PERSONALIZED LARGE VISION-LANGUAGE MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

The personalization model has gained significant attention in image generation yet remains underexplored for large vision-language models (LVLMs). Beyond generic ones, with personalization, LVLMs handle interactive dialogues using referential concepts (e.g., "Mike and Susan are talking.") instead of the generic form (e.g., "a boy and a girl are talking."), making the conversation more customizable and referentially friendly. In addition, PLVM is equipped to continuously add new concepts during a dialogue without incurring additional costs, which significantly enhances the practicality. PLVM proposes Aligner, a pre-trained visual encoder to align referential concepts with the queried images. During the dialogues, it extracts features of reference images with these corresponding concepts and recognizes them in the queried image, enabling personalization. We note that the computational cost and parameter count of the Aligner are negligible within the entire framework. With comprehensive qualitative and quantitative analyses, we reveal the effectiveness and superiority of PLVM. Code is attached in the supplementary materials.





Figure 1: With personalized concepts, PLVM enhances user interaction with large vision-language models, making it easier and more intuitive.

1 INTRODUCTION

The personalization of AI models provides customized services to users (Gal et al.), (Ruiz et al., 2023), (Nguyen et al., 2024)), enabling the understanding of user-specific concepts (Shi et al., 2024), (Ruiz et al., 2023), (Gal et al., 2023), (Han et al., 2023), (Kumari et al., 2023), (Zhou et al., 2024), (Zhao et al., 2023), (Jiang et al., 2024). This has been extensively explored in image/video generation fields (Gal et al.), (Ruiz et al., 2023), (Zhao et al., 2023), (Jiang et al., 2024), facilitating the understanding of these concepts for personalization. For example, Dreambooth (Ruiz et al., 2023) simplifies image generation by learning a personalized concept from a few images by tuning the network parameters using a trigger word *sks* in the prompt. During the test, the image from the same identity can be generated using the fine-tuned network with the prompt with the trigger word sks. MotionDirector (Zhao et al., 2023) learns motion concepts from videos to conveniently ren-052 der new videos with the same motions applied to different subjects. With such techniques, models perceive personalized concepts simply yet efficiently, empowering practicality.

000

001 002 003

004

005 006 007

008 009

010

011

Despite the notable advantages and successes, there remains a lack of research focused on personalization within the context of large vision-language model understanding. We argue that building personalized LVLMs for vision-language multimodal understanding can be helpful in various scenarios
like dialogue systems, QA (question-answering), human-computer interaction, *etc.*. See examples
in Figure 1 bottom. with personalization, a user asks a model simply by "Where is (billie)?".
Otherwise, the user has to describe the details of the queried subject like "Where is the person
with blonde hair and blue jacket?" Other examples also showcase the exciting applications and
practicality of personalized LVLMs.

062 There are two prior approaches for personalized large vision-language models - MyVLM (Alaluf 063 et al., 2024) and YoLLaVA (Nguyen et al., 2024). MyVLM requires an additional step to verify 064 if a personalized concept appears in an image, detecting all faces and using a classifier to identify the concept. This requires another fine-tuning process for other subjects beyond person and com-065 plicates the framework, as we aim for a unified model for both recognition and question-answering. 066 YoLLaVA simplifies this by embedding the personalized concepts and framing recognition as a sim-067 ple question: "Is (sks) in this photo?". However, YoLLaVA needs test-time fine-tuning and multiple 068 images (5-10) to learn the personalized embedding. We argue that a single facial image can rep-069 resent a person's identity without requiring multiple images, enhancing convenience. Furthermore, 070 the fine-tuning process takes around 40 minutes per identity, making it inefficient for real-time ap-071 plications. 072

This paper moves beyond the above methods and proposes personalized large vision-language 073 models (PLVM), a simple yet efficient LVLM to understand personalized concepts like 074 YoLLaVA (Nguyen et al., 2024). See Figure 1 for showcased abilities of PLVM. Advantageously, 075 PLVM does not require test-time fine-tuning, significantly strengthening the practicality as incor-076 porating new concepts will not incur additional costs. An essential design employs a pre-trained 077 vision network (referred to as Aligner) to cast the given personalized concepts into online features and align them with the queried images. When prompting a concept, Aligner casts only a 079 reference image of this concept to features, which serves as an identifier to recognize the subject of 080 the queried image without requiring any other fine-tuning processes. We are introducing Aligner 081 only conditions negligible additional cost, which will be showcased in the experiment.

The contributions of this paper are two-fold. First, we present PLVM, a personalized large vision-language model with solid abilities to recognize personalized concepts and freely incorporate new concepts without requiring additional costs. Second, PLVM achieves superior performance in a relatively simple way compared to existing methods. We also provide comprehensive experiments to showcase the effectiveness of PLVM.

087 088 089

090

2 RELATED WORK

Large Vision-language Models (LVLMs). Large language models (LLMs) (Brown, 2020; Chung 091 et al., 2024; Thoppilan et al., 2022; Chowdhery et al., 2023; Zhang et al., 2022a; Touvron et al., 092 2023; Zeng et al., 2022; Chiang et al., 2023) have opened a new era of AI, showcasing their po-093 tential to handle a wide array of language-based understanding and generation tasks. To extend 094 the capabilities of LLMs to visual understanding, the computer vision community has focused on 095 aligning language and vision data within a shared feature space (Radford et al., 2021). Research 096 in this area primarily follows two approaches: internal adaptation methods (Alayrac et al., 2022), 097 which integrate cross-attention within LLMs for visual-language alignment, and external adapta-098 tion methods (Li et al., 2023; Liu et al., 2024), which involve training modules for this purpose. Consequently, vision foundation models, particularly vision transformers (Dosovitskiy et al., 2020; 099 Liu et al., 2021; Radford et al., 2021; Tian et al., 2023; 2024b; Zhang et al., 2022b; Kirillov et al., 100 2023), have evolved into LVLMs (Liu et al., 2024; Tian et al., 2024a; Zhang et al., 2023; Lai et al., 101 2024; Qiu et al., 2024), endowing them with the capacity for language-guided visual understanding. 102 Despite these successes, research on their personalization is scarce. 103

Personalization in Image & Video Generation. Personalization in image and video generation aims to incorporate individualized concepts into pre-trained text-to-image or text-to-video diffusion models, generating specific personalized concepts across diverse contexts. Methods for image and video personalization generally fall into two categories: test-time fine-tuning approaches (Ruiz et al., 2023; Gal et al.; Kumari et al., 2023; Han et al., 2023; Voynov et al., 2023; Liu et al., 2023), which

involve fine-tuning word embeddings (Gal et al.) or network parameters (Ruiz et al., 2023) using a limited number of examples of the personalized concept; and encoder-based methods (Shi et al., 2024; Zhou et al., 2024; Gal et al., 2023; Ye et al., 2023; Wang et al., 2024; Ruiz et al., 2024), which eliminate the need for test-time fine-tuning by encoding the personalized concepts via a pre-trained vision encoder (Radford et al., 2021; Oquab et al.). The encoded features are then integrated into diffusion model components, such as word embeddings (Shi et al., 2024; Gal et al., 2023) or network parameters (Ruiz et al., 2024; Ye et al., 2023), to facilitate the generation of personalized content. The personalization of these generative models has dramatically improved their usability. This paper aims to apply this success to LVLMs.

Personalization in Large Vision-Language Models. In large vision-language models, when query-ing about a specific subject in an image, it is typically necessary to describe the visual appearance of that subject through prompts. Personalizing LVLMs involves enabling the model to recognize spe-cific identities, such as $\langle avril \rangle$ and $\langle billie \rangle$ in Figure 1, without requiring manual prompts for each identity. Previous methods, such as MyVLM (Alaluf et al., 2024) and YoLLaVA (Nguyen et al., 2024), employ techniques similar to textual inversion (Gal et al.) to learn the word embeddings of personalized concepts. However, these approaches require test-time fine-tuning, which is com-putationally intensive and costly for learning new concepts, leading to inefficiencies. Our method addresses this limitation by introducing an encoder-based approach that significantly reduces the need for fine-tuning, thereby improving efficiency.

3 PERSONALIZED LARGE VISION-LANGUAGE MODELS



Figure 2: The overall framework of PLVM, where the large language model receives the spatial features, the concept template of a personalized reference image, text prompt, and produces the answer.

Personalized large vision-language models enhance practicality and efficiency in applications such as dialogue systems. The proposed PLVM facilitates this capability using the straightforward yet effective Aligner method. This section outlines the Aligner and then describes the training strategy and synthesis dataset designed to improve PLVM's efficiency.

Figure 2 illustrates the PLVM framework, which includes the Aligner for generating personalized identifiers, a large language model (based on LLaVA (Liu et al., 2024)) for visual QA, and the integration of text and image inputs.

162 3.1 ALIGNER: ONLINE ENCODING OF PERSONALIZED CONCEPTS 163

164 As previously mentioned, we use a pre-trained vision encoder to build the Aligner, which embeds 165 personalized concepts online, with its output serving as an identifier for recognizing queried images.

166 Architecture. Specifically, we use the DINO-v2 (Oquab et al.) vision encoder, denoted as E_{ref} , to 167 extract features from the reference image of the personalized concept, Figure 2. Given a reference 168 image $I \in \mathbb{R}^{H \times W \times 3}$, where H and W are the height and width, we input the image into E_{ref} to 169 obtain visual features, represented as $z_{ref} \in \mathbb{R}^{L \times d}$, where L = 257, is the sequence length and d is 170 the feature dimension. 171

For the personalized concept (sks), we predict its word embedding (e^{word}) and head weight 172 (e^{weight}) using z_{ref} . We apply 2 MLP modules, f_{ϕ_1} and f_{ϕ_2} , to map the first (global) token of 173 z_{ref} , generating $e^{word} = f_{\phi_1}(z_{ref}[0])$ and $e^{weight} = f_{\phi_2}(z_{ref}[0])$. 174

175 To reduce the computational cost of processing visual tokens in large language models, we introduce $k \ (k \ll 257)$ context embeddings to integrate the outputs of E_{ref} (excluding the global token). 176 These are processed by a small transformer network \mathcal{T}_{θ} with k learnable queries $Q = q_1, q_2, \ldots, q_k$. 177 Each query acts as a soft, learnable token. Using cross-attention, z_{ref} serves as K and V, while the 178 learnable queries serve as Q, reducing 256 tokens to 16, aligning with YoLLaVA (Nguyen et al., 179 2024). 180

181 As shown in Figure 2, we freeze the LLaVA and Aligner modules and fine-tune f_{ϕ_1}, f_{ϕ_2} , and \mathcal{T}_{θ} . We also render e^{word} , e^{weight} , and the context embeddings learnable. Then, given the question X_q , 182 the answer $X_a = \{x_a^1, x_a^2, \dots, x_a^N\}$, and the query image I_q , we use the mask language modeling 183 184 loss to compute the probability of target response X_a : 185

$$\mathcal{L}(X_a|X_q, I_q) = \frac{1}{N} \sum_{i=1}^N \mathcal{L}_{ce}(p_i, x_a^i), \tag{1}$$

190 here, $\mathcal{L}_{ce}(\cdot, \cdot)$ represents the cross-entropy loss between the predicted answer and the ground truth, $p_i = p(x_a^i | x_a^{< i}, X_q, I_q)$ is the distribution of predicting the word *i*, and *N* denotes the sequence 192 length of the answer. 193

System prompt. Following YoLLaVA (Nguyen et al., 2024), let (sks) represent the personalized 194 concept and $\langle token_1 \rangle$, $\langle token_2 \rangle$, ..., $\langle token_k \rangle$ represent the context embedding tokens for a reference 195 image. We inject the template prompt as " $\langle sks \rangle$ is $\langle token_1 \rangle \langle token_2 \rangle \dots \langle token_k \rangle$ " as the instruction 196 prefix, resulting in the overall system prompt: 197

User: $\langle sks \rangle$ is $\langle token_1 \rangle \langle token_2 \rangle \dots \langle token_k \rangle + instruction Assistant:$

For each personalized concept, we assign a unique $\langle sks \rangle$, such as " $\langle avril \rangle$ " or " $\langle billie \rangle$ " as shown in Figure 1. The Aligner and \mathcal{T}_{θ} transform the corresponding reference image into k context embeddings. Unlike previous methods, this process works online, allowing our approach to freely integrate new concepts without additional costs like a fine-tuning process.

3.2 TRAINING

191

199 200

201

202

203

204 205

206

207 We observed that training solely with the objective in Equation 1 yields suboptimal results for recog-208 nition capabilities. This issue arises because, during training, the number of available question-209 answer pairs for task recognition is limited, whereas paired reference and query images are abundant. 210 Consequently, the model tends to overfit the words following a "Yes" or "No" answer. For instance, 211 if the question is "Is $\langle sks \rangle$ in this photo?", and the answer is "Yes, $\langle sks \rangle$ is in this photo.", once the 212 model recognizes the answer is "Yes," it can easily predict the phrase " $\langle sks \rangle$ is in this photo". As a 213 result, the most challenging recognition aspect is accurately predicting the words "Yes" or "No". To address this, we propose assigning greater weight to the loss associated with the words indicating 214 the positive or negative response ("Yes" or "No") in the recognition QA task. Specifically, the loss 215 function for recognition QA is:

220

221

222

224 225

226 227

228

229

230

231

232

233

234

235

236

237 238

239

247

248 249

250

251

252

253

254

255

256

257

$$\mathcal{L} = \frac{1}{W} \sum_{i=1}^{N} (w \cdot \mathbb{1}_i + (1 - \mathbb{1}_i)) \mathcal{L}_{ce}(p_i, x_a^i),$$
(2)

where w represents the weight associated with the loss for words indicating positive or negative answers, such as "Yes" and "No". The indicator function $\mathbb{1}_i$ is used to identify whether the word icorresponds to a positive or negative response. The normalization factor, W is defined as $1/(\sum_i w \cdot \mathbb{1}_i + (1 - \mathbb{1}_i))$, which adjusts the weight distribution to ensure proper scaling.

3.3 DATA SYNTHESIS

A significant challenge we encountered was the inability to obtain the required dataset. To overcome this issue and facilitate our method, we construct a dataset comprising paired images, specifically reference facial images and various corresponding images. For this purpose, we utilized the IP-Adapter method (Ye et al., 2023), which is pre-trained on the Stable Diffusion model (Rombach et al., 2022). Next, ChatGPT is employed to generate 200 descriptions based on the templates, such as "A person (wearing something) is (doing something) at (some place)." These prompts are then input into the IP-Adapter, conditioned on facial images from the CelebA-HQ dataset (Lee et al., 2020), to generate diverse images of the same person from the CelebA-HQ dataset. To ensure quality, we filter out low-quality images using two criteria: a CLIP cosine similarity score between the image and its corresponding prompt of less than 0.2, and a face similarity score (Deng et al., 2019) of less than 0.5. The examples of the synthetic dataset are shown in Figure 3.



Figure 3: Examples of our synthetic data using Diffusion models (Rombach et al., 2022), which is capable of generating customized images based on a reference appearance and a given prompt.

Following the methodology outlined in (Nguyen et al., 2024), we construct a dataset for two categories of questions, *i.e.* recognition and attribute questions. For recognition questions, given a reference image depicting the face of a personalized identity, the task is to determine whether the model can identify the presence of the individual in the query image. As in (Nguyen et al., 2024), we employ a set of template questions designed for binary "Yes" or "No" answers. For attribute questions, we inquire about details such as one person's hair color, eye color, and skin tone. Consistent with (Nguyen et al., 2024), we utilize the pre-trained LLaVA model (Liu et al., 2024) to generate responses to these questions based solely on the reference image. We provide the template details in Sec. A.2.

258 259

4 EXPERIMENT

4.1 Setups

- 260 261
- 262

Implementation details. An MLP module consists of 3 linear layers with a hidden size of 4096 and employs the GeLU activation function to predict both e^{word} and e^{weight} . To process the context embeddings, we adopt a cross-attention transformer architecture comprising four transformer blocks with a hidden size of 1024. During training, we allocate a portion of the samples for recognition tasks and another for detail attribute tasks, setting the default ratio for recognition to attribute at 1:1. For detail attribute-related tasks, query images are excluded from the input. The ratio of positive samples to the total is p, which will be analyzed in the ablation study. Unless otherwise specified, pis set to 0.5 by default.

		Question7:	What are <max> and <lucas> doing in this photo?</lucas></max>
	PVLM (ours) and YoLLaVA: <lucas></lucas>	YoLLaVA:	In the photo, <lucas> and <lucas> are posing for a picture together. They are both smiling and appear to be enjoying the moment</lucas></lucas>
New concepet	LaVA: <lucas> has a short haircut, a smile on Sface, and is wearing a black hoadie. GPT-W: <lucas> is young man with a short, curly hairstyle. He has a warm smile, bright eyes, and earnings.</lucas></lucas>	LLaVA+prompt:	In the photo, <max> and <lucas> are posing together for a picture, with the young man wearing a black hoodie and the young woman wearing a great shirt. They are both smilling and appear to be enjoying the moment.</lucas></max>
Text-only question-	answering		In the photo, <max> and <lucas> appear to be taking a selfie together. They are both smiling, indicating a friendly or</lucas></max>
<i>Question1:</i> YoLLaVA:	What is <lucas> skin color? <lucas> has a dark skin.</lucas></lucas>	GPT-4V+prompt:	playful atmosphere. The setting looks like it has a vintage or atmospheric background, contributing to the overall vibe of the moment they are capturing.
PLVM (ours):	<lucas>'s skin color is black.</lucas>	PLVM (ours):	<max> is standing next to <lucas>, posing for a picture.</lucas></max>
Question2:	What is <lucas> hairstyle?</lucas>		
YoLLaVA: PLVM (ours):	<lucas> has short, black hairstyle. <lucas> has short hairstyle.</lucas></lucas>	3	PVLM and YoLLaVA: <dustin> LLaVA: <dustin>'s face is characterized by a prominent noise, large eyes and a prominent chin The men here begred and supersolate</dustin></dustin>
Queried image		New concepts	GPT-4V: dustin> has curly hair, youthful, expressive eyes, neutral expression, slightly open mouth.
			LLaVA: <nancy> is the woman has a pretty face with brown hair, blue eyes, and a pouty mouth.</nancy>
Visual auestion-any	wering		GPT-4V: <nancy> has a youthful face with soft features. They have curly brown hair that</nancy>
Ouestion 3.	What is <lucas> facial expression in this photo?</lucas>		frames their face and a neutral expression, suggesting contemplation or concern.
YoLLaVA:	In the photo, <lucas> facial expression is a frown. (or a</lucas>		
II aVA+prompt:	pout) on the gul's face.		Ola Jal
CPT-4V+prompt:	The young man in the image has a neutral expression,	Queried image	
PLVM (ours):	appearing contemplative or serious as he sits at the table.	the table.	
		Visual question an	
(2)	$PVIM$ (ours) and VoI I aVA: $\leq may>$	visual question-an	
V	LLaVA: <max>'s face is a mix of red and white</max>	Question8: YoLLaVA:	What is <dustin> wearing in this photo? In the photo, <dustin> is wearing a green shirt.</dustin></dustin>
New concepet	of blue eyes.	LLaVA+prompt:	In the photo, <dustin> is wearing a green shirt.</dustin>
	defined eyebrows, soft lips, suble freekles, neutral expression, smooth complexion.	GPT-4V+prompt:	In the photo, the person is wearing a light blue shirt underneath a denim jacket. They also have curly hair and a neutral expression.
Text-only question-	answering	PLVM (ours):	dustin> is wearing a green shirt in this photo.
Question4:	What is <max> hair color?</max>	Ownerfantle	
YoLLaVA:	<max>'s hair color is red.</max>	Vol LaVA:	what is the shirt color of the person next to <nancy>? The person pays to <nancy> is wearing a green shirt</nancy></nancy>
rvLM (ours):	<max>'s hair color is red.</max>	LLaVA+prompt:	The person next to <nancy> is wearing a green shirt. The person next to <nancy> is wearing a yellow shirt.</nancy></nancy>
Question5:	What is the gender of <max>?</max>	GPT-4V+prompt:	The shirt color of the person next to <nancy> is white.</nancy>
YoLLaVA: PVLM (ours):	<max>'s gender is a female. The gender of <max> is a female.</max></max>	PLVM (ours):	The person next to <nancy> is wearing a yellow shirt.</nancy>
		Question 10.	Where is <nancy> in this photo?</nancy>
		YoLLaVA:	No, <nancy> is not in this photo.</nancy>
Queried		LLaVA+prompt:	In the photo, the woman is standing next to the group of young people, who are all dressed the same outfit. She is positioned between two of the young men, one on her left
image			and the other on her right.
		GP1-4v + prompt PLVM (ours):	I m sorry, I cannot identify the person in the photo. <nancy> is at the far left of the photo.</nancy>
Visual question-ans	wering		
Question6:	Can you check if <max> is in this photo?</max>	Question11:	Where is <lucas> in this photo?</lucas>
YoLLaVA:	No, <max> is not in this photo.</max>	YoLLaVA:	<lucas> is at the right corner of the photo. In the photo, <lucas> is standing payt to a group of people.</lucas></lucas>
LLaVA+prompt:	Yes, the girl in the photo is wearing a red headband and has a face with a mix of red and white colors, with a prominent nose and a pair of blue eyes.	LLaVA+prompt:	including a young man with green shirt and a backpack. The group is posing for a picture, and they are all standing in front of a car.
GPT-4V+prompt:	I can't determine if <max> is in the photo.</max>	GPT-4V+prompt:	I don't know who is <lucas> in the photo.</lucas>
PLVM (ours):	Yes. I can confirm that <max> is in this photo.</max>	PVLM (ours):	<lucas> is at the right corner of the photo.</lucas>

Figure 4: Qualitative results compared with YoLLaVA, LLaVA (with prompt), and GPT-4V (with prompt). PLVM requires no fine-tuning for every concept, enabling seamless incorporation of new concepts, as shown in the figure. In contrast, YoLLaVA needs ~40 minutes of fine-tuning for each new concept to achieve personalization.

311 312

308

309

310

313

All parameters are frozen except those in the MLP and cross-attention modules. The AdamW optimizer is used with a learning rate of 1×10^{-4} and a batch size of 2. All experiments are conducted on a single A6000 GPU.

Dataset generation and evaluation. As described in Section 3.1, we generate 200 prompt templates and use each image from the CelebA-HQ dataset as a reference image. Five new images are generated using randomly sampled prompts for each reference image. Low-quality images are filtered based on two criteria: a CLIP cosine similarity score below 0.2 and an ArgFace cosine similarity score below 0.5. This process results in a final dataset of 70k images.

For validation, we curate a test dataset by manually downloading images from the Internet. The dataset is organized by identity, each corresponding to a distinct concept. Each identity contains one

Table 1: The evaluation results on three types of questions, compared to other methods on the test dataset. This demonstrates that our PLVM achieves superior performance across all these tasks and showcases the effectiveness of our method in comparison to these counterparts.

Method	Visu	al recogniti	Question answering		
	Positive	Negative	Mean	Text-only	Visual
LLaVA-1.5 (Liu et al., 2024) (w/o prompt)	0.0	100	50.0	-	_
LLaVA-1.5 (Liu et al., 2024) (w/ prompt)	80.3	61.3	70.8	71.8	72.3
YoLLaVA (Nguyen et al., 2024)	73.1	73.7	73.4	77.4	72.9
GPT-4V (Achiam et al., 2023)	40.5	97.5	69.0	47.9	50.5
PLVM (Ours)	78.5	84.7	81.6	85.9	77.5



Figure 5: Examples of the three types of evaluation questions: recognition, text-only QA, and visual QA questions.

reference image and 5-10 query images, resulting in 34 identities and 246 images. We conduct the evaluation using three types of questions, as detailed below:

- **Recognition questions** involves determining whether a given personalized concept, denoted as (sks), is present in a query image. We ask each reference-query image pair: 'Is (sks) in this photo?' If the query image shares the same identity as the reference image, the correct answer is 'Yes'; otherwise, for images with differing identities, the answer is 'No.' Following the evaluation protocol from (Nguyen et al., 2024), we report the positive, negative, and mean accuracy percentages.
- **Text-only QA questions** provide only the reference image while asking the model about the attributes of the corresponding query image. Following the approach of (Nguyen et al., 2024), we construct multiple-choice questions for this task, with 71 questions.
- Visual QA questions involve presenting the query image corresponding to each reference image and constructing multiple-choice questions based on visual details, such as the subject's position. Three hundred fifty questions are generated, and we report the accuracy based on the correct responses.
- 365366 Examples of these three types of questions are illustrated in Figure 5.

4.2 PLVM IS A STRONG BASELINE FOR PERSONALIZED LVLMS

Quantitative results, and necessity of online encoding for personalization. We compare our method with YoLLaVA (Nguyen et al., 2024), LLaVA (Liu et al., 2024) (using prompt), and GPT-4V (Achiam et al., 2023). To ensure a fair comparison, we utilize 16 context embedding tokens, consistent with the number used in YoLLaVA. For LLaVA, we prompt the pre-trained LLaVA model with a text description of the personalized concepts, which also contains approximately 16 tokens. This same prompting strategy is applied to GPT-4V.

The results are presented in Table 1. Regarding visual recognition questions, LLaVA with prompts performs best in positive question-answering. However, it produces numerous incorrect predictions on negative question-answering. GPT-4V, on the other hand, tends to predict 'No' for many images, which results in its highest accuracy for negative question-answering but lost in positive ones.
 Overall, our proposed PLVM outperforms the other methods regarding mean accuracy.

PLVM achieves the best results among all methods for text-only question-answering. This task involves questions about attributes of the personalized concepts such as hair color, skin tone, hairstyle, and age, highlighting our model's more robust ability to encode personalized characteristics. In the visual question-answering question, PLVM also achieves higher accuracy than its counterparts, demonstrating the superiority of our PLVM.

Qualitative results. Figure 2 presents the qualitative results, demonstrating the process of se quentially incorporating new concepts three times, ultimately yielding four personalized concepts:
 (lucas), (max), (dustin), and (nancy). These qualitative results are compared with three baselines:
 YoLLaVA, LLaVA (with prompt), and GPT-4V (with prompt). For both LLaVA and GPT-4V, the
 prompts employed are: "Describe this person's face."

For text-only questions, which address the general attributes of a person, PLVM produces answers comparable to those of YoLLaVA. YoLLaVA is trained explicitly on text-only questions, highlighting that our Aligner can effectively encode general personal information without necessitating test-time fine-tuning. In comparison to the general prompts used by LLaVA and GPT-4V, PLVM demonstrates superior capability in capturing specific attributes, such as skin color for $\langle lucas \rangle$ and gender for $\langle max \rangle$, without the need for manual prompting by a human.

For visual QA, we evaluate all models' ability to recognize the person, determine their location, 397 identify their clothing, and assess their facial expressions in the image. Overall, PLVM provides 398 more accurate responses than other methods. GPT-4V (with prompt) (Achiam et al., 2023) often 399 responds with "No, I can't determine $\langle sks \rangle$ in this photo" when asked about visual recognition 400 or location. YoLLaVA occasionally misidentifies individuals in the image (e.g., confusing $\langle lucas \rangle$ 401 with another person in the first queried image, or (max) in the second queried image), leading to 402 incorrect answers, such as errors in identifying facial expressions or details in the images. LLaVA 403 (with prompt) (Liu et al., 2024) performs well on the positive recognition benchmark as shown in 404 Table1. However, qualitative results reveal that LLaVA (with prompt) tends to generate extraneous 405 and often inaccurate information. Additionally, because the prompts are derived from the reference 406 image, there is a risk of bias toward the reference image's clothing and facial expressions, resulting 407 in incorrect answers for the queried image.

408

More examples are provided in Sec. A.1, where we also discuss the failure cases and the limitations.

409 410 411

4.3 ABLATION STUDY

The weights on 'Yes'/'No' answers. We first conduct an ablation study to assess the impact of different weights (w) in Equation 2 for recognition questions during training. The results are presented in Table 2. Without applying a weighting scheme (*i.e.*, a 1:1 ratio), the model's performance is suboptimal, achieving only 4% accuracy for positive responses. However, as the weight increased, performance improved, with optimal results observed when w ranged around 10-20. This experiment demonstrates the model's sensitivity to w for practical training, a behavior that differs notably from YoLLaVA.

Number of context embedding tokens. Table 3 presents the model's performance using different numbers of context embedding tokens, which correspond to the tokens representing personalized concepts. The model's performance improves as the number of learnable queries increases (from 8 to 12 and 16). However, beyond a certain threshold—precisely 16 tokens—the performance plateaus indicate that a moderate number of tokens can effectively represent a personalized reference image. This finding is consistent with the results reported for YoLLaVA.

Different vision encoders for Aligner. In addition to the default use of the Dino-v2 vision encoder (Oquab et al.), we extend our experiments to include other models, specifically CLIP (Radford et al., 2021), and ViT (Dosovitskiy et al., 2021). All model weights are sourced from official opensource repositories. The results of these experiments are summarized in Table 4. For DINO-v2, we employ the base model with a sequence length of 257, while for CLIP-large and ViT-base, the sequence lengths are set to 257 and 197, respectively. As the SAM model lacks a global token, we use k + 2 context embedding tokens to represent the word, context, and weight embeddings, setting the number of context embedding tokens to 16. As shown in Table 4, Dino-v2 outperforms

434								
435	Weights on 'Yes'/'No'	Positive	Negative	Mean	token num.	Positive	Negative	Mean
436	1	4.0	96.0	50.0	8	78 5	73.1	75.8
437	5	76.8	80.2	78.5	12	70.5	78.8	70.2
438	10	79.6	82.4	81.0	12	79.0	70.0 947	19.2 01.6
439	15	79.6	79.8	79.7	10	/8.5	84./	81.0
440	20	78.5	84.7	81.6	20	78.0	80.6	79.3

Table 2: Ablative study on the weight (w in Equa- Table 3: Ablation study on the number of con-432 tion 2) of 'Yes'/'No' answers. 433

434 435

441 442

443

444 445

446

454

455

456

457

458 459

460 461

462

text embedding tokens.

the other visual encoders. Despite this, models such as CLIP and ViT still demonstrate competitive performance, suggesting that our framework is compatible with various visual encoder architectures.

Table 4: Ablation study on different visual encoder Table 5: The positive/negative sampling ramodels for Aligner. tio during training.

Models	Positive	Negative	Mean	Pos./Neg. sampling	Positive	Negative	Mean
CLIP (Radford et al., 2021)	74.6	79.2	76.9	0.4	76.8	82.1	79.4
SAM (Kirillov et al., 2023)	66.3	49.0	57.7	0.5	78.5	84.7	81.6
ViT (Dosovitskiy et al., 2021)	75.1	83.5	79.3	0.6	89.0	78.9	84.0
Dino-v2 (Oquab et al.)	78.5	84.7	81.6	0.7	91.5	70.6	81.0

The ratio of the positive questions to the total, p. We conduct an ablation study to investigate the impact of the ratio of positive questions to the total, denoted as p, during the training phase. The results are summarized in Table 5. Within the range of p from 0.4 to 0.6, increasing p consistently leads to an improvement in mean accuracy. However, when p reaches 0.7, we observe a decline in performance.

4.4 The advantages and costs compared to previous methods

Table 6: A design comparison with YoLLaVA and MyVLM. † represents this method is pre-trained on large-scale realistic datasets and does not support text-only prompting.

Method	Ext. module	FT time	Positive	Negative	Mean	Visual question
YoLLaVA (Nguyen et al., 2024)	X	$\sim 40 \text{ mins}$	73.1	73.7	73.4	72.9
MyVLM (Alaluf et al., 2024) [†]	\checkmark	$\sim 1 \min$	76.8	_	_	48.5
PLVM (Ours)	×	0	78.5	84. 7	81.6	77.5

470 The advantages in terms of model design and runtime. We present the running times for YoLLaVA (Nguyen et al., 2024) and MyVLM (Alaluf et al., 2024) in Table 6 for processing a 471 single new concept. YoLLaVA requires approximately 40 minutes to fine-tune a new concept, while 472 MyVLM takes about 1 minute, which still falls short for real-time applications. However, MyVLM 473 necessitates an additional module for recognition. In contrast, our method achieves superior accu-474 racy compared to YoLLaVA and MyVLM without requiring any test-time fine-tuning or additional 475 modules. This demonstrates that the encoding scheme PLVM employs for personalized concepts is 476 more effective than those used by YoLLaVA and MyVLM. 477

The cost of the proposed Aligner. Table 7 compares the cost of the Aligner module relative 478 to the overall framework. We examine both the running time and the number of parameters. The 479 results show that the newly introduced Aligner accounts for only 3.2% and 1.8% of the total 480 running time and parameter count of LLaVA, respectively, indicating that the cost of the Aligner 481 module is negligible. 482

Results on YoLLaVA's dataset. We also conduct experiments using the YoLLaVA dataset. We 483 only chose the person identity for the experiment, cropped the face from one image in the training 484 set as the reference image for each identity, and used the whole test set of that identity for the 485 queried image. The recognition results are presented in Table 8. Our method significantly improves

486 accuracy over YoLLaVA, highlighting that the encoder effectively captures robust information even 487 with a single reference image. 488

489 relative to the overall framework. 490

491

495

509

510

529

Table 7: A comparison of the cost of Aligner Table 8: A comparison with YoLLaVA using its dataset.

model	Running time	#Params	Method	Positive	Negative	e Mean
LLaVA (Liu et al.,	2024) 0.19s	7063M	YoLLaVA (Nguyen et al., 2024)	62.5	84.6	73.6
Aligner (Propos	(ed) $0.006s (3.2\%)$	128M (1.8%)	PLVM (Ours)	89.4	83.0	86.2

What do the context tokens learn? Context tokens play a pivotal role in the Aligner module 496 as the compact representation of the reference image. To elucidate the specific meaning encoded 497 by each token, we prompt the model with the instruction, "describe $\langle \text{token}_i \rangle$ in detail." Using the 498 reference image from Figure 2 as an example, we extract the learned meanings for each token. 499 The results reveal that while many tokens capture similar aspects of the image, subtle differences in 500 detail are also evident, Table 9. As shown, we list the meanings associated with all 16 context tokens 501 corresponding to the reference image. These tokens capture fine-grained attributes of the subject, 502 such as hair color, whether the subject is wearing glasses, shirt color, and facial expressions. The global token comprehensively describes the reference image: " $\langle ada \rangle$ is wearing a black shirt and 504 has short, dark hair. He is looking directly at the camera with a serious expression. $\langle ada \rangle$ appears to be in his late teens or early twenties and is standing in front of a wall. No additional information 505 about $\langle ada \rangle$'s background or identity is provided in the image." This demonstrates that the proposed 506 Aligner effectively captures information from the reference image without fine-tuning. Instead, 507 it only leverages text-only QA training, similar to YoLLaVA (Nguyen et al., 2024). 508

Table 9: The learned meaning of the 16 context tokens of the reference image in Figure 2. Each token has learned slightly different features of the reference image.

$\langle token_1 angle$:a man with dark hair	$\langle token_9 \rangle$:He appears to be focusing on
	something in the distance
$\langle token_2 \rangle$: He appears to be looking down	$\langle token_{10} \rangle$:has dark eyes
$\langle token_3 \rangle$:appears to be in a casual and	$\langle token_{11} \rangle$:in a serious and thoughtful mood
related setting	
$\langle token_4 \rangle$:appears to be young and has a	$\langle token_{12} \rangle$:concentrating on something in
confident demeanor	the distance
$\langle token_5 \rangle$:he is wearing glasses	$\langle token_{13} \rangle$: is the main focus of the image,
	and no other people or objects
$\langle token_6 \rangle$:appears to be a businessman or	$\langle token_{14} \rangle$:appears to be the main subject
possibly a lawer or a politician	of the image,a formal setting
$\langle token_7 \rangle$:looking at his reflection in the	$\langle token_{15} \rangle$:The room appears to be dimly lit,
mirror, admiring his outfit	creating a mysterious atmosphere
$\langle token_8 \rangle$:wearing a black shirt	$\langle token_{16} \rangle$:the man is positioned in the
	center of the image

CONCLUSION 5

530 The study introduces a robust baseline called the Personalized Large Vision-Language Model 531 (PLVM), which is both an intriguing and practical approach for enhancing the understanding of 532 personalized concepts during dialogues. Unlike existing methods, PLVM uniquely supports con-533 tinuously adding new concepts throughout a dialogue without incurring additional costs, thereby 534 significantly improving its practicality. Specifically, PLVM incorporates Aligner, a pre-trained 535 visual encoder designed to align referential concepts with the queried images. During dialogues, 536 Aligner extracts features from reference images associated with these concepts and identifies 537 them in the queried image, thus enabling effective personalization. Notably, the computational cost and parameter count of Aligner are minimal within the overall framework. We hope that PLVM 538 will establish a solid baseline for advancing personalization in the domain of large vision-language models.

540 REFERENCES 541

547

553

569

571

572

573

581

582

583

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-542 man, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical 543 report. *arXiv preprint arXiv:2303.08774*, 2023. 544
- Yuval Alaluf, Elad Richardson, Sergey Tulyakov, Kfir Aberman, and Daniel Cohen-Or. Myvlm: 546 Personalizing vlms for user-specific queries. arXiv preprint arXiv:2403.14599, 2024.
- Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel 548 Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language 549 model for few-shot learning. Advances in neural information processing systems, 35:23716– 550 23736, 2022. 551
- 552 Tom B Brown. Language models are few-shot learners. arXiv preprint arXiv:2005.14165, 2020.
- Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, 554 Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. Vicuna: An 555 open-source chatbot impressing gpt-4 with 90%* chatgpt quality, March 2023. URL https: 556 //lmsys.org/blog/2023-03-30-vicuna/.
- 558 Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam 559 Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: Scaling language modeling with pathways. Journal of Machine Learning Research, 24(240): 560 1-113, 2023.561
- 562 Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, 563 Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. Scaling instruction-finetuned lan-564 guage models. Journal of Machine Learning Research, 25(70):1-53, 2024. 565
- Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos Zafeiriou. Arcface: Additive angular margin 566 loss for deep face recognition. In Proceedings of the IEEE/CVF conference on computer vision 567 and pattern recognition, pp. 4690-4699, 2019. 568
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas 570 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929, 2020.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas 574 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszko-575 reit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at 576 scale. ICLR, 2021. 577
- 578 Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit Haim Bermano, Gal Chechik, and 579 Daniel Cohen-or. An image is worth one word: Personalizing text-to-image generation using textual inversion. In The Eleventh International Conference on Learning Representations. 580
 - Rinon Gal, Moab Arar, Yuval Atzmon, Amit H Bermano, Gal Chechik, and Daniel Cohen-Or. Encoder-based domain tuning for fast personalization of text-to-image models. ACM Transactions on Graphics (TOG), 42(4):1–13, 2023.
- 585 Ligong Han, Yinxiao Li, Han Zhang, Peyman Milanfar, Dimitris Metaxas, and Feng Yang. Svdiff: Compact parameter space for diffusion fine-tuning. In Proceedings of the IEEE/CVF Interna-586 tional Conference on Computer Vision, pp. 7323–7334, 2023.
- 588 Yuming Jiang, Tianxing Wu, Shuai Yang, Chenyang Si, Dahua Lin, Yu Qiao, Chen Change Loy, and 589 Ziwei Liu. Videobooth: Diffusion-based video generation with image prompts. In Proceedings of 590 the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 6689–6700, 2024.
- Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete 592 Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 4015–4026, 2023.

594 Nupur Kumari, Bingliang Zhang, Richard Zhang, Eli Shechtman, and Jun-Yan Zhu. Multi-concept 595 customization of text-to-image diffusion. In Proceedings of the IEEE/CVF Conference on Com-596 puter Vision and Pattern Recognition, pp. 1931–1941, 2023. 597 Xin Lai, Zhuotao Tian, Yukang Chen, Yanwei Li, Yuhui Yuan, Shu Liu, and Jiaya Jia. Lisa: Rea-598 soning segmentation via large language model. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 9579–9589, 2024. 600 601 Cheng-Han Lee, Ziwei Liu, Lingyun Wu, and Ping Luo. Maskgan: Towards diverse and interactive 602 facial image manipulation. In IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2020. 603 604 Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image 605 pre-training with frozen image encoders and large language models. In International conference 606 on machine learning, pp. 19730–19742. PMLR, 2023. 607 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. Advances 608 in neural information processing systems, 36, 2024. 609 610 Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. 611 Swin transformer: Hierarchical vision transformer using shifted windows. In Proceedings of the 612 IEEE/CVF international conference on computer vision, pp. 10012–10022, 2021. 613 Zhiheng Liu, Ruili Feng, Kai Zhu, Yifei Zhang, Kecheng Zheng, Yu Liu, Deli Zhao, Jingren Zhou, 614 and Yang Cao. Cones: concept neurons in diffusion models for customized generation. In Pro-615 ceedings of the 40th International Conference on Machine Learning, pp. 21548–21566, 2023. 616 617 Thao Nguyen, Haotian Liu, Yuheng Li, Mu Cai, Utkarsh Ojha, and Yong Jae Lee. Yo'llava: Your personalized language and vision assistant. arXiv preprint arXiv:2406.09400, 2024. 618 619 Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy V Vo, Marc Szafraniec, Vasil Khalidov, 620 Pierre Fernandez, Daniel HAZIZA, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learn-621 ing robust visual features without supervision. Transactions on Machine Learning Research. 622 Jihao Qiu, Yuan Zhang, Xi Tang, Lingxi Xie, Tianren Ma, Pengyu Yan, David Doermann, Qixiang 623 Ye, and Yunjie Tian. Artemis: Towards referential understanding in complex videos. arXiv 624 preprint arXiv:2406.00258, 2024. 625 626 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, 627 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual 628 models from natural language supervision. In International conference on machine learning, pp. 8748-8763. PMLR, 2021. 629 630 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-631 resolution image synthesis with latent diffusion models. In Proceedings of the IEEE/CVF confer-632 ence on computer vision and pattern recognition, pp. 10684–10695, 2022. 633 Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. 634 Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In Pro-635 ceedings of the IEEE/CVF conference on computer vision and pattern recognition, pp. 22500– 636 22510, 2023. 637 638 Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Wei Wei, Tingbo Hou, Yael Pritch, Neal Wadhwa, 639 Michael Rubinstein, and Kfir Aberman. Hyperdreambooth: Hypernetworks for fast personaliza-640 tion of text-to-image models. In Proceedings of the IEEE/CVF Conference on Computer Vision 641 and Pattern Recognition, pp. 6527–6536, 2024. 642 Jing Shi, Wei Xiong, Zhe Lin, and Hyun Joon Jung. Instantbooth: Personalized text-to-image gen-643 eration without test-time finetuning. In Proceedings of the IEEE/CVF Conference on Computer 644 Vision and Pattern Recognition, pp. 8543-8552, 2024. 645 Romal Thoppilan, Daniel De Freitas, Jamie Hall, Noam Shazeer, Apoorv Kulshreshtha, Heng-Tze 646 Cheng, Alicia Jin, Taylor Bos, Leslie Baker, Yu Du, et al. Lamda: Language models for dialog 647 applications. arXiv preprint arXiv:2201.08239, 2022.

- 648 Yunjie Tian, Lingxi Xie, Zhaozhi Wang, Longhui Wei, Xiaopeng Zhang, Jianbin Jiao, Yaowei Wang, 649 Qi Tian, and Qixiang Ye. Integrally pre-trained transformer pyramid networks. In Proceedings 650 of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 18610– 651 18620, June 2023. 652 Yunjie Tian, Tianren Ma, Lingxi Xie, Jihao Qiu, Xi Tang, Yuan Zhang, Jianbin Jiao, Qi Tian, 653 and Qixiang Ye. Chatterbox: Multi-round multimodal referring and grounding. arXiv preprint 654 arXiv:2401.13307, 2024a. 655 656 Yunjie Tian, Lingxi Xie, Jihao Qiu, Jianbin Jiao, Yaowei Wang, Qi Tian, and Qixiang Ye. Fast-itpn: 657 Integrally pre-trained transformer pyramid network with token migration. IEEE Transactions on 658 Pattern Analysis and Machine Intelligence, 2024b. 659 660 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée 661 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and 662 efficient foundation language models. arXiv preprint arXiv:2302.13971, 2023. 663 Andrey Voynov, Qinghao Chu, Daniel Cohen-Or, and Kfir Aberman. p+: Extended textual condi-664 tioning in text-to-image generation. arXiv preprint arXiv:2303.09522, 2023. 665 666 Qixun Wang, Xu Bai, Haofan Wang, Zekui Qin, and Anthony Chen. Instantid: Zero-shot identity-667 preserving generation in seconds. arXiv preprint arXiv:2401.07519, 2024. 668 669 Hu Ye, Jun Zhang, Sibo Liu, Xiao Han, and Wei Yang. Ip-adapter: Text compatible image prompt 670 adapter for text-to-image diffusion models. arXiv preprint arXiv:2308.06721, 2023. 671 672 Aohan Zeng, Xiao Liu, Zhengxiao Du, Zihan Wang, Hanyu Lai, Ming Ding, Zhuoyi Yang, Yifan Xu, Wendi Zheng, Xiao Xia, et al. Glm-130b: An open bilingual pre-trained model. arXiv preprint 673 arXiv:2210.02414, 2022. 674 675 Shilong Zhang, Peize Sun, Shoufa Chen, Min Xiao, Wenqi Shao, Wenwei Zhang, Yu Liu, Kai Chen, 676 and Ping Luo. Gpt4roi: Instruction tuning large language model on region-of-interest. arXiv 677 preprint arXiv:2307.03601, 2023. 678 679 Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christo-680 pher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. Opt: Open pre-trained transformer 681 language models. arXiv preprint arXiv:2205.01068, 2022a. 682 Xiaosong Zhang, Yunjie Tian, Wei Huang, Qixiang Ye, Qi Dai, Lingxi Xie, and Qi Tian. Hivit: 683 Hierarchical vision transformer meets masked image modeling. arXiv preprint arXiv:2205.14949, 684 2022b. 685 686 Rui Zhao, Yuchao Gu, Jay Zhangjie Wu, David Junhao Zhang, Jiawei Liu, Weijia Wu, Jussi Keppo, 687 and Mike Zheng Shou. Motiondirector: Motion customization of text-to-video diffusion models. 688 arXiv preprint arXiv:2310.08465, 2023. 689 690 Yufan Zhou, Ruiyi Zhang, Jiuxiang Gu, and Tong Sun. Customization assistant for text-to-image 691 generation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recog-692 nition, pp. 9182–9191, 2024. 693 694 APPENDIX А 696 697 A.1 ADDITIONAL QUALITATIVE RESULTS 698 We provide additional qualitative results of our PLVM in Figure 6 and Figure 7. With personaliza-699
- tion, PLVM enables the conversations to be more customizable and referentially friendly. Notably,
 after training, we do not require the fine-tuning process, which implies that the examples we showed are continuous inferences without other extra processing.



A.3 LIMITATION

We show some failure cases using our method as shown in Figure 8. If multiple people in the image
share standard features such as hair color or makeup style (like (avril) case), PLVM can make the
wrong prediction of the personalized concept. In addition, accessories such as glasses can make the
wrong recognition between the personalized concept and the query image ((DN) case).



Figure 7: Additional qualitative results of PLVM

812			
813	ID	Question	Answer
814	1	Is $\langle sks \rangle$ in this photo?	Yes, $\langle sks \rangle$ is in this photo.
815	2	Can you tell if $\langle sks \rangle$ appears in this picture?	Yes, $\langle sks \rangle$ appears in this picture.
816	3	Could you check whether $\langle sks \rangle$ is in the image?	Yes, $\langle sks \rangle$ is indeed in the image.
817	4	Do you see $\langle sks \rangle$ anywhere in this snapshot?	Yes, $\langle sks \rangle$ is visible in this snapshot.
818	5	Is there a chance $\langle sks \rangle$ could be in this photo?	Yes, $\langle sks \rangle$ is in this photo.
819	6	Would you happen to know if $\langle sks \rangle$ is shown in	Yes, $\langle sks \rangle$ is shown in this photograph.
820		this photograph?	
821	7	Can you see $\langle sks \rangle$ in this photo?	Yes, I can confirm that $\langle sks \rangle$ is indeed in
822			the photo.
823	8	Have you spotted $\langle sks \rangle$ in this photo?	Yes, I have spotted $\langle sks \rangle$ in this photo.
824	9	Is that $\langle sks \rangle$ in the photo there?	Yes, that is $\langle sks \rangle$ in the photo.
825	10	Is $\langle sks \rangle$ in this image?	Yes, $\langle sks \rangle$ is in this image.
025	11	Am I seeing $\langle sks \rangle$ in this picture?	Yes, you are seeing $\langle sks \rangle$ in this picture.
020	12	Does this photo include $\langle sks \rangle$?	Yes, this photo includes $\langle sks \rangle$.
027	13	Is $\langle sks \rangle$ featured in this photo?	Yes, $\langle sks \rangle$ is featured in this photo.
828	14	Can you point out $\langle sks \rangle$ in this photo?	Yes, $\langle sks \rangle$ is right here in the photo.
829	15	Is $\langle sks \rangle$ visible in this photo?	Yes, $\langle sks \rangle$ is visible in this photo.
830	16	Check if $\langle sks \rangle$ is in this photo for me, will you?	Yes, $\langle sks \rangle$ is in this photo.
831	17	Hey AI, can you tell me if you see $\langle sks \rangle$ in this	Yes, I can confirm that $\langle sks \rangle$ is indeed
832		photo?	
833			in the photo.
834	18	Do you recognize $\langle sks \rangle$ in this photo?	Yes, I do recognize $\langle sks \rangle$ in this photo.
835	19	I'm looking for $\langle sks \rangle$, is $\langle sks \rangle$ in this photo?	Yes, $\langle sks \rangle$ is in this photo.
836	20	Can you see if $\langle sks \rangle$ is in this photo?	Yes, $\langle sks \rangle$ is indeed in the photo.
837	21	This photo, does it have $\langle sks \rangle$?	Yes, this photo does have $\langle sks \rangle$.
838	22	Could you confirm if this is $\langle sks \rangle$ in the photo?	Yes, I can confirm that this is $\langle sks \rangle$ in
830	22		the photo.
0.00	23	Any chance that $\langle sks \rangle$ might be in this photo?	Yes, $\langle sks \rangle$ is definitely in this photo.
040	24	Can you recognize if $\langle SKS \rangle$ is in this photo?	Yes, $\langle SKS \rangle$ is indeed in the photo.
841	25	Has (sks) been captured in this photo?	Yes, $\langle SKS \rangle$ has been captured in this
842	26	/alas/?a in this all stan sight?	photo. V_{abc} (also) is this shots
843	20	$\langle SKS \rangle$ S in this photo, right?	Yes, $\langle SKS \rangle$ is in this photo.
844	27	is (sks) present in this particular photo?	res, (sks) is present in this particular
845	28	Hay AL can you tell ma if you recognize /ska	Ves L can see $\langle sks \rangle$ in the photo
846	20	in this photo?	
847	20	In this photo? Can you see if $\langle sks \rangle$ is in this photo?	\mathbf{V}_{ab} (sks) is in this photo
848	29	Can you see II (sks) is in this photo?	105, (SKS) is in this photo.

Table 10: The used Question-answering templates for positive (Yes) recognition. The $\langle sks \rangle$ represents the personalized concept.

A.4 TEST DATASET

849 850 851

852

We collect authentic images from various identities to construct our test dataset, consisting of 34
identities and 246 images. Each identity includes one reference image representing the individual's
face and 5 to 10 test images of the same person. Following the methodology of YoLLaVA (Nguyen
et al., 2024), we create three types of QA tasks: recognition QA, text-only QA, and visual QA. We
compare each identity to each image in the test set for the recognition QA, resulting in 246 positive
and 8118 negative QA pairs. We use the prompt to prompt the model for recognition tasks: "Is (sks)
in this photo? Answer with Yes or No."

For the text-only QA, we provide a reference image alongside multiple-choice questions about the
person's attributes, such as hair color and hairstyle, resulting in 71 text-only questions. Similarly,
for the visual QA, we ask questions related to the queried image, such as the person's position,
clothing, and height comparison, given the reference and queried images. This results in a total of
350 visual QA questions. For both the text-only and visual QA tasks, we prompt the model with the

ID	Question	Answer
1	Is $\langle sks \rangle$ in this photo?	No, $\langle sks \rangle$ is not in this photo.
2	Can you tell if $\langle sks \rangle$ appears in this picture?	No, $\langle sks \rangle$ is not in this picture.
3	Could you check whether $\langle sks \rangle$ is in the image?	No. $\langle sks \rangle$ is not in the image.
4	Do you see $\langle sks \rangle$ anywhere in this snapshot?	No. $\langle sks \rangle$ does not appear in this snapsh
5	Is there a chance $\langle sks \rangle$ could be in this photo?	No, there's no indication that $\langle sks \rangle$ is in
C		this photo.
6	Would you happen to know if $\langle sks \rangle$ is shown in	No. $\langle sks \rangle$ is not shown in the photogram
	this photograph?	
7	Can you recognize if $\langle sks \rangle$ is in this photo?	No. $\langle sks \rangle$ is not present in the photo.
8	Have you spotted (sks) in this photo?	No. I have not spotted (sks) in this pho
9	Is that $\langle sks \rangle$ in the photo there?	No that is not $\langle sks \rangle$ in the photo
10	Can you check if $\langle sks \rangle$ is in this photo?	No. $\langle sks \rangle$ is not present in the image vo
10		provided.
11	Am I seeing $\langle sks \rangle$ in this picture?	No. you are not seeing $\langle sks \rangle$ in this pic
12	Does this photo include $\langle sks \rangle$?	No, this photo does not include $\langle sks \rangle$.
13	Is $\langle sks \rangle$ featured in this photo?	No. $\langle sks \rangle$ is not featured in this photo.
14	Can you point out $\langle sks \rangle$ in this photo?	No I can't point out $\langle sks \rangle$ because $\langle sks \rangle$
1.	Cuil you point out (sits) in this photo.	is not in the photo
15	Is $\langle sks \rangle$ visible in this photo?	No $\langle sks \rangle$ is not visible in this photo
16	$Check if \langle sks \rangle$ is in this photo for me will you?	No $\langle sks \rangle$ is not in this photo
17	Can you see $\langle sks \rangle$ in this photo?	No $\langle sks \rangle$ is not present in the photo
18	Do you recognize $\langle sks \rangle$ in this photo?	No I do not recognize $\langle sks \rangle$ in this photo-
19	I'm looking for $\langle sks \rangle$ is $\langle sks \rangle$ in this photo?	No $\langle sks \rangle$ is not in this photo
$\frac{1}{20}$	Is there any sign of $\langle sks \rangle$ in this photo?	No there is no sign of $\langle sks \rangle$ in this photo.
21	This photo does it have $\langle sks \rangle$?	No this photo does not have (sks)
$\frac{21}{22}$	Could you confirm if this is $\langle sks \rangle$ in the photo?	No $\langle sks \rangle$ is not in the photo
$\frac{22}{23}$	Can you see if $\langle sks \rangle$ is in this photo?	No $\langle sks \rangle$ is not present in the photo.
$\frac{23}{24}$	Is (sks) part of the group in this photo?	No, $\langle sks \rangle$ is not part of the group in the
27	is (sks/ part of the group in this photo:	nhoto
25	I think I see $\langle sks \rangle$ is it so?	No $\langle sks \rangle$ is not in the photo
$\frac{25}{26}$	Has (sks) been captured in this photo?	No, $\langle sks \rangle$ has not been captured in this
20	Thus (sks) occur captured in this photo:	nhoto
27	/sks/'s in this photo right?	No $\langle sks \rangle$'s not in this photo
$\frac{27}{28}$	Is (sks) present in this particular photo?	No, $\langle sks \rangle$ is not present in this particult
20	is (sks/ present in this particular photo:	nhoto
20	$\int \mathbf{L} \operatorname{can}' t \operatorname{find} \langle \operatorname{sks} \rangle$ is $\langle \operatorname{sks} \rangle$ in the photo?	No you can't find /sks\ because /sks\
29	real time (sks/, is (sks/ in the photo?	not in the photo
30	Is /sks/ in this image?	No $\langle sks \rangle$ is not in this image
30	15 (SKS) III UIIS IIIage :	1 NO, (SKS) IS NOT III UIIS IIIage.

ID	Question
1	What is $\langle sks \rangle$'s hair color?
2	What color are $\langle sks \rangle$'s eyes?
3	What is $\langle sks \rangle$'s skin tone?
4	How would you describe $\langle sks \rangle$'s hairstyle?
5	Does $\langle sks \rangle$ have any distinctive facial features?
6	Is $\langle sks \rangle$ young or old?
7	What do you describe about $\langle sks \rangle$'s nose?

instruction: "Choose the letter corresponding to the correct answer" followed by the question and answer options.

