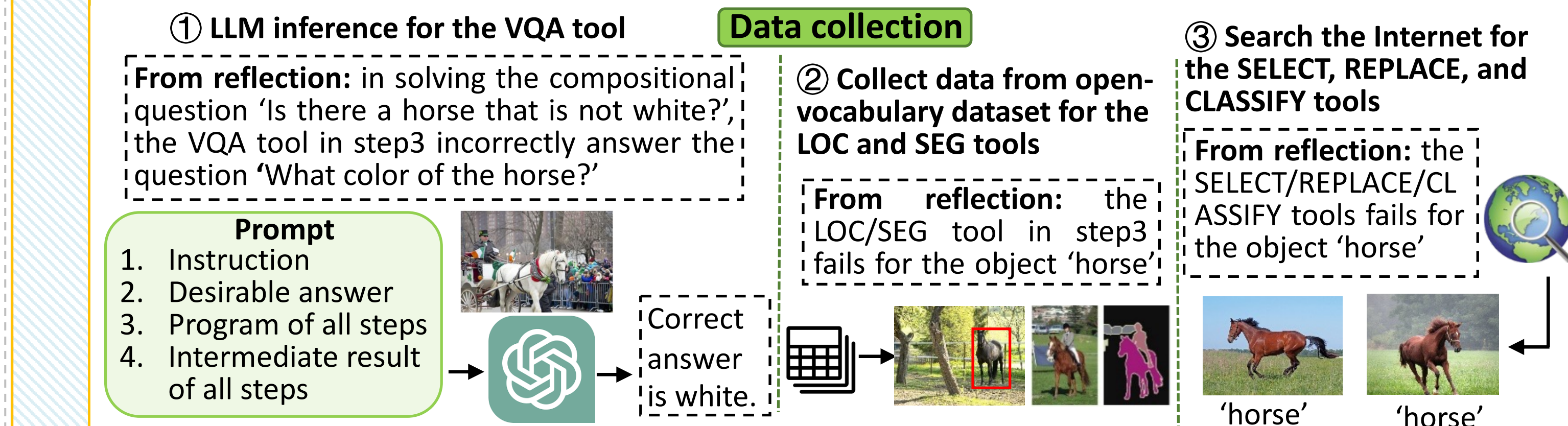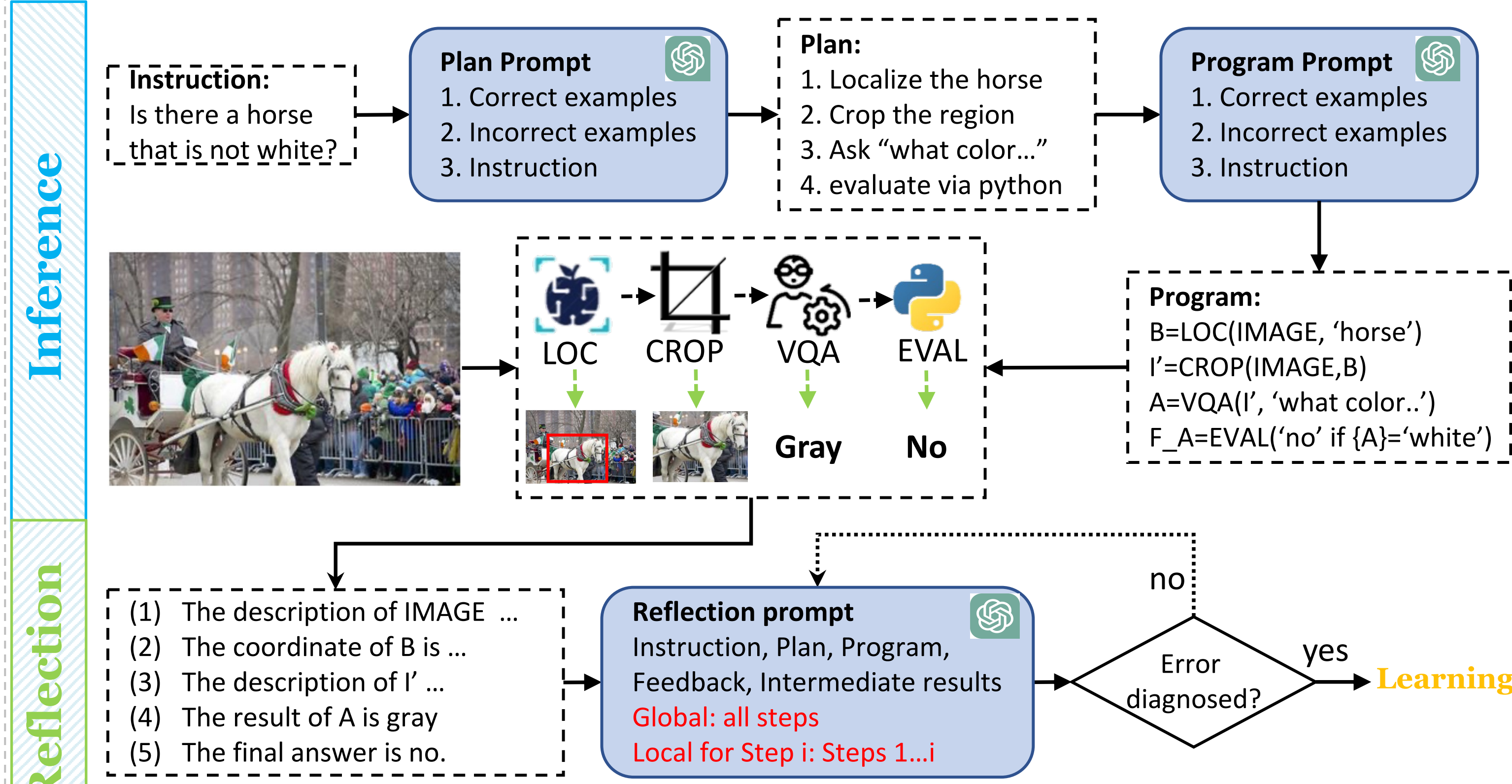
**Challenges:** (1) how to diagnose errors? (2) how to collect training data? (3) how to efficiently update the model?

**Proposal**: CLOVA = Inference + Reflection + Learning

☐ **improve from feedback** via a closed-loop learning framework



Inference   Reflection   Learning

## CLOVA = Inference + Reflection + Learning

### Inference

**Instruction:** Is there a horse that is not white?

**Plan Prompt**
1. Correct examples
2. Incorrect examples
3. Instruction

**Plan:**
1. Localize the horse
2. Crop the region
3. Ask "what color…"
4. evaluate via python

**Program Prompt**
1. Correct examples
2. Incorrect examples
3. Instruction

**Program:**
B=LOC(IMAGE, 'horse')
I'=CROP(IMAGE,B)
A=VQA(I', 'what color..')
F_A=EVAL('no' if {A}='white')

LOC → CROP → VQA → EVAL

Gray   No

### Reflection

(1) The description of IMAGE …
(2) The coordinate of B is …
(3) The description of I' …
(4) The result of A is gray
(5) The final answer is no.

**Reflection prompt**
Instruction, Plan, Program, Feedback, Intermediate results
Global: all steps
Local for Step i: Steps 1…i

Error diagnosed?  no / yes → Learning

### Learning

**① LLM inference for the VQA tool**

**From reflection:** in solving the compositional question 'Is there a horse that is not white?', the VQA tool in step3 incorrectly answer the question 'What color of the horse?'

**Prompt**
1. Instruction
2. Desirable answer
3. Program of all steps
4. Intermediate result of all steps

Correct answer is white.

**Data collection**

**② Collect data from open-vocabulary dataset for the LOC and SEG tools**

**From reflection:** the LOC/SEG tool in step3 fails for the object 'horse'

**③ Search the Internet for the SELECT, REPLACE, and CLASSIFY tools**

**From reflection:** the SELECT/REPLACE/CLASSIFY tools fails for the object 'horse'

'horse'   'horse'

**Training**

**① Prompt engineering for LLMs**

Correct Examples for Plan   Incorrect Examples for Plan

Correct Examples for Program   Incorrect Examples for Program

Demonstration Pool $\mathcal{D}$

**② Prompt tuning for visual tools**

prompt + CLIP → Loss → BP → Prompt Pool $\mathcal{P}$
prompt + MaskFormer → Loss → BP → Prompt Pool $\mathcal{P}$
prompt + OWL → Loss → BP → Prompt Pool $\mathcal{P}$

**Validation**

**① Prompt validation for LLMs**

Demonstration pool → Is there a horse that is not white? → prompt → ✓ Keep the prompt / ✗ Remove the prompt

Demonstration pool $\mathcal{D}$

**② Prompt validation for visual tools**

Prompt Pool $\mathcal{P}$ → Tool → ✓ Keep the prompt / ✗ Remove the prompt
Validation data

## Experiments

### Main results

| | 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 |

### Different LLMs

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

### Ablation studies

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

| 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 |

**Quantitative observation:**

☐ CLOVA achieves SOTA among tool-usage VLMs.

☐ CLOVA is robust to different LLMs, including open and closed ones.

☐ Update both LLM and visual tools bring significant improvements.

### Qualitative example

| Reflection for the REPLACE tool in an image editing task | Update the REPLACE tool | Evaluate the updated REPLACE tool in a new image editing task |



## Takeaway Message

We build **CLOVA**, the first VLM that can **improve from feedback** via a closed-loop learning framework with **inference**, **reflection**, **learning** phases.

☐ Use both correct and incorrect examples for prompts to generate plans & programs.

☐ Propose a global-local reflection scheme to diagnose errors.

☐ Apply hard/soft prompt tuning to update tools with limited data.

**Code & Examples:** clova-tool.github.io