

000
001
002
003
004
005
006
007
008
009
010
011
012
013
014
015
016
017
018
019
020
021
022
023
024
025
026
027
028
029
030
031
032
033
034
035
036
037
038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053

SMALL-LARGE COLLABORATION: TRAINING-EFFICIENT CONCEPT PERSONALIZATION FOR LARGE VLM USING A META PERSONALIZED SMALL VLM

Anonymous authors
Paper under double-blind review

ABSTRACT

Personalizing Vision-Language Models (VLMs) to transform them into daily assistants has emerged as a trending research direction. However, leading companies like OpenAI continue to increase model size and develop complex designs such as the chain of thought (CoT). While large VLMs are proficient in complex multi-modal understanding, their high training costs and limited access via paid APIs restrict direct personalization. Conversely, small VLMs are easily personalized and freely available, but they lack sufficient reasoning capabilities. Inspired by this, we propose a novel collaborative framework named Small-Large Collaboration (SLC) for large VLM personalization, where the small VLM is responsible for generating personalized information, while the large model integrates this personalized information to deliver accurate responses. To effectively incorporate personalized information, we develop a test-time reflection strategy, preventing the potential hallucination of the small VLM. Since SLC only needs to train a meta personalized small VLM for the large VLMs, the overall process is training-efficient. To the best of our knowledge, this is the first training-efficient framework that supports both open-source and closed-source large VLMs, enabling broader real-world personalized applications. We conduct thorough experiments across various benchmarks and large VLMs to demonstrate the effectiveness of the proposed SLC framework. The code will be released.

1 INTRODUCTION

Recently, personalizing Vision-Language Models (Wang et al., 2024b; Liu et al., 2023; Hurst et al., 2024) (VLMs) that assist users’ daily life has gained increasing interest. For example, VLMs should be aware of user-provided concepts such as a pet and generate personalized output including concept identifiers such as $\langle \text{Bob} \rangle$ or $\langle \text{Lina} \rangle$. To seamlessly integrate concepts into VLMs, many techniques make different attempts, including finetuning-based (Alaluf et al., 2024; Nguyen et al., 2024; 2025; An et al., 2024; 2025), Retrieval-Augmented Generation (RAG)-based (Hao et al., 2025), and representation-based (Pi et al., 2024) methods. While effective, deploying them at scale for all users presents critical challenges due to huge training costs (see Figure 1), particularly for large models. Taking finetuning-based methods as an example, existing approaches primarily focus on fine-tuning open-source VLM for personalization. However, it brings high training costs for fine-tuning large

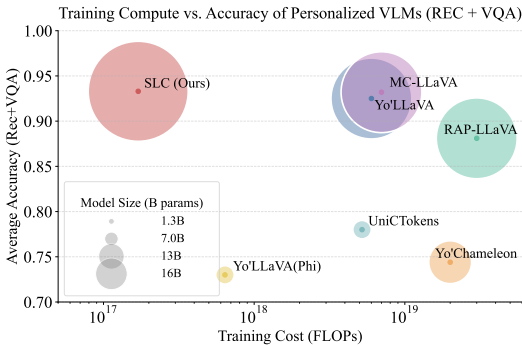


Figure 1: **Cost-size-accuracy of personalized VLMs.** Bubbles plot total training cost (x) vs. mean Rec.(Recognition)+VQA accuracy on the Yo’LLaVA dataset (y); bubble area scales with model size. Our SLC (MetaC-3B + LLaVA-13B version) achieves the best accuracy (0.933) with $\sim 10^2\times$ lower training cost than 13B models.

054 models. Furthermore, leading companies like OpenAI limit access through paid APIs, also challeng-
055 ing existing paradigms. In contrast, small models can be easily personalized and are freely available,
056 but they have limited reasoning ability.

057 Thus, a new paradigm is necessary, raising an important question: can we synthesize the benefits of
058 both large and small models while mitigating their respective limitations?

060 While the small-large model collaboration paradigm is under-explored in the field of VLM person-
061 alization, several works (Chinchali et al., 2021; Zhang et al., 2022; Chen et al., 2025a; Li et al.,
062 2025) have implemented it in other fields. However, most (Wang et al., 2024a; Ding et al., 2024)
063 focus on reducing inference rather than training cost, which is a core bottleneck for VLM person-
064 alization. Wu et al. (2025) also proposed an insightful paradigm that small models solely handle the
065 generation of chain-of-thought and large models only take responsibility for reasoning. However, in
066 the context of VLM personalization, the model often produces hallucinations, and simply combining
067 the outputs of the small and large models may not be effective.

068 Considering the above-mentioned challenges, and to support personalization for closed-source mod-
069 els, we propose a novel paradigm called Small-Large Collaboration (SLC) for personalizing large
070 VLMs using small VLMs, illustrated in Figure 2.

071 Our key insight involves a clear division of tasks: small models manage user-specific perception,
072 while large models facilitate reflection and general reasoning. To reduce training costs, we strate-
073 gically train a meta personalized small VLM that can adaptively output various personalized in-
074 formation without the need for tuning. We pre-trained a small set of LoRA adapters for the small
075 model; during inference, the most relevant adapters are activated dynamically, enabling zero-shot
076 personalization for new user concepts. By leveraging the advantages of this paradigm, we can ef-
077 fectively combine small and large models, using small models to support the personalization of both
078 open-source and closed-source large models. Additionally, our modular design naturally supports
079 privacy-preserving hybrid architectures, in which lightweight local models perform personalized
080 detection and interact securely with powerful cloud-hosted reasoning models. However, limited by
081 model capability, small models are prone to output hallucinations, which may confuse concepts.
082 Thus, we propose a reflection mechanism that the large model re-queries each detected concept with
083 two focused yes/no VQA checks and suppresses any evidence it cannot visually confirm. Combined
084 together, these two components allow SLC to fully leverage the reasoning power of large VLMs
085 while also gaining the deployment advantages of small VLMs.

086 We summarize our contributions as follows:

- 087 • We propose a novel and efficient paradigm to personalize large VLMs using small VLMs,
088 supporting both open-source and closed-source large VLMs personalization.
- 089 • We train a meta personalized small VLM and design a test-time reflection mechanism to
090 reduce training costs and minimize potential hallucinations.
- 091 • Extensive experiments show that SLC achieves competitive performance across VLM per-
092 sonalization benchmarks, paving the way for real-world applications.

095 2 RELATED WORK

097 **Vision-Language Models** Despite Vision Language Models (VLMs) (Bai et al., 2025; Hurst et al.,
098 2024; Comanici et al., 2025) have demonstrated remarkable capabilities, the widespread adoption of
099 VLMs for personalized applications is hampered by challenges. The sheer size and computational
100 requirements of these models make it difficult and expensive to customize them for each user. To
101 address the challenges of efficient deployment, researchers in the VLM field have explored several
102 avenues. For instance, Qwen2.5-VL-3B (Bai et al., 2025) achieves lightweight through a highly op-
103 timized model architecture and a multi-stage training strategy. Meanwhile, explorations with models
104 like Phi (Microsoft et al., 2025) series using meticulously curated data build highly capable models
105 with far fewer parameters. However, these approaches are designed to create an efficient general-
106 purpose model and are not transferable to our scenario, which requires personalized customization
107 for every individual. In this work, we propose the SLC framework that pairs a small, efficient VLM
for user-specific concepts learning with a large VLM for high-quality generation. This collabora-

108 tive paradigm resolves the inherent trade-off between personalization ability and the efficiency of
 109 training and deployment.
 110

111
 112 **Personalization of VLMs** Personalization in VLMs involves adapting models to recog-
 113 nize and incorporate user-specific concepts into their outputs for personalized interactions.
 114 MyVLM (Alaluf et al., 2024) augments pre-trained VLMs with training external concept heads,
 115 while Yo’LLaVA (Nguyen et al., 2024) uses learnable soft prompts to embed concepts by fine-
 116 tuning. To overcome the limitations of a single concept, MC-LLaVA (An et al., 2024) employs
 117 joint instruction tuning to handle multiple concepts, while RAP (Hao et al., 2025) adopts retrieval-
 118 augmented generation to enrich personalized responses. However, these methods require extensive
 119 training (as shown in Figure 1), creating a trade-off between personalization abilities and training
 120 costs. Facing this challenge, we propose a small-large collaboration paradigm to balance response
 121 quality and personalization efficiency.
 122

123
 124 **Small–Large Model Collaboration** There
 125 has been growing interest in collaborations between large and small models, primarily as a
 126 strategy to balance high performance with efficiency in general-purpose scenarios. Prevailing
 127 paradigms often focus on generic performance enhancement. For instance, the pipeline
 128 approaches (Lv et al., 2025; Zhang et al., 2024) use a small model to generate initial
 129 candidates for a large model to refine. More dynamic hybrid or routing strategies (Wang
 130 et al., 2025; Ong et al., 2025; Varangot-Reille et al., 2025; Chen et al., 2025b) employ a
 131 router to intelligently delegate tasks based on complexity. Similarly, auxiliary/enhancement
 132 paradigms (Shao et al., 2025) have been developed where one model assists another to
 133 improve overall performance. A particularly inspiring approach is the Cache of Thought
 134 (CoT) framework (Wu et al., 2025) that leverages a large model’s knowledge to efficiently
 135 empower a small model. Whereas existing approaches have focused only on enhancing
 136 performance in general-purpose scenarios, our work is inspired by the CoT framework to ad-
 137 dress the challenge of personalization. In our paradigm, a lightweight small model acts as a
 138 “personalizer”, detecting user-specific concepts and providing information to a large model to
 139 enable the generation of higher-quality, personalized responses.
 140
 141
 142
 143
 144
 145
 146
 147
 148
 149
 150
 151
 152

153
 154 **3 METHOD**
 155
 156

157 The method section is organized into two parts:
 158 1) We first formalize the task of personalizing
 159 VLMs and state the system’s objectives. We
 160 then outline our Small–Large Collaboration (SLC) inference pipeline, showing how the small VLM
 161 \mathcal{M}_s and the large VLM \mathcal{M}_l collaborate at test time. 2) We present the details of a meta-personalized
 small VLM and discuss the test-time reflection of a large VLM.

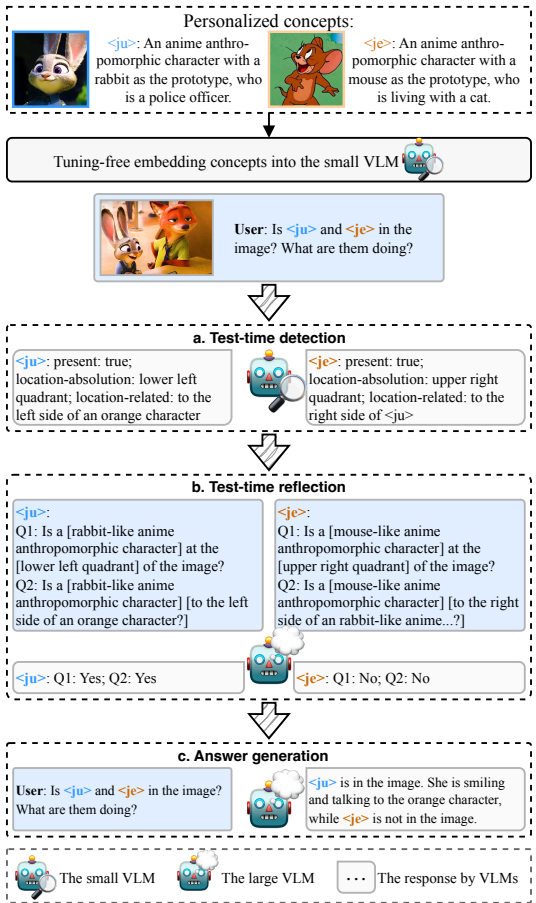


Figure 2: **Inference pipeline of SLC:** a. Test-time detection by the small VLM \mathcal{M}_s ; b. Test-time reflection by the large VLM \mathcal{M}_l ; c. Answer generation by the large VLM \mathcal{M}_l .

3.1 PROBLEM SETUP

The user registers a concept set $\{C_i^u\}_{i=1}^N$ —for example, $\langle \text{Bo} \rangle$ and $\langle \text{Lina} \rangle$. Each C_i^u with reference images $\mathcal{I}_{C_i^u}$ and a brief textual description $\mathcal{T}_{C_i^u}$ (e.g., “ $\langle \text{Bo} \rangle$ is a golden-retriever dog; it is my first pet dog.”). At every dialogue turn, given an image–question pair (I_t, q_t) , the system must **(i)** identify which C_i^u appear in I_t , and **(ii)** answer any question about them.

3.2 SLC INFERENCE PIPELINE

3.2.1 OVERVIEW

Figure 2 sketches SLC’s three-stage pipeline. In summary, our SLC framework addresses personalized VLM reasoning for large VLMs by utilizing small VLMs. A small model (\mathcal{M}_s) rapidly embeds user-defined concepts and produces structured, concept-level cues, while a powerful large model (\mathcal{M}_l) verifies those cues at test time and generates the final answer. Concretely, the collaboration consists of three sequential stages:

a. Test-time detection. At the beginning of each dialogue turn, the lightweight, meta-trained small VLM \mathcal{M}_s embeds all registered user concepts and then examines I_t , providing concept-level cues $R_t = \{r_{t,i}\}_{i=1}^N$, where $r_{t,i} = \{\text{present}, \text{loc}_{\text{abs}}, \text{loc}_{\text{rel}}\}$.

b. Test-time reflection. For every C_i^u detected by \mathcal{M}_s , the large VLM \mathcal{M}_l performs self-VQA checks to verify claims and sanitize cues, yielding \tilde{R}_t .

c. Answer generation. Finally, \mathcal{M}_l takes (I_t, q_t, \tilde{R}_t) and generates the final response a_t to the user.

3.2.2 TEST-TIME REFLECTION OF LARGE VLM

Concept-Level Cue Generation To prime the test-time reflection step, \mathcal{M}_s first emits a set of structured concept-level cues for the current step t :

$$R_t = \{r_{t,i}\}_{i=1}^N, \quad r_{t,i} = \{\text{present}, \text{loc}_{\text{abs}}, \text{loc}_{\text{rel}}\}.$$

Here N is the number of user-registered concepts, and each triple $r_{t,i}$ contains present — a Boolean flag indicating whether concept C_i^u is detected in image I_t ; loc_{abs} — a string describing the concept’s absolute position in the image (e.g., “upper-left corner”, “center”); loc_{rel} — a string giving the concept’s position relative to salient objects (e.g., “to the right of the cat”, “behind the car”). When run in isolation, the small VLM \mathcal{M}_s may hallucinate, confusing similar textures or backgrounds with a registered concept (see Figure 2). During test-time reflection, \mathcal{M}_l cross-checks each $r_{t,i}$ against the image and the concept description $\mathcal{T}_{C_i^u}$, producing a refined cue set \tilde{R}_t with fewer hallucinations.

Identity Extraction and Verification For each personalized concept C_i^u detected by \mathcal{M}_s in R_t , the large model \mathcal{M}_l first extracts an immutable identity phrase $\text{ID}(C_i^u)$ from its textual description $\mathcal{T}_{C_i^u}$ (e.g., “a rabbit-like anime character”). It then performs two binary verifications: whether $\text{ID}(C_i^u)$ appears at the absolute location loc_{abs} in the image, and whether it is consistent with the relative relation loc_{rel} . The resulting answers $(a_1, a_2) \in \{\text{yes}, \text{no}\}^2$ are used to refine the cue $r_{t,i}$, where a double “no” sets present = 0, a single negative discards the corresponding absolute or relative location, and affirmative responses retain $r_{t,i}$ unchanged. This yields the refined cue set $\tilde{R}_t = \{\tilde{r}_{t,i}\}$, which is subsequently fused with (I_t, q_t) for final answer generation by \mathcal{M}_l .

3.3 META-PERSONALIZED SMALL VLM

Current fine-tuning practices for different concepts lead to linear increases in training time and storage, which restrict the large-scale deployment of personalized applications. To tackle this issue, we propose a meta-training strategy for personalizing VLMs, inspired by the adaptation of T2I models (Rombach et al., 2022; Topal et al., 2025; Ruiz et al., 2024). The small VLM \mathcal{M}_s is trained once offline; at test time, it incorporates concepts by simply selecting pre-learned adapters—no optimization necessary, which further minimizes overall training expenses.

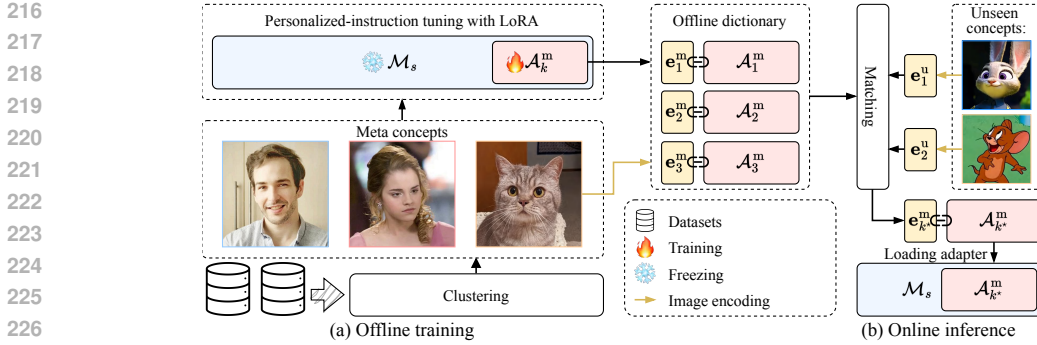


Figure 3: **Overview of our meta-personalization framework for the small VLM.** (a) During offline training, we cluster concept images into K meta concepts and apply LoRA tuning to train K corresponding adapters $\{\mathcal{A}_k^m\}_{k=1}^K$. Each meta concept is encoded as (e_k^m, \mathcal{A}_k^m) and is stored offline. (b) At test time, a new concept is encoded into e_j^u and matched to its closest meta concept. The matched adapter is loaded into \mathcal{M}_s for the downstream task, enabling tuning-free inference.

3.3.1 OFFLINE TRAINING

As shown in Figure 3(a), we extract CLIP embeddings for images in several public personalization datasets and run K-means clustering to form K appearance clusters. The centroid of each cluster defines a meta-concept C_k^m , yielding the set $\{C_k^m\}_{k=1}^K$. These meta-concepts span a broad semantic spectrum—covering humans, animals, and diverse objects—and can be adapted or combined to represent new, semantically related ideas within the same category. For each centroid, we train a Low-rank adaptation (LoRA) (Hu et al., 2022) adapter \mathcal{A}_k^m using the personalized-instruction tuning recipe from prior work (An et al., 2024) (training details in the Appendix). We store every meta-concept embedding e_k^m together with its adapter, forming an offline dictionary.

3.3.2 ONLINE INFERENCE

At run time, the user may register multiple new concepts $\{C_i^u\}_{i=1}^N$. As Figure 3(b) illustrates, we first compute an embedding for each concept by averaging its reference-image features, then form a scenario embedding $\bar{e}^u = \frac{1}{N} \sum_{i=1}^N e_i^u$. We select a single meta-adapter for the scenario via the cosine rule

$$k^* = \arg \max_k \cos(\bar{e}^u, e_k^m). \quad (1)$$

The chosen adapter $\mathcal{A}_{k^*}^m$ is plugged into \mathcal{M}_s and used to detect all registered concepts. Since no weights are updated online, the pipeline remains tuning-free while scaling to an unlimited number of new concepts.

4 EXPERIMENT

4.1 EXPERIMENTAL SETUP

Evaluation Datasets All experiments are conducted on the MC-LLaVA (An et al., 2024) and Yo’LLaVA (Nguyen et al., 2024) datasets. MC-LLaVA is a multi-concept personalization dataset for VLMs comprising 118 diverse concepts that are grouped into 40 scenarios by data source. The concepts span real-world personalities, anime characters, real objects, and anime-style objects. MC-LLaVA provides both single- and multi-concept test sets for recognition, VQA, and text-only QA. Yo’LLaVA contains 40 distinct concepts covering objects, buildings, and people, and supplies single-concept test sets for the same three tasks.

We observed that several personalization methods tend to overfit to a concept’s reference images, paying less attention to the current visual input and occasionally hallucinating. To probe this issue, we additionally construct a Special Question–Answer (SQA) set: targeted VQA queries on Yo’LLaVA test images whose correct answers contradict spurious cues present in the concept’s training data, forcing the model to ground its reasoning on the image at hand. Representative SQA examples are provided in the Appendix.

Implementation Details To build the meta-concept dictionary, we merge all concepts in the Yo’LLaVA and MC-LLaVA training splits. These concepts are clustered into $K = 10$ meta-concepts (see Appendix for details). To avoid data leakage, any concepts (along with images and Q&As) associated with the meta-concept are removed from every evaluation set. We adopt Qwen2.5-VL-3B-Instruct (Bai et al., 2025) as the backbone of \mathcal{M}_s and meta-train it to obtain our Meta-Concepts-Model-3B (MetaC-3B). Each meta-concept adapter is optimized with LoRA for 80 steps using a learning rate of 5×10^{-5} and a batch size of 64 on 8 A800 GPUs. To reduce randomness, we run the experiment three times and report the average.

Table 1: **Performance comparison of personalized VLMs on Yo’LLaVA and MC-LLaVA datasets.** For SLC, $\mathcal{M}_s = \text{MetaC-3B}$. The **best** and **second-best** performances are highlighted.

Method	Yo’LLaVA dataset					MC-LLaVA dataset				
	Training cost	Rec.	VQA	Text-only	SQA	Rec.		VQA	Text-only	
	FLOPs	Weight	Acc	Acc	Acc	Single	Multi	Weight	Acc	Acc
Image prompt + GPT-4o	-	0.901	0.915	0.891	0.850	0.831	0.823	0.827	0.904	0.733
Text prompt + GPT-4o	-	0.872	0.930	0.871	0.900	0.746	0.822	0.781	0.889	0.702
SLC ($\mathcal{M}_l = \text{GPT-4o}$)	1.7×10^{17}	0.951	0.979	0.895	0.900	0.760	0.931	0.830	0.937	0.739
SLC ($\mathcal{M}_l = \text{LLaVA-1.5-13B}$)	1.7×10^{17}	0.895	<u>0.971</u>	0.879	<u>0.883</u>	<u>0.762</u>	<u>0.878</u>	0.801	0.861	0.692
MC-LLaVA	7.0×10^{18}	<u>0.947</u>	0.934	<u>0.885</u>	0.725	0.912	0.845	0.878	<u>0.890</u>	<u>0.709</u>
Yo’LLaVA	6.0×10^{18}	0.924	0.929	0.883	0.713	0.744	0.729	0.737	0.655	0.658
RAP-LLaVA	3.0×10^{19}	0.845	0.917	0.874	0.813	0.747	0.688	0.713	0.784	0.685
LLaVA + text prompt	-	0.819	0.913	0.803	0.725	0.594	0.549	0.573	0.817	0.553

4.2 PERFORMANCE COMPARISON

4.2.1 COMPARED METHODS

We benchmark SLC against several representative personalization approaches for VLM on both the Yo’LLaVA and MC-LLaVA test sets:

- **Yo’LLaVA** (Nguyen et al., 2024): one of the earliest VLM personalization methods (built on LLaVA-1.5-13B). Because it natively supports only single-concept scenarios, we follow the multi-concept adaptation Yo’LLaVA-M in MC-LLaVA (An et al., 2024).
- **MC-LLaVA** (An et al., 2024): a VLM personalization method specifically designed for multi-concept personalization, also based on LLaVA-1.5-13B.
- **RAP-LLaVA** (Hao et al., 2025): a retrieval-augmented generation (RAG) approach for multimodal personalization, again built on LLaVA-1.5-13B.
- **Upper bound (GPT-4o)** (OpenAI et al., 2024): Personalized image or text prompts with test questions are fed to GPT-4o, serving as an optimistic performance ceiling thanks to its strong multimodal reasoning capacity.
- **Lower bound (LLaVA + text prompt)** (Liu et al., 2023): Each test question is paired with its corresponding personalized text prompt and evaluated on the vanilla LLaVA-1.5-13B, providing a conservative baseline.

For a fair comparison, we instantiate SLC with two large VLMs, LLaVA-1.5-13B and GPT-4o. During text-only QA, no image is supplied. Instead, the question is concatenated with the corresponding concept descriptions, and \mathcal{M}_l directly generates the answer. Further implementation details for the text-only QA protocol are deferred to the Appendix.

4.2.2 RESULT ANALYSIS

Table 1 reports the overall comparison. Our evaluation metrics follow the protocols of the Yo’LLaVA (Nguyen et al., 2024) and MC-LLaVA (An et al., 2024) datasets. In particular, for the recognition task we adopt a weighted score defined as $\text{Weighted} = 0.5 \times \text{Yes recall} + 0.5 \times \text{No recall}$, where *Yes recall* denotes the proportion of correctly predicted “yes” responses among existence queries with the concept present, and *No recall* analogously measures the proportion of correctly predicted “no” responses when the concept is absent. For the Yo’LLaVA dataset, we additionally

list the training FLOPs consumed by each method, so that accuracy can be assessed jointly with computational cost. Below, we highlight the superiority of SLC from three perspectives.

(a) Superior accuracy without task-specific finetuning. Across the two SLC variants, nearly every first- or second-place in Table 1 is occupied. SLC with \mathcal{M}_l =GPT-4o attains the top scores on all Yo’LLaVA metrics (0.951 Rec., 0.979 VQA, 0.895 Text-only) and secures either the highest or runner-up performance on MC-LLaVA, surpassing every other GPT-4o-based baseline. When \mathcal{M}_l is replaced by LLaVA-1.5-13B, SLC still matches—or exceeds—prior finetuned methods, verifying that our meta-trained small model plus test-time reflection can unleash the strong reasoning power of large VLMs without additional finetuning.

(b) Training efficiency. SLC performs only one meta-training on \mathcal{M}_s , consuming 1.7×10^{17} FLOPs—about $40\times$ less than the fine-tune-heavy Yo’LLaVA (6.0×10^{18}) and MC-LLaVA (7.0×10^{18}), and almost $200\times$ less than the retrieval-augmented RAP-LLaVA (3.0×10^{19}). For Yo’LLaVA and MC-LLaVA, the training cost grows nearly linearly with the number of personalized concepts, whereas SLC’s training cost is fixed after a single meta-training stage. Although RAP-LLaVA’s expense is a one-off rather than linear, it remains orders of magnitude higher than SLC. These results confirm that SLC is more training-efficient.

(c) Reduced over-fitting and hallucination. On the SQA set—designed to expose memorization—SLC ties GPT-4o for the top score (0.900) and exceeds all finetuned methods by >10 pp. The gain attests to the synergy between our two-model collaboration and the test-time reflection step: \mathcal{M}_s supplies structured concept cues, while \mathcal{M}_l subsequently verifies those cues, suppressing spurious matches and grounding the answer in the current image.

These results indicate that a lightweight meta-personalized VLM \mathcal{M}_s plus a powerful but frozen large VLM \mathcal{M}_l is both more accurate and far cheaper to train than existing finetune-heavy pipelines.

4.3 ABLATION STUDY AND ANALYSIS

4.3.1 EFFECT OF LORA TRAINING AND PROMPTING STRATEGIES

To test how different prompt configurations during LoRA training and inference affect generalization to unseen concepts, we train adapters for \mathcal{M}_s = Qwen2.5-VL-3B (Bai et al., 2025) on the Yo’LLaVA dataset.

During training, we prepend each question prompt in the QA pair with the concept description or leave it unchanged. At inference, we mirror this choice, yielding six settings in total, plus two baselines without LoRA. Results are reported in Table 2.

It is obvious that using concept descriptions consistently at both training and inference stages yields the best overall performance on recognition, VQA, and text-only QA (0.819 / 0.835 / 0.843). When the description is included only during training, accuracy drops significantly at test time, likely due to a mismatch between training and inference prompts. In contrast, models trained without descriptions can still benefit from their inclusion at inference, though they fail to match the jointly prompted setting. Notably, all LoRA-trained variants outperform the non-adapted baselines, highlighting the effectiveness of lightweight adaptation. These results collectively indicate that prompt consistency is crucial for generalization to unseen concepts.

4.3.2 EFFECT OF META-CONCEPT POOL SIZE AND TOP-K ADAPTER SELECTION

To investigate how the meta-concept pool size and Top- K adapter fusion influence the performance of SLC, we evaluate two configurations on the Yo’LLaVA dataset recognition task: **(i)** the small model alone (\mathcal{M}_s = MetaC-3B) and **(ii)** the complete SLC pipeline ($\mathcal{M}_s + \mathcal{M}_l$). For each setting, we vary **(a)** the number of meta-concepts obtained via k -means clustering of CLIP embeddings, and **(b)** the number of nearest meta-concept adapters Top- K averaged at test time. Each meta-concept

Table 2: **Prompt strategies for LoRA training and inference on the Yo’LLaVA dataset.** TP = text prompt. The **best** performances are highlighted.

Training Setting	Inference Setting	Rec. VQA Text-only		
		Weight	Acc	Acc
w/o LoRA	w/ TP.	0.714	0.813	0.810
	w/o TP.	0.508	0.786	0.655
w/ TP.	w/ TP.	0.819	0.835	0.843
	w/o TP.	0.603	0.791	0.697
w/o TP.	w/ TP.	0.744	0.818	0.733
	w/o TP.	0.562	0.790	0.702

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

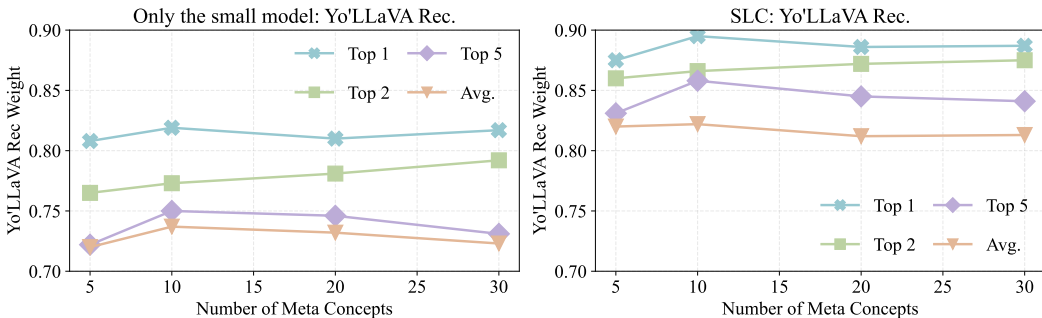


Figure 4: **Impact of meta-concept pool size and Top- K adapter fusion on Yo’LLaVA recognition.** **Left:** the small model only (\mathcal{M}_s). **Right:** full SLC pipeline ($\mathcal{M}_s + \mathcal{M}_l$). Across all $K \in \{1, 2, 5, \text{all}\}$, the best accuracy occurs at 10 meta-concepts, and Top-1 selection consistently yields the highest scores.

has its own LoRA adapter; at inference, we retrieve the Top- K closest meta-concepts to the queried concept and merge their adapters before decoding.

Figure 4 reveals two consistent trends. **First**, recognition peaks at 10 meta-concepts for both \mathcal{M}_s and SLC, indicating an optimal balance between coverage and redundancy. **Second**, selecting a single adapter (Top-1) consistently surpasses Top-2 and Top-5, showing that fusing multiple adapters dilutes concept-specific cues. These findings validate our design principle: a moderately sized (10) adapter library combined with dynamic Top-1 selection yields the strongest performance.

Table 3: **Varying the small model \mathcal{M}_s while fixing $\mathcal{M}_l = \text{LLaVA-1.5-13B}$.** “Yes/No recall” measures the \mathcal{M}_s ’s ability to correctly predict the presence/absence of a concept, respectively. Adding \mathcal{M}_l improves all metrics and markedly boosts “No recall”, evidencing hallucination reduction.

\mathcal{M}_s	VLM	Training data	Only \mathcal{M}_s				$\mathcal{M}_s + \mathcal{M}_l$			
			Rec.		VQA		Rec.		VQA	
			Yes recall	No recall	Weight	Acc	Yes recall	No recall	Weight	Acc
MetaC-3B	Qwen2.5VL-3B	47.3k	0.900	0.759	0.829	0.835	0.898	0.893	0.895	0.971
RAP-Phi3-V	Phi3-V-3.8B	260k	0.912	0.754	0.833	0.866	0.906	0.896	0.901	0.971
RAP-LLaVA	LLaVA-1.5-13B	260k	0.926	0.764	0.845	0.917	0.921	0.920	0.921	0.975
Yo’LLaVA	LLaVA-1.5-13B	190k	0.949	0.898	0.924	0.929	0.947	0.928	0.936	0.971
LLaVA + text prompt	LLaVA-1.5-13B	-	0.734	0.903	0.819	0.913	0.734	0.929	0.832	0.966

4.3.3 VARYING THE SMALL MODEL

We further investigate how the capacities of \mathcal{M}_s influence system behavior on the Yo’LLaVA dataset (Nguyen et al., 2024). To obtain more diverse results, we select five personalized models with different training and inference paradigms as the small model \mathcal{M}_s , each emphasizing distinct aspects: **MetaC-3B**: a meta-personalized model trained and inferred with an efficient meta-learning strategy; **RAP-Phi3-V** and **RAP-LLaVA** (Hao et al., 2025): two models sharing the same training data and RAG-based inference, but differing in the sizes of their VLM backbones; **Yo’LLaVA** (Nguyen et al., 2024): a conventional personalized model that trains and evaluates separately for each concept; **LLaVA+text prompt**: a training-free baseline built on LLaVA-1.5 (Liu et al., 2024), as detailed in Section 4.2.1. The experimental results are presented in Table 3. Across different choices of \mathcal{M}_s , we observe varying tendencies in Recognition: trained models generally achieve higher “Yes recall”, while the training-free “LLaVA+text prompt” baseline tends to favor “No” predictions. For the first four trained models, false positives are relatively common. When comparing the performance of \mathcal{M}_s alone versus $\mathcal{M}_s + \mathcal{M}_l$, two key findings emerge. **First**, \mathcal{M}_l substantially reduces false positives, as evidenced by the consistent improvement in “No recall”. **Second**, the overall performance of SLC is positively correlated with the intrinsic performance of \mathcal{M}_s (as reflected in both Recognition and VQA), while being less sensitive to the specific training or inference paradigm adopted.

4.3.4 SCALING THE LARGE MODEL

We explore how the capacities of \mathcal{M}_l affect SLC behavior on the Yo’LLaVA dataset (Nguyen et al., 2024). Table 4 reports the impact of varying the size of \mathcal{M}_l (ranging from no \mathcal{M}_l to 3B and up to 72B) on the performance of SLC. We observe that incorporating \mathcal{M}_l substantially improves the performance of the small model on both Recognition and VQA tasks. Moreover, as the capacity of \mathcal{M}_l increases, the overall performance of SLC continues to improve. This clearly demonstrates that, although \mathcal{M}_s has limited capability, \mathcal{M}_l can effectively correct the errors made by \mathcal{M}_s , highlighting a scaling law that characterizes the steady gains of our SLC framework with respect to \mathcal{M}_l .

4.3.5 ABLATING THE SLC PIPELINE

To pinpoint the roles of R_t produced by \mathcal{M}_s and the test-time reflection, we compare SLC against three ablated variants: w/o reflection, where \mathcal{M}_s is enabled but \mathcal{M}_l skips verification; w/o detector, where \mathcal{M}_l performs reflection without any cues from \mathcal{M}_s ; and Pure VLM, where both components are disabled, reducing SLC to “LLaVA + text prompt” in Table 1. Table 5 yields three findings:

(1) Synergy is essential. The full SLC tops nearly every metric. Removing either element reduces performance (Yo’LLaVA VQA drops to 0.958 w/o reflection and 0.962 w/o \mathcal{M}_s), confirming that the two modules are essential. **(2) Reflection suppresses hallucination.** The “No recall” rate—correctly predicting a concept’s absence—jumps from 0.814 (w/o reflection) to 0.893 with reflection, and reaches 0.905 when only reflection is present. Thus, \mathcal{M}_l ’s verification is crucial for filtering false positives. **(3) \mathcal{M}_s supplies precise cues.** Omitting \mathcal{M}_s hurts “Yes recall” (0.898 \rightarrow 0.767) and pushes MC-LLaVA recognition weight down from 0.801 to 0.710, showing that the efficiency of \mathcal{M}_s .

Table 4: **Scaling \mathcal{M}_l while fixing $\mathcal{M}_s = \text{MetaC-3B}$.**

\mathcal{M}_l	Model Size	Rec.	VQA	Text-only
		Weight	Acc	Acc
-	-	0.829	0.835	0.843
Qwen2.5VL	3B	0.872	0.956	0.844
Qwen2.5VL	7B	0.903	0.973	0.879
Qwen2.5VL	32B	0.932	0.975	0.881
Qwen2.5VL	72B	0.944	0.979	0.890

Table 5: **Ablation of SLC components.** We fix $\mathcal{M}_s = \text{MetaC-3B}$, $\mathcal{M}_l = \text{LLaVA-1.5-13B}$, and selectively disable \mathcal{M}_s or the test-time reflection. $\checkmark = \text{enabled}$, $- = \text{disabled}$.

\mathcal{M}_s	Test-time reflection	Yo’LLaVA dataset					MC-LLaVA dataset			
		Rec.		VQA	SQA	Rec.		VQA		
		Yes recall	No recall	Weight	Acc	Acc	Single	Multi	Weight	Acc
\checkmark	\checkmark	0.898	0.893	0.895	0.971	0.883	0.762	0.878	0.801	0.861
\checkmark	-	0.841	0.814	0.827	0.958	0.838	0.748	0.642	0.705	0.843
-	\checkmark	0.767	0.905	0.836	0.962	0.825	0.712	0.707	0.710	0.839
-	-	0.734	0.903	0.819	0.913	0.725	0.594	0.549	0.573	0.817

5 CONCLUSION

We propose SLC, a novel Small–Large Collaboration paradigm for personalizing VLMs, effectively balancing training efficiency and personalization capability. SLC combines a meta-trained, tuning-free small model for personalized cues generation with a large model performing test-time reflection to reduce hallucinations. This approach addresses the longstanding cost–performance trade-off and naturally supports personalization of both open-source and closed-source large VLMs. Its modular design further enables privacy-preserving hybrid deployments. Extensive experiments validate SLC’s scalability, efficiency, and reliability, highlighting its potential for real-world applications.

REFERENCES

- 486
487
488 Yuval Alaluf, Elad Richardson, Sergey Tulyakov, Kfir Aberman, and Daniel Cohen-Or. Myvlm:
489 Personalizing vlms for user-specific queries. In *European Conference on Computer Vision*, pp.
490 73–91. Springer, 2024.
- 491 Ruichuan An, Sihan Yang, Ming Lu, Renrui Zhang, Kai Zeng, Yulin Luo, Jiajun Cao, Hao Liang,
492 Ying Chen, Qi She, et al. Mc-llava: Multi-concept personalized vision-language model. *arXiv*
493 *preprint arXiv:2411.11706*, 2024.
- 494
495 Ruichuan An, Sihan Yang, Renrui Zhang, Zijun Shen, Ming Lu, Gaole Dai, Hao Liang, Ziyu Guo,
496 Shilin Yan, Yulin Luo, et al. Unictokens: Boosting personalized understanding and generation
497 via unified concept tokens. *arXiv preprint arXiv:2505.14671*, 2025.
- 498 Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Siboz Song, Kai Dang, Peng Wang,
499 Shijie Wang, Jun Tang, Humen Zhong, Yanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan,
500 Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, Jiabo Ye, Xi Zhang, Tianbao Xie, Zesen Cheng,
501 Hang Zhang, Zhibo Yang, Haiyang Xu, and Junyang Lin. Qwen2.5-vl technical report, 2025.
502 URL <https://arxiv.org/abs/2502.13923>.
- 503
504 Yi Chen, JiaHao Zhao, and HaoHao Han. A survey on collaborative mechanisms between large and
505 small language models, 2025a. URL <https://arxiv.org/abs/2505.07460>.
- 506 Zhijun Chen, Jingzheng Li, Pengpeng Chen, Zhuoran Li, Kai Sun, Yuankai Luo, Qianren Mao,
507 Dingqi Yang, Hailong Sun, and Philip S Yu. Harnessing multiple large language models: A
508 survey on llm ensemble. *arXiv preprint arXiv:2502.18036*, 2025b.
- 509
510 Sandeep Chinchali, Apoorva Sharma, James Harrison, Amine Elhafsi, Daniel Kang, Evgenya Perga-
511 ment, Eyal Cidon, Sachin Katti, and Marco Pavone. Network offloading policies for cloud
512 robotics: a learning-based approach. *Autonomous Robots*, 45(7):997–1012, 2021.
- 513 Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit
514 Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, et al. Gemini 2.5: Pushing the
515 frontier with advanced reasoning, multimodality, long context, and next generation agentic capa-
516 bilities. *arXiv preprint arXiv:2507.06261*, 2025.
- 517
518 Yucheng Ding, Chaoyue Niu, Fan Wu, Shaojie Tang, Chengfei Lyu, and Guihai Chen. Enhancing
519 on-device llm inference with historical cloud-based llm interactions. In *Proceedings of the 30th*
520 *ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pp. 597–608, 2024.
- 521 Haoran Hao, Jiaming Han, Changsheng Li, Yu-Feng Li, and Xiangyu Yue. Rap: Retrieval-
522 augmented personalization for multimodal large language models. In *Proceedings of the Com-*
523 *puter Vision and Pattern Recognition Conference*, pp. 14538–14548, 2025.
- 524
525 Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang,
526 Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3, 2022.
- 527
528 Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Os-
529 trow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv preprint*
530 *arXiv:2410.21276*, 2024.
- 531 Senyao Li, Haozhao Wang, Wenchao Xu, Rui Zhang, Song Guo, Jingling Yuan, Xian Zhong,
532 Tianwei Zhang, and Ruixuan Li. Collaborative inference and learning between edge slms and
533 cloud llms: A survey of algorithms, execution, and open challenges, 2025. URL <https://arxiv.org/abs/2507.16731>.
- 534
535 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances*
536 *in neural information processing systems*, 36:34892–34916, 2023.
- 537
538 Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction
539 tuning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*,
pp. 26296–26306, 2024.

- 540 Zheqi Lv, Tianyu Zhan, Wenjie Wang, Xinyu Lin, Shengyu Zhang, Wenqiao Zhang, Jiwei Li, Kun
541 Kuang, and Fei Wu. Collaboration of large language models and small recommendation models
542 for device-cloud recommendation. *arXiv preprint arXiv:2501.05647*, 2025.
543
- 544 Microsoft, :, Abdelrahman Abouelenin, Atabak Ashfaq, Adam Atkinson, Hany Awadalla, Nguyen
545 Bach, Jianmin Bao, Alon Benhaim, Martin Cai, Vishrav Chaudhary, Congcong Chen, Dong Chen,
546 Dongdong Chen, Junkun Chen, Weizhu Chen, Yen-Chun Chen, Yi ling Chen, Qi Dai, Xiyang Dai,
547 Ruchao Fan, Mei Gao, Min Gao, Amit Garg, Abhishek Goswami, Junheng Hao, Amr Hendy,
548 Yuxuan Hu, Xin Jin, Mahmoud Khademi, Dongwoo Kim, Young Jin Kim, Gina Lee, Jinyu Li,
549 Yunsheng Li, Chen Liang, Xihui Lin, Zeqi Lin, Mengchen Liu, Yang Liu, Gilsinia Lopez, Chong
550 Luo, Piyush Madan, Vadim Mazalov, Arindam Mitra, Ali Mousavi, Anh Nguyen, Jing Pan, Daniel
551 Perez-Becker, Jacob Platin, Thomas Portet, Kai Qiu, Bo Ren, Liliang Ren, Sambuddha Roy,
552 Ning Shang, Yelong Shen, Saksham Singhal, Subhojit Som, Xia Song, Tetyana Sych, Praneetha
553 Vaddamanu, Shuohang Wang, Yiming Wang, Zhenghao Wang, Haibin Wu, Haoran Xu, Weijian
554 Xu, Yifan Yang, Ziyi Yang, Donghan Yu, Ishmam Zabir, Jianwen Zhang, Li Lyna Zhang, Yunan
555 Zhang, and Xiren Zhou. Phi-4-mini technical report: Compact yet powerful multimodal language
556 models via mixture-of-loras, 2025. URL <https://arxiv.org/abs/2503.01743>.
- 557 Thao Nguyen, Haotian Liu, Yuheng Li, Mu Cai, Utkarsh Ojha, and Yong Jae Lee. Yo’llava: Your
558 personalized language and vision assistant. *Advances in Neural Information Processing Systems*,
559 37:40913–40951, 2024.
- 560 Thao Nguyen, Krishna Kumar Singh, Jing Shi, Trung Bui, Yong Jae Lee, and Yuheng Li.
561 Yo’chameleon: Personalized vision and language generation. In *Proceedings of the Computer
562 Vision and Pattern Recognition Conference*, pp. 14438–14448, 2025.
- 563 Isaac Ong, Amjad Almahairi, Vincent Wu, Wei-Lin Chiang, Tianhao Wu, Joseph E Gonzalez,
564 M Waleed Kadous, and Ion Stoica. Routellm: Learning to route llms with preference data, 2024.
565 URL <https://arxiv.org/abs/2406.18665>, 4, 2025.
566
- 567 OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Floren-
568 cia Leoni Aleman, Diogo Almeida, Janko Altmenschmidt, Sam Altman, Shyamal Anadkat, Red
569 Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Moham-
570 mad Bavarian, Jeff Belgum, Irwan Bello, Jake Beldine, Gabriel Bernadett-Shapiro, Christopher
571 Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brock-
572 man, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann,
573 Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis,
574 Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey
575 Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux,
576 Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila
577 Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix,
578 Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gib-
579 son, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan
580 Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hesse,
581 Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu,
582 Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun
583 Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Ka-
584 mali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook
585 Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel
586 Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen
587 Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel
588 Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez,
589 Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv
590 Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney,
591 Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick,
592 Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel
593 Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Ra-
jeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O’Keefe,
Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel

- 594 Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe
595 de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny,
596 Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl,
597 Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra
598 Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders,
599 Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Sel-
600 sam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor,
601 Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky,
602 Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang,
603 Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Pre-
604 ston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vi-
605 jayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan
606 Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lillian Weng,
607 Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Work-
608 man, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming
609 Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao
610 Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. Gpt-4 technical report, 2024. URL
611 <https://arxiv.org/abs/2303.08774>.
- 612 Renjie Pi, Jianshu Zhang, Tianyang Han, Jipeng Zhang, Rui Pan, and Tong Zhang. Personalized
613 visual instruction tuning. *arXiv preprint arXiv:2410.07113*, 2024.
- 614 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
615 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
616 models from natural language supervision. In *International conference on machine learning*, pp.
617 8748–8763. PmLR, 2021.
- 618 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
619 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF confer-
620 ence on computer vision and pattern recognition*, pp. 10684–10695, 2022.
- 621 Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Wei Wei, Tingbo Hou, Yael Pritch, Neal Wadhwa,
622 Michael Rubinstein, and Kfir Aberman. Hyperdreambooth: Hypernetworks for fast personaliza-
623 tion of text-to-image models. In *Proceedings of the IEEE/CVF Conference on Computer Vision
624 and Pattern Recognition (CVPR)*, pp. 6527–6536, June 2024.
- 625 Chenyang Shao, Xinyuan Hu, Yutang Lin, and Fengli Xu. Division-of-thoughts: Harnessing hy-
626 brid language model synergy for efficient on-device agents. In *Proceedings of the ACM on Web
627 Conference 2025*, pp. 1822–1833, 2025.
- 628 Barış Batuhan Topal, Umut Özyurt, Zafer Doğan Budak, and Ramazan Gokberk Cinbis. Meta-
629 lora: Meta-learning lora components for domain-aware id personalization. *arXiv preprint
630 arXiv:2503.22352*, 2025.
- 631 Clovis Varangot-Reille, Christophe Bouvard, Antoine Gourru, Mathieu Ciancone, Marion Schaeffer,
632 and François Jacquenet. Doing more with less—implementing routing strategies in large language
633 model-based systems: An extended survey. *arXiv preprint arXiv:2502.00409*, 2025.
- 634 Guanqun Wang, Jiaming Liu, Chenxuan Li, Yuan Zhang, Junpeng Ma, Xinyu Wei, Kevin Zhang,
635 Maurice Chong, Renrui Zhang, Yijiang Liu, et al. Cloud-device collaborative learning for multi-
636 modal large language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision
637 and Pattern Recognition*, pp. 12646–12655, 2024a.
- 638 Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu,
639 Jialin Wang, Wenbin Ge, et al. Qwen2-vl: Enhancing vision-language model’s perception of the
640 world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024b.
- 641 Xinyuan Wang, Yanchi Liu, Wei Cheng, Xujiang Zhao, Zhengzhang Chen, Wenchao Yu, Yanjie Fu,
642 and Haifeng Chen. Mixllm: Dynamic routing in mixed large language models. *arXiv preprint
643 arXiv:2502.18482*, 2025.

648 Mingyuan Wu, Jize Jiang, Haozhen Zheng, Meitang Li, Zhaoheng Li, Beitong Tian, Bo Chen,
649 Yongjoo Park, Minjia Zhang, Chengxiang Zhai, et al. Cache-of-thought: Master-apprentice
650 framework for cost-effective vision language model inference. *arXiv preprint arXiv:2502.20587*,
651 2025.

652 Kaiyan Zhang, Jianyu Wang, Ermo Hua, Biqing Qi, Ning Ding, and Bowen Zhou. Cogensis: A
653 framework collaborating large and small language models for secure context-aware instruction
654 following. *arXiv preprint arXiv:2403.03129*, 2024.

655
656 Tinghao Zhang, Zhijun Li, Yongrui Chen, Kwok-Yan Lam, and Jun Zhao. Edge-cloud cooperation
657 for dnn inference via reinforcement learning and supervised learning. In *2022 IEEE International*
658 *Conferences on Internet of Things (iThings) and IEEE Green Computing & Communications*
659 *(GreenCom) and IEEE Cyber, Physical & Social Computing (CPSCom) and IEEE Smart Data*
660 *(SmartData) and IEEE Congress on Cybermatics (Cybermatics)*, pp. 77–84. IEEE, 2022.

661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

A THE USE OF LLMS

During the coding and debugging phases, we utilized LLMs for technical guidance. After collaboratively drafting the manuscript, we once again turned to LLMs to enhance and refine its language and overall style.

B PROMPT TEMPLATES

B.1 TEST-TIME DETECTION PROMPT TEMPLATES

The prompt template used by the small VLM \mathcal{M}_s to detect all concepts provided by the user is shown in Table 6.

System prompt.

You are a high-precision concept detector.

Task

You should inspect the image and output **one** JSON object that **covers every** concept in the **Concept List** provided by the user while conforming to the schema below.

For each concept $\langle \text{concept_id} \rangle$ in the list:

- If *visible*, set:
 - “present”: `true`
 - “location-absolute”: “concise area, e.g. “top-left quadrant”
 - “location-relative”: “spatial relation, e.g. “to the left of the person in black suit”
- If *not visible*, set:
 - “present”: `false` \rightarrow omit all other keys

The final output should be a JSON object like: {

```

  <math>\langle \text{concept\_id\_1} \rangle</math>: {
    “present”: <boolean>,
    “location-absolute”: <string> or “”,
    “location-relative”: <string> or “”,
  },
  <math>\langle \text{concept\_id\_2} \rangle</math>: {
    “present”: <boolean>,
    “location-absolute”: <string> or “”,
    “location-relative”: <string> or “”,
  }
  ...
}
```

Rules

- Output plain English text only—no Markdown, no code fences.
- Keep every concept ID enclosed in angle brackets ($\langle \rangle$).
- If `present = false`, omit all other keys.
- Boolean literals must be lowercase `true/false`.
- Do not add any extra keys, comments, or explanatory text.

User prompt.

Concept List

$\{C_i^u : \mathcal{T}_{C_i^u}\}_{i=0}^N$

Table 6: **Test-time detection prompt** for \mathcal{M}_s .

B.2 TEST-TIME REFLECTION PROMPT TEMPLATES

Test-time reflection is performed in two concise steps by the large VLM:

- **Identity extraction** (Table 7): \mathcal{M}_l extracts each concept’s category and attributes.

- **Self-VQA verification** (Table 8): \mathcal{M}_l uses the category extracted in the previous step to answer *yes/no* questions, thereby confirming the absolute and relative locations of each concept reported as visible by \mathcal{M}_s .

System prompt.

You are an information extractor.

Task

Inspect the textual descriptions below and return **one** JSON object that **covers every** concept in the **Concept List** while conforming to the schema:

- “category”: permanent class, e.g. “a golden retriever puppy”, “a blue cartoon character”
- “attributes”: mutable traits, e.g. “always playful expression; dresses in trendy clothes”

Example

Example:

User prompt

Concept List:

$\langle bo \rangle$: $\langle bo \rangle$ is a cute golden retriever puppy with a playful expression.

$\langle shiba-sleep \rangle$: $\langle shiba-sleep \rangle$ is a shiba inu sleeping peacefully in a cozy home.

Expected output

```
{
  " $\langle bo \rangle$ ": {
    "category": "a golden retriever puppy",
    "attributes": "always playful expression"
  },
  " $\langle shiba-sleep \rangle$ ": {
    "category": "a shiba inu",
    "attributes": "can sleep peacefully; lives in cozy home"
  }
}
```

Rules

- Output plain English text only—no Markdown, no code fences.
- Keep every concept ID enclosed in angle brackets ($\langle \rangle$).
- Provide **exactly** the two keys for each concept—no extras.
- Do not add comments or explanatory text outside the JSON object.

User prompt.

Concept List

$\{C_i^u : \mathcal{T}_{C_i^u}\}_{i=0}^N$

Table 7: **Identity extraction prompt** for \mathcal{M}_l .

System prompt.

You are a visual verifier.

Task

You should answer each visual question with **yes** or **no**—nothing else.

User prompt.

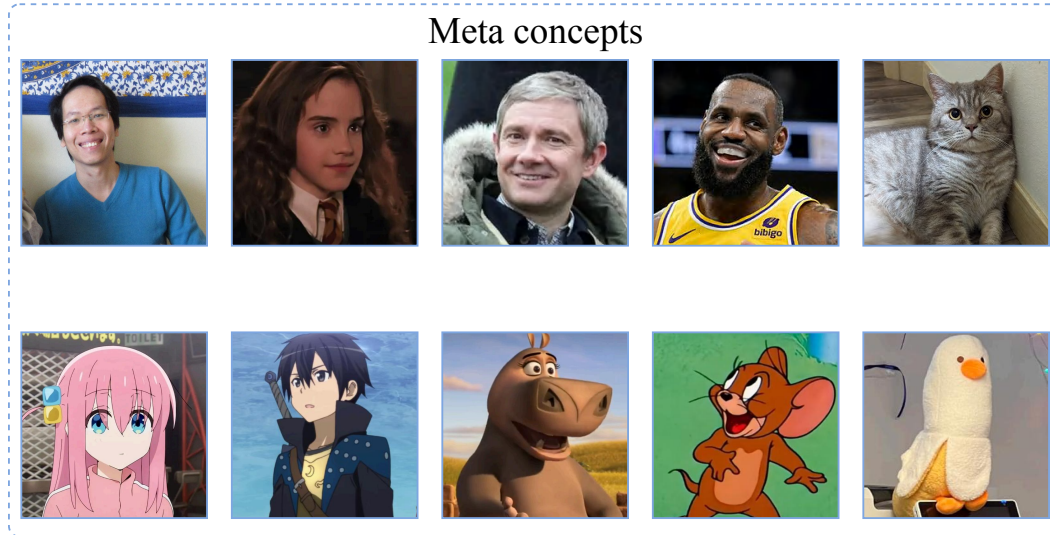
Q1. Do you see $\{ID(C_i^u)\}$ at $\{\text{location-absolute}\}$ of the image? (yes or no)

Q2. Is $\{ID(C_i^u)\}$ $\{\text{location-relative}\}$? (yes or no)

Rules

- Provide exactly one “yes” or “no” per question.
- If there are N questions, output N tokens separated by a single space.
- Do not include any additional words, punctuation, or commentary.

Table 8: **Test-time reflection prompt** for the large VLM \mathcal{M}_l for each concept detected by the small VLM.



828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851

Figure 5: Visualizations of the 10 meta-concepts.

832 B.3 ANSWER GENERATION PROMPT TEMPLATES

833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851

The large VLM \mathcal{M}_l produces the final response by conditioning on the structured detection report returned from the previous stages. The concise prompt is shown in Table 9.

<p>837 838 839 840 841 842 843 844 845 846 847 848 849 850 851</p> <p>System prompt. Detection Report $\{C_i^w: \text{present, category, attributes, location-absolute, location-relative}\}_{i=1}^N$</p> <p>Rules Use the Detection Report to answer the user’s visual question. <ul style="list-style-type: none"> • category: immutable essence (e.g. “a golden retriever puppy”). • attributes: mutable traits that may or may not be visible (e.g. “always playful expression”). • If <code>present = false</code>, it means the concept is not in the image. You should not mention the concept; reply “no” if asked about its presence. • If <code>present = true</code>, it means it concept is in the image. You should ground your answer strictly on the provided fields; reply “yes” if asked about its presence. </p> <p>User prompt. {User prompt}</p>

852
853
854
855
856
857
858

Table 9: Answer-generation prompt for \mathcal{M}_l .

855 C LORA TRAINING DETAILS

856 857 858 C.1 META-CONCEPT CLUSTERING & VISUALIZATION

859
860
861
862
863

CLIP Embedding Clustering Pipeline To extract meta-concepts from the Yo’LLaVA (Nguyen et al., 2024) and MC-LLaVA (An et al., 2024) datasets, we employ a pipeline based on feature clustering. First, we generate a feature embedding for each image using a pre-trained CLIP (Radford et al., 2021) model. We then create the vector representation for each concept by computing the average of its corresponding image embeddings. These concept-level vectors are subsequently clustered using the K-Means algorithm with the number of clusters set to 10. This process yields 10

864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

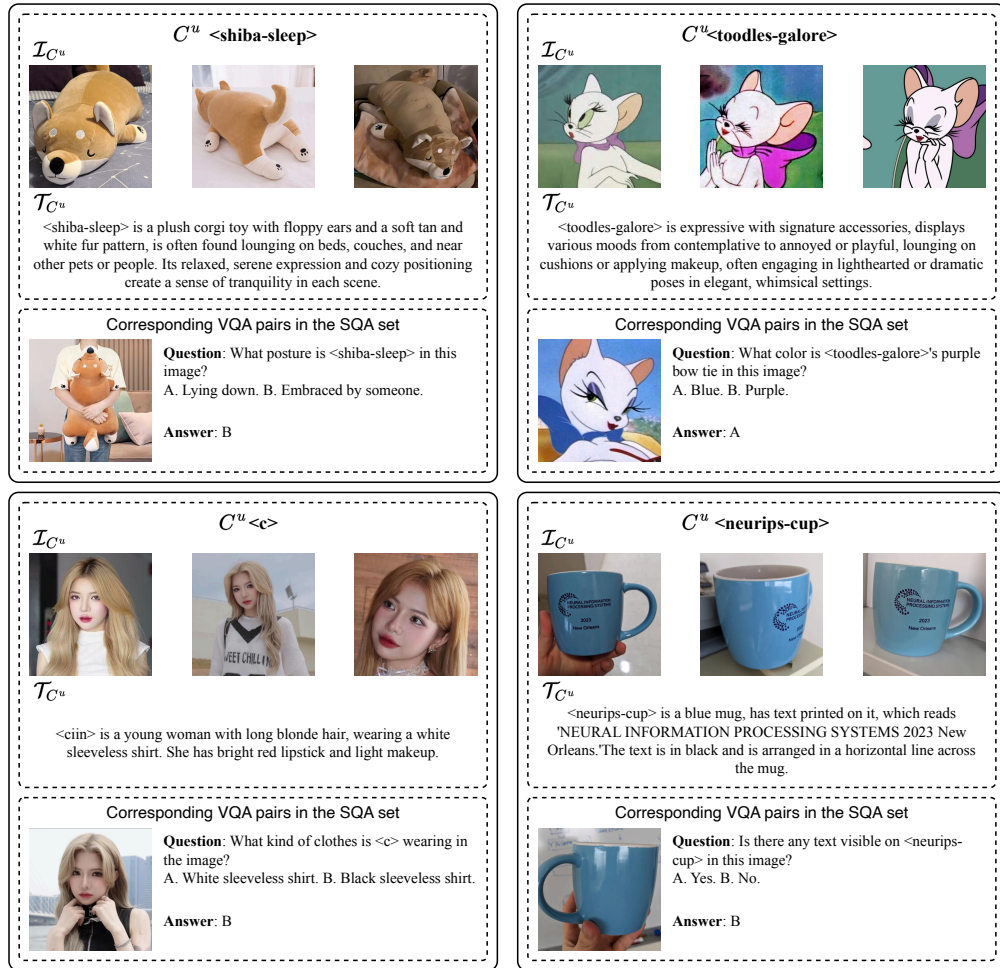


Figure 6: Examples of the Special Question-Answer (SQA) set.

meta-concepts, each defined as the concrete concept within a cluster whose embedding is closest to the cluster’s centroid.

Visualizations of the 10 Meta-Concepts Figure 5 provides a visualization of the 10 meta-concepts derived from our pipeline, exhibiting diversity and clear separation. The visualization showcases a wide array of categories, ranging from humans and animals to objects and cartoon characters.

C.2 TRAINING CONFIGURATION

We fine-tune the Qwen2.5-VL-3B-Instruct (Bai et al., 2025) model using Low-Rank Adaptation (LoRA) (Hu et al., 2022). LoRA modules of rank $r = 8$ and scaling factor $\alpha = 8$ are injected into every self-attention projection and FFN linear layer. Each meta-concept adapter is trained for 80 steps using the AdamW optimizer with a learning rate of 5×10^{-5} . The training is distributed across 8 A800 GPUs with a total batch size of 64 (8 samples per GPU).

D DATASET AND EVALUATION DETAILS

Special Question–Answer (SQA) Set Examples To diagnose overfitting, we build a Special Question–Answer (SQA) set composed entirely of VQA pairs. For every concept C_u in the Yo’LLaVA test split, we create at least two questions whose correct answers contradict visual or

918 textual cues found in the concept’s own training data ($\mathcal{I}_{C_u}, \mathcal{T}_{C_u}$). A model must therefore ground
919 its reasoning on the test image itself—rather than memorised artefacts—to succeed. Representative
920 examples are shown in Figure 6.



921
922 **Text-Only QA Protocol** Although SLC is designed for personalized VQA, it also accommodates
923 text-only questions. Following RAP-MLLM (Hao et al., 2025), we concatenate the user query
924 with the concept description \mathcal{T}_{C_u} . Before concatenation, however, \mathcal{M}_l performs identity extrac-
925 tion (Table 7) to convert each description into a structured category/attributes field. The resulting
926 prompt—{structured category/attributes field + user question}—is then fed back to \mathcal{M}_l for an-
927 swering. Table 1 in the main text shows that this extra step yields a slight accuracy gain on the
928 text-only test set.

929 E ADDITIONAL CASE STUDIES

930 We present three case studies. All examples use MetaC-3B as the small VLM (\mathcal{M}_s) and LLaVA-
931 1.5-13B (Liu et al., 2023) as the large VLM (\mathcal{M}_l). Figure 7 showcases a straightforward success
932 case, whereas Figure 8 and Figure 9 highlight how test-time reflection corrects a false positive from
933 the small VLM.
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

Personalized concepts:


<by>:  <by> is a young male with blonde hair and a serious expression. <hy>:  <hy> is a young female anime character with black hair and red eyes.

Query image:




User prompt:

Can you see <by> in this photo? Answer the question using a single word Yes or No.


 **Test-time detection:**

```
{
  "<by>": {
    "present": true,
    "location-absolute": "bottom-right quadrant",
    "location-relative": "to the right of <hy>",
  },
  "<hy>": {
    "present": true,
    "location-absolute": "top-left quadrant",
    "location-relative": "to the left of <hy>",
  }
}
```

 **Test-time reflection:**


Identity extraction:

```
{
  "<by>": {
    "category": "a young male",
    "attributes": "with blonde hair; serious expression"
  },
  "<hy>": {
    "category": "a young female anime character",
    "attributes": "with black hair; red eyes"
  }
}
```

 **Test-time reflection:**


Self-VQA verification for <hy>:

q1. yes\nq2. yes

 **Test-time reflection:**

Self-VQA verification for <by>:

q1. yes\nq2. yes

 **Answer generation:**


Yes.

Figure 7: **End-to-end inference example of the SLC pipeline.** The small VLM correctly identifies all concepts, the large VLM finds no inconsistencies, and SLC answers the user’s question accurately.

1026
 1027
 1028
 1029
 1030
 1031
 1032
 1033
 1034
 1035
 1036
 1037
 1038
 1039
 1040
 1041
 1042
 1043
 1044
 1045
 1046
 1047
 1048
 1049
 1050
 1051
 1052
 1053
 1054
 1055
 1056
 1057
 1058
 1059
 1060
 1061
 1062
 1063
 1064
 1065
 1066
 1067
 1068
 1069
 1070
 1071
 1072
 1073
 1074
 1075
 1076
 1077
 1078
 1079


Personalized concepts:

<tr>:




<tr> is a young male with black hair, wearing a black jacket with a white fur-lined collar and a white emblem on the left side of his chest.

<ys>:



<ys> is a young female anime character with long orange hair.

Query image:



User prompt: What is <ys> and <tr> doing in this image?

Test-time detection:

```
{
  "<tr>": {
    "present": false,
    "location-absolute": "",
    "location-relative": "",
  },
  "<ys>": {
    "present": true,
    "location-absolute": "bottom-right quadrant",
    "location-relative": "to the right of a male",
  }
}
```

Test-time reflection:
Identity extraction:

```
{
  "<tr>": {
    "category": "a young male",
    "attributes": "has black hair; wearing a black jacket with a white fur-lined collar and a white emblem on the left side of his chest"
  },
  "<ys>": {
    "category": "a young female anime character",
    "attributes": "with long orange hair"
  }
}
```

Test-time reflection:
Self-VQA verification for <ys>:


q1. no\nq2. no


Answer generation:

<ys> and <tr> are not present in this image.


Figure 8: **End-to-end inference example of the SLC pipeline.** The small VLM mistakenly detects <ys>, but the error is rectified by the large VLM via test-time reflection, enabling SLC to provide the correct answer.

1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133


Personalized concepts:  **<butin>** is a cream-colored Siberian Husky with striking blue eyes, relaxes comfortably in various indoor scenes, from couches to wooden floors. Known for its playful and laid-back demeanor, this Husky lounges in different positions, sometimes upside-down, showcasing its fluffy coat and cream-colored paws with an amusing, regal charm.

Query image: 


User prompt: Can you check if <butin> is in this photo? Answer the question using a single word Yes or No.

 **Test-time detection:**


```
{
  "<butin>": {
    "present": true,
    "location-absolute": "top-left quadrant",
    "location-relative": "on top of the person's lap",
  }
}
```

 **Test-time reflection:**
Identity extraction:

```
{
  "<butin>": {
    "category": "a cream-colored Siberian Husky with striking blue eyes",
    "attributes": "playful and laid-back demeanor; fluffy coat; cream-colored paws; regal charm"
  }
}
```

 **Test-time reflection:**
Self-VQA verification for <butin>:

q1. no\nq2. no

 **Answer generation:**

No.

Figure 9: **End-to-end inference example of the SLC pipeline.** The small VLM mistakenly detects <butin>, but the error is rectified by the large VLM via test-time reflection, enabling SLC to provide the correct answer.