# ConsistentID🥗: Portrait Generation with Multimodal Fine-Grained Identity Preserving

Jiehui Huang[1], Xiao Dong[2], Wenhui Song[1], Hanhui Li[1], Jun Zhou[1], Yuhao Cheng[3], Shutao Liao[1], Long Chen[3], Yiqiang Yan[3], Shengcai Liao[4], and Xiaodan Liang[1]*

[1]Shenzhen Campus of Sun Yat-sen University, [2]Zhuhai Campus of Sun Yat-sen University, [3]Lenovo Research, [4]Inception Institute of Artificial Intelligence

{xdliang328}@gmail.com
https://ssugarwh.github.io/consistentid.github.io/

**Fig. 1:** Given some images of input IDs, our ConsistentID can generate diverse personalized ID images based on text prompts using only a single image.

**Abstract.** Diffusion-based technologies have made significant strides, particularly in personalized and customized facial generation. However, existing methods face challenges in achieving high-fidelity and detailed identity (ID) consistency, primarily due to insufficient fine-grained control over facial areas and the lack of a comprehensive strategy for ID preservation by fully considering intricate facial details and the overall face. To address these limitations, we introduce **ConsistentID**, an innovative method crafted for diverse identity-preserving portrait generation under fine-grained multimodal facial prompts, utilizing only a single reference image. ConsistentID comprises two key components: a multimodal facial prompt generator that combines facial features, corresponding facial descriptions and the overall facial context to enhance precision in facial details, and an ID-preservation network optimized through the facial attention localization strategy, aimed at preserving ID consistency in facial regions. Together, these components significantly enhance the accuracy of ID preservation by introducing fine-grained multimodal ID information from facial regions. To facilitate training of ConsistentID, we present a fine-grained portrait dataset, FGID, with

---

over 500,000 facial images, offering greater diversity and comprehensiveness than existing public facial datasets. Experimental results substantiate that our ConsistentID achieves exceptional precision and diversity in personalized facial generation, surpassing existing methods in the MyStyle dataset. Furthermore, while ConsistentID introduces more multimodal ID information, it maintains a fast inference speed during generation. Our codes and pre-trained checkpoints will be available at https://github.com/JackAILab/ConsistentID.

**Keywords:** Portrait generation · fine-grained conditions · identity preservation

# 1   Introduction

Recently, image-generation technology  [14, 16, 25, 41, 46] has undergone significant evolution, driven by the emergence and advancement of diffusion-based [11, 45] text-to-image large models like GLIDE [30], DALL-E 2 [36], Imagen [42], Stable Diffusion (SD) [37], eDiff-I [1] and RAPHAEL [55]. This progress has given rise to a multitude of application approaches across diverse scenarios. Positioned as the central focus of these application approaches, personalized and customized portrait generation has attracted widespread attention in both academic and industrial domains, owing to its extensive applicability in downstream tasks such as E-commerce advertising, personalized gift customization and virtual try-ones.

The primary challenge in customized facial generation lies in maintaining facial image consistency across different attributes based on one or multiple reference images, leading to two key issues: ensuring accurate identity (ID) consistency and achieving high-fidelity, diverse facial details. Current text-to-image models [41, 49, 51, 56, 60], despite incorporating structural and content guidance, face limitations in accurately controlling personalized and customized generation, particularly concerning the fidelity of generated images to reference images.

To improve the precision and diversity of personalized portrait generation with reference images, numerous customized methodologies have emerged, meeting users' demands for high-quality customized images. These personalized approaches are categorized based on whether fine-tuning occurs during inference, resulting in two distinct types: test-time fine-tuning and direct inference. **Test-time fine-tuning**: This category includes methods such as Textual Inversion [8], HyperDreambooth [41], and CustomDiffusion [20]. Users can achieve personalized generation by providing a set of target ID images for post-training. Despite achieving commendable high-fidelity results, the quality of the generated output depends on the quality of manually collected data. Additionally, the manual collection of customized data for fine-tuning introduces a labor-intensive and time-consuming aspect, limiting its practicality. **Direct inference**: Another category of models, including IP-Adapter [56], Fastcomposer [54], Photomaker [22], and InstantID [53], adopts a single-stage inference approach. These models enhance global ID consistency by either utilizing the image as a conditional input or manipulating image-trigger words. However, most methods frequently overlook fine-grained information, such as landmarks and facial features. Although InstantID

improves ID consistency to some extent by introducing landmarks, the visual prompt landmark restricts the diversity and variability of key facial regions, leading to stiff generated facial features. In summary, **two pivotal challenges** requiring meticulous consideration persist in personalized portrait generation: 1) neglect of fine-grained facial information and 2) identity inconsistency between facial areas and the whole face, as illustrated in Figure 5.

To address these challenges, we introduce a novel method, ConsistentID, crafted to maintain identity consistency and capture diverse facial details through multimodal fine-grained ID information, employing only a single facial image while ensuring high fidelity. Figure 2 provides the overview of our ConsistentID. ConsistentID comprises two key modules: 1) a multimodal facial prompt generator and 2) an ID-preservation network. The former component includes a fine-grained multimodal feature extractor and a facial ID feature extractor, enabling the generation of more detailed facial ID features using multi-conditions, incorporating facial images, facial regions, and their corresponding textual descriptions extracted from the multimodal large language model LLaVA1.5 [23]. Utilizing the facial ID features obtained from the initial module, we feed them into the latter module, promoting ID consistency across each facial region via the facial attention localization strategy. Additionally, we recognize the limitations of existing portrait datasets [3, 27, 31, 52, 61], particularly in capturing diverse and fine-grained identity-preserving facial details, crucial for the effectiveness of ConsistentID. To address this, we introduce the inaugural Fine-Grained ID Preservation (FGID) dataset, along with a fine-grained identity consistency metric, providing a unique and comprehensive evaluation approach to enhance our training and performance evaluation in facial details.

In summary, our contributions are as follows.

- We introduce ConsistentID to improve fine-grained customized facial generation by incorporating detailed descriptions of facial regions and local facial features. Experimental results showcase the superiority of ConsistentID in terms of ID consistency and high fidelity, even with just one reference image. Simultaneously, despite the introduction of more detailed multimodal fine-grained ID information in ConsistentID, the inference speed remains relatively efficient, as shown in Table 1

- We devise an ID-preservation network optimized by facial attention localization strategy, enabling more accurate ID preservation and more vivid facial generation. This mechanism ensures the preservation of ID consistency within each facial region by preventing the blending of ID information from different facial regions.

- We introduce the inaugural fine-grained facial generation dataset, FGID, addressing limitations in existing datasets for capturing diverse identity-preserving facial details. This dataset includes facial features and descriptions of both facial regions and the entire face, complemented by a novel fine-grained identity consistency metric, establishing a comprehensive evaluation framework for fine-grained facial generation performance.

## 2   Related Work

**Text-to-image Diffusion Models.** Diffusion models have made notable advancements, garnering significant attention from both industry and academia, primarily due to their exceptional semantic precision and high fidelity. The success of these models can be attributed to the utilization of high-quality image-text datasets, continual refinement of foundation modality encoders, and the iterative enhancement of controlled modules. In the domain of text-to-image generation, the encoding of text prompts involves utilizing a pretrained language encoder, such as CLIP [35], to transform it into a latent representation, subsequently inserted into the diffusion model through the cross-attention mechanism. Pioneering models in this realm encompass GLIDE [30], SD [37], DiT [32], among others, with further developments and innovations continuing to emerge [13,55]. A notable advancement in this lineage is SDXL [33], which stands out as the most powerful text-to-image generation model. It incorporates a larger Unet [38] model and employs two text encoders for enhanced semantic control and refinement. As a follow-up, we use SD [37] model as our base model to achieve personalized portrait generation.

**Personalization in Diffusion Models.** Due to the potent generative capability of the text-to-image diffusion model, many personalized generation models are constructed based on it. The mainstream personalized image synthesis methods are categorized into two groups based on whether fine-tuning occurs during test time. One group relies on optimization during test-time, with typical methods including Dreambooth [40], Textual Inversion [8], IP-Adapter [53], ControlNet [60], Custom Diffusion [20], and LoRA [12]. Dreambooth and Textual Inversion fine-tune a special token S* to learn the concept during the fine-tuning stage. In contrast, IP-Adapter, ControlNet, and LoRA insert image semantics using an additional learned module, such as cross-attention, to imbue a pre-trained model with visual reasoning understanding ability.

Despite their advancements, these methods necessitate resource-intensive backpropagation during each iteration, making the learning process time-consuming and limiting their practicality. Recently, researchers have focused more on methods bypassing additional fine-tuning or inversion processes, mainly including IP-Adapter [56], FastComposer [54], PhotoMaker [22], and InstantID [53]. This type of method performs personalized generation using only an image with a single forward process, which is more advantageous in calculation efficiency compared to the former type. However, we observe that fine-grained facial features are not fully considered in the training process, easily leading to ID inconsistency or lower image quality, as shown in Figure 4. To address these limitations, we introduce ConsistentID, aiming to mitigate the ID-preserving issue and enhance fine-grained control capabilities while reducing data dependency. Our approach incorporates a specially designed facial encoder using detailed descriptions of facial features and local image conditions as inputs. Additionally, we contribute to a new landmark in the facial generation field by proposing: 1) the introduction of the first fine-grained facial generation datasets and 2) the presentation of a new metric that redefines the performance evaluation of facial generation.
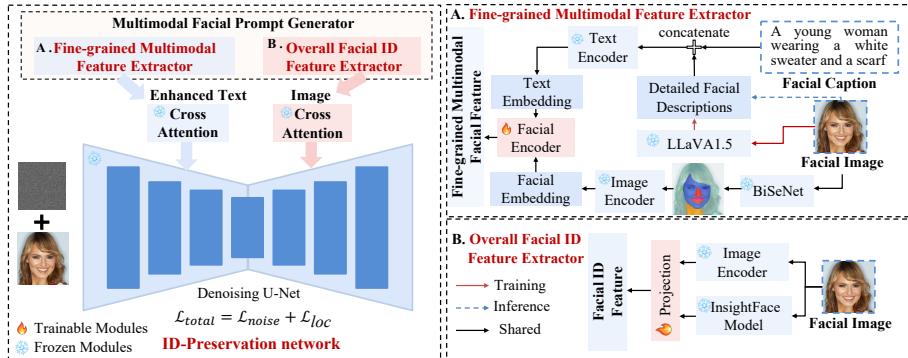
**Fig. 2:** The overall framework of our proposed ConsistentID. The framework comprises two key modules: a multimodal facial ID generator and a purposefully crafted ID-preservation network. The multimodal facial prompt generator consists of two essential components: a fine-grained multimodal feature extractor, which focuses on capturing detailed facial information, and a facial ID feature extractor dedicated to learning facial ID features. On the other hand, the ID-preservation network utilizes both facial textual and visual prompts, preventing the blending of ID information from different facial regions through the facial attention localization strategy. This approach ensures the preservation of ID consistency in the facial regions.

## 3   Method

### 3.1   Multimodal Facial Prompt Generator

**Fine-grained Multimodal Feature Extractor.** In this module, we independently learn fine-grained facial visual and textual embeddings and feed them into the designed lightweight facial encoder to generate fine-grained multi-modal facial features. Three key components are used in the module, including text embedding, facial embedding and facial encoder.

1) **Text Embedding.** Motivated by recent works [25, 26, 29, 50] in personalized facial generation, our goal is to introduce more detailed and accurate facial descriptions. To achieve this, we input the entire facial image into the Multimodal Large Language Model (MLLM) LLaVA1.5 [23] using the command prompt 'Describe this person's facial features, including face, ears, eyes, nose, and mouth' to obtain a description at the facial feature level. Subsequently, we replace the words 'face, ears, eyes, nose, and mouth' in these feature-level descriptions with the delimiter '<facial>' and concatenate them with the captions of the entire facial image. Finally, the concatenated descriptions are fed into the pre-trained text encoder to learn fine-grained multimodal facial features. With both visual and textual descriptions containing more precise ID information, our ConsistentID effectively mitigates ID inconsistency issues in facial details.

2) **Facial Embedding.** In contrast to existing methods [17, 28], [8, 9, 17, 22, 24, 28, 39, 53, 54]. that rely on brief textual descriptions or coarse-grained visual prompts, our goal is to integrate more fine-grained multimodal control

information at the facial region level. This aims to achieve improved accuracy in identity (ID) consistency and diverse facial generation. To enrich the ID-preservation information, we delve into more fine-grained facial features, including eye gaze, earlobe characteristics, nose shape, and others. Following the previous method [15, 18, 21, 48, 58, 59], we employ the pre-trained face model BiSeNet [57] to extract segmentation masks of facial areas, encompassing eyes, nose, ears, mouth, and other regions, from the entire face. Subsequently, the facial region images obtained from these masks are fed into the pre-trained image encoder to learn fine-grained facial embeddings. The inclusion of facial regions' features results in fine-grained facial embeddings containing more abundant ID-preservation information compared to features learned from the entire face. More detailed processes are outlined in the supplementary materials.

3) **Facial Encoder.** Previous studies [8, 35, 40, 54, 57] have demonstrated that relying solely on visual or textual prompts cannot comprehensively maintain ID consistency both in appearance and semantic details. While IP-Adapter makes the initial attempt to simultaneously inject multimodal information through two distinct decoupled cross-attention mechanisms, it overlooks ID information from crucial facial regions, rendering it susceptible to ID inconsistency in facial details.

To cultivate the potential of image and text prompts, inspired by the token fusion approach of multimodal large language models, we design a facial encoder to seamlessly integrate visual prompts with text prompts along the dimension of the text sequence, as depicted in Figure 3. Specifically, given a facial embedding and a caption embedding, the facial encoder initially employs a self-attention mechanism to align the entire facial features with facial areas' features, resulting in aligned features denoted as $\widehat{f^i} \in \mathbb{R}^{N \times D}$, where $N = 5$ represents the number of facial feature areas, including eyes, mouth, ears, nose, and other facial regions, and $D$ represents the dimension of text embeddings. In cases where face images lack a complete set of N facial features, the missing features are padded using an all-zero matrix. Subsequently, we replace the text features at the position of the delimiter '<facial>' with $\widehat{f^i}$, and then employ two multi-layer perceptron (MLP) to learn the text conditional embeddings.

**Facial ID Feature Extractor.** Except for the input condition of fine-grained facial features, we also inject the character's overall ID information into our ConsistentID as a visual prompt. This process relies on the pre-trained CLIP image encoder and the pre-trained face model from the specialized version of the IP-Adapter [56] model, IPA-FaceID-Plus [56]. Specifically, the complete facial images are concurrently fed into both encoders for visual feature extraction. After these two encoders, a lightweight projection module, with parameters initialized by IPA-FaceID-Plus, is used to generate the face embedding of the whole image.

## 3.2   ID-Preservation network

The effectiveness of image prompts in pre-trained text-to-image diffusion models [4, 22, 34, 47, 54, 60] significantly enhances textual prompts, especially for content that is challenging to describe textually. However, visual prompts alone often provide only coarse-grained control due to the semantic fuzziness of visual
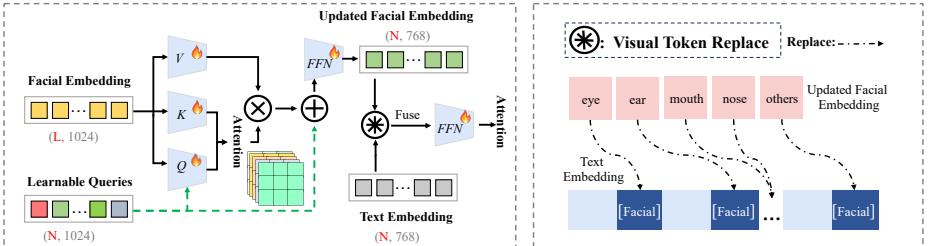
**Fig. 3:** The framework of our facial encoder for generating fine-grained multimodal facial features.

tokens. To solve this, we follow the IP-Adapter and integrate fine-grained multimodal ID prompts and overall ID prompts into the UNet model through the cross-attention module to achieve precise ID preservation.

Motivated by Fastcomposer [54], we introduce an ID-consistent network to maintain consistency in local ID features by directing the attention of facial features to align with the corresponding facial regions. This optimization strategy is derived from the observation that traditional cross-attention maps tend to simultaneously focus on the entire image, posing a challenge in maintaining ID features during the generation of facial regions. To address this issue, we introduce facial segmentation masks during training to obtain attention scores learned from the enhanced text cross-attention module for facial regions.

Let $P \in [0, 1]^{h \times w \times n}$ represent the cross-attention map that connects latent pixels to multimodal conditional embeddings at each layer, where $P[i, j, k]$ signifies the attention map from the $k$-th conditional token to the $(i, j)$ latent pixel. Ideally, the attention maps of facial region tokens should focus exclusively on facial feature areas, preventing the blending of identities between facial features and averting propagation to the entire facial image. To achieve this goal, we propose localizing the cross-attention map using the segmentation mask of reference facial regional features.

Let $M = \{m_1, m_2, m_3, ..., m_N\}$ represent the segmentation masks of the reference portrait, $I = \{i_1, i_2, i_3, ..., i_N\}$ as the index list indicating which facial feature corresponds to visual and textual tokens in the multimodal prompt, and $P_i = P[:, :, i] \in [0, 1]^{h \times w}$ denote the cross-attention map of the $i$-th facial region's token, where $I$ is generated using the special token '<facial>'. Given the cross-attention map $P_{i_j}$, it should closely correspond to the facial region identified by the $j$-th multimodal token, segmented by $m_j$. To achieve this, we introduce $m_j$ and apply it to $P_{i_j}$ to obtain its corresponding activation region, aligning with the segmentation mask $m_j$ of the $j$-th facial feature token. For achieving this correspondence, a balanced L1 loss is employed to minimize the distance between the cross-attention map and the segmentation mask:

$$\mathcal{L}_{\text{loc}} = \frac{1}{N} \sum_{j=1}^{N} \left( \text{mean} \left( P_{i_j} [1 - m_j] \right) - \text{mean} \left( P_{i_j} [m_j] \right) \right), \tag{1}$$

where N denotes the number of the segmentaion masks. This loss formulation aims to ensure that each facial feature token's attention map aligns closely with

its corresponding segmentation mask, promoting precise and localized attention during the generation process.

### 3.3   Training and Inference Details

During the training process, we optimize only the parameters of the facial encoder and the projection module within the overall facial ID feature extractor, while maintaining the parameters of the pre-trained diffusion model in a frozen state. The training data for ConsistentID consists of facial image-text pairs. In ConsistentID, to enhance text controllability, we prioritize the caption as the primary prompt and concatenate it with more detailed descriptions of facial regions extracted from LLaVA1.5, forming the ultimate textual input. Regarding training loss functions, they align with those used in the original stable diffusion models and are expressed as:

$$\mathcal{L}_{noise} = \mathbb{E}_{z_t,t,C_f,C_l,\epsilon\sim\mathcal{N}(0,1)} \left[ \|\epsilon - \epsilon_\theta \left(z_t, t, C_f, C_i\right)\|_2^2 \right], \tag{2}$$

where $C_l$ denotes the facial ID feature, and $C_f$ is the fine-grained multimodal facial feature.

The total loss function is $\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{noise}} + \mathcal{L}_{\text{loc}}$.

During the inference process, we employ a straightforward delayed primacy condition as similar to Fastcomposer. This allows the use of a separate text representation initially, followed by enhanced text representation after a specific step, effectively balancing identity preservation and editability. More discussions are provided in the supplementary materials.

### 3.4   Fine-grained Human Dataset Construction

Our ConsistentID necessitates detailed facial features and corresponding textual prompts to address issues like deformation, distortion, and blurring prevalent in current facial generation methods. However, existing datasets [2, 5, 17, 61] predominantly focus on local facial areas and lack fine-grained ID annotations for specific features such as the nose, mouth, eyes, and ears.

To tackle this issue, we introduce a dataset pipeline outlined in the supplementary materials to create our dataset FGID. This dataset encompasses comprehensive fine-grained ID information and detailed facial descriptions, which are crucial for training the ConsistentID model.

Our FGID dataset, comprising 525,258 images is a compilation from diverse public datasets such as FFHQ [17], CelebA [27], SFHQ [2], *etc*, including about 524,000 training data. This dataset is tailored to encompass both whole-face and facial feature information, providing richer textual and visual details for model training[1]. More details about the FGID dataset are available in the supplementary materials, providing information on its statistics, characteristics, and

---

[1] In the future, we intend to introduce additional datasets, including higher-level datasets like LAION-Face [61] and self-collected multi-ID data, to further augment diversity and information content.

related ethical considerations. In the following, we elaborate on the multi-modal data handling process. **Textual Data**: For textual facial descriptions, MLLM LLaVA1.5 is utilized to extract detailed information using an embedded prompt 'Please describe the people in the image, including their gender, age, clothing, facial expressions, and any other distinguishing features'. **Visual Data**: BiSeNet [57] and InsightFace [7] models are deployed to capture both whole-face ID information and facial feature information, ensuring comprehensive identity details.

## 4 Experiment

**Experimental Implements.** In ConsistentID, we employ the SD1.5 model as the foundational text-to-image model. For the fine-grained multimodal feature extractor, we initialize the parameters of all text encoders and image encoders with CLIP-ViT-H [44] and utilize its image projection layers to initialize the learnable projection in the overall facial ID feature extractor. The entire framework is optimized using Adam [19] on 8 NVIDIA 3090 GPUs, with a batch size of 16. We set the learning rate for all trainable modules to $1 \times 10^{-4}$. During training, we probabilistically remove 50% of the background information from the characters with a 50% probability to mitigate interference. Additionally, to enhance generation performance through classifier guidance, there is a 10% chance of replacing the original updated text embedding with a zero text embedding. In inference, we employ delayed topic conditioning [17, 54] to resolve conflicts between text and ID conditions. We utilize a 50-step DDIM [45] sampler, and the scale for classifier-guided settings is set to 5.

**Experimental Metrics.** To assess the effectiveness and efficiency of ConsistentID, we employ six widely used metrics [40]: CLIP-I [8], CLIP-T [35], DINO [6], FaceSim [43], FID [10], and inference speed. CLIP-T measures the average cosine similarity between prompt and image CLIP embeddings, evaluating ID fidelity. CLIP-I calculates the average pairwise cosine similarity between CLIP embeddings of generated and real images, assessing prompt fidelity. DINO represents the average cosine similarity between ViT-S/16 embeddings of generated and real images, indicating fine-grained image-level ID quality. FaceSim determines facial similarity between generated and real images using FaceNet [43]. FID gauges the quality of the generated images [10]. Inference speed denotes the calculation time under the same running environment. Additionally, we introduce a novel metric, FGIS (fine-grained identity similarity), to assess ID quality at the region level. FGIS is computed as the average cosine similarity between DINO embeddings of the generated facial regions in reference and generated images. A higher FGIS value indicates increased ID fidelity in the generated facial regions.

### 4.1 Comparation Results

To demonstrate the effectiveness of ConsistentID, we conduct a comparative analysis against state-of-the-art methods, including Fastcomposer [54], IP-Adapter
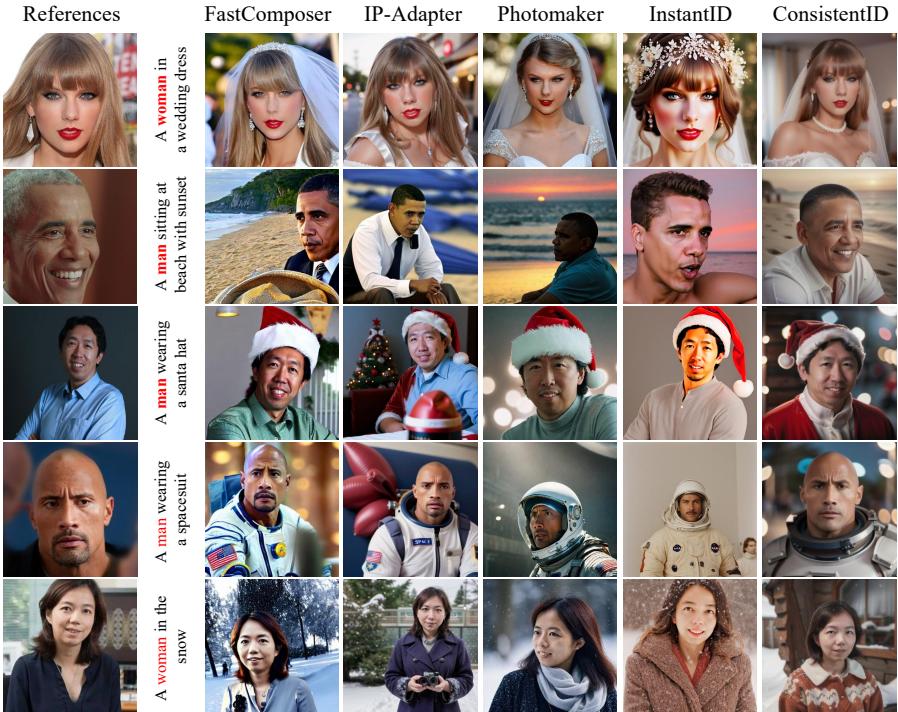
| References | FastComposer | IP-Adapter | Photomaker | InstantID | ConsistentID |

**Fig. 4:** Qualitative comparison of universal recontextualization samples is conducted, comparing our approach with other methods using five distinct identities and their corresponding prompts. Our ConsistentID exhibits a more powerful capability in high-quality generation, flexible editability, and strong identity fidelity.

[56], Photomaker [22], and InstantID [53]. Our focus is on personalized generation utilizing only one reference image. We utilize the officially provided models, use the default parameters for each method, and restrict the inference to a single reference image. In alignment with the Photomaker methodology, we employ the Mystyle [31] dataset for quantitative assessment and incorporate over ten identity datasets for visualization.

**Quantitative results:** Following Photomaker [22], we use the test dataset from Mystyle [31], using MLLM LLaVA1.5 to obtain facial descriptions during inference. The quantitative comparison is conducted under the universal recontextualization setting, utilizing a set of metrics to benchmark various aspects.

The results are showcased in Table 1. A thorough analysis of the table demonstrates that ConsistentID consistently outperforms other methods across most evaluated metrics, and surpasses other IP-Adapter-based methods in terms of generation efficiency. This is attributed to ConsistentID's fine-grained ID preservation capability and the efficiency of the lightweight multimodal facial prompt generator. Regarding the FID metric, the lower performance could be primarily attributed to the limited generative capability of the base model SD1.5.

**Visualized comparisons under different scenarios:** To visually demonstrate the advantages of ConsistentID, we present the text-edited generation results of all methods using reference images of five distinct identities in Fig-

| | CLIP-T ↑ | CLIP-I ↑ | DINO↑ | FaceSim ↑ | FGIS ↑ | FID ↓ | Speed (s) |
|---|---|---|---|---|---|---|---|
| Fastcomposer [54] | 27.8 | 67.0 | 68.4 | 75.2 | 77.7 | 372.8 | **10** |
| IP-Adapter [56] | 27.6 | <u>75.0</u> | 74.5 | 75.6 | 73.4 | 320.0 | 13 |
| Photomaker [22] | 30.7 | 71.7 | 72.6 | 69.3 | 73.2 | 336.5 | 17 |
| InstantID [53] | <u>30.3</u> | 68.2 | <u>77.6</u> | <u>76.5</u> | <u>78.3</u> | **271.9** | 19 |
| ConsistentID | **31.1** | **76.7** | **78.5** | **77.2** | **81.4** | <u>312.4</u> | 16 |

**Table 1:** Quantitative comparison of the universal recontextualization setting on the MyStyle test dataset. The benchmark metrics assessed text consistency (CLIP-T), the preservation of coarse- and fine-grained ID information (CLIP-I, DINO, FaceSIM, and FGIS), generation quality (FID), and inference efficiency (speed in seconds).
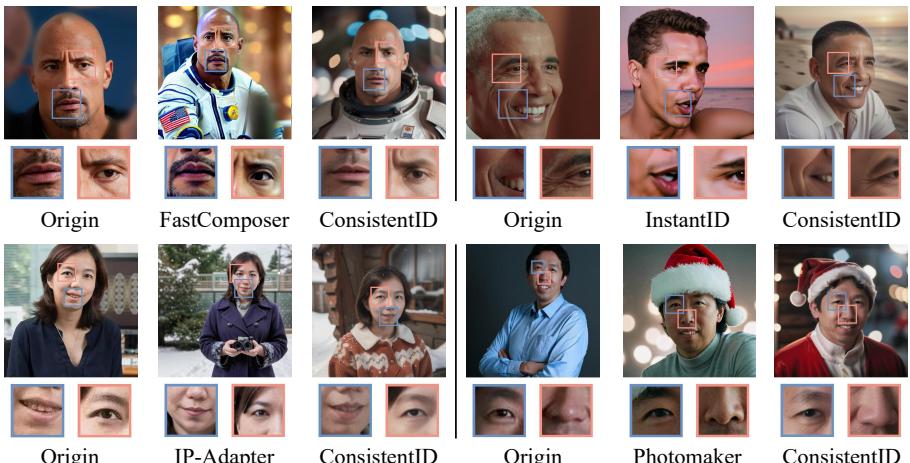


**Fig. 5:** Comparison of facial feature details between our method and existing approaches. Notably, the characters generated by our method exhibit superior ID consistency in facial features such as eyes, nose, and mouth.

ure 4. This visualization highlights ConsistentID's capability to produce vibrant and realistic images, with a particular emphasis on facial features. To further elucidate this observation, we selectively magnify and compare specific facial details across all methods in four identities, as depicted in Figure 5. Our model showcases exceptional ID preservation capabilities in facial details, especially in the eyes and nose, attributed to fine-grained multimodal prompts and facial regions' ID information.

To validate the accurate text understanding capability, we additionally show style-based and action-based text-edited results in Figure 6. We observe that the generated images from InstantID lack sufficient flexibility in facial poses. This limitation is likely attributed to Controlnet-based prompt insertion methods, which may easily overlook textual prompts. Simultaneously, we notice that, while Photomaker can accurately comprehend textual and visual prompts, it lacks the ID consistency of facial regions. In contrast, our ConsistentID achieves optimal generation results due to its precise understanding of textual and visual prompts. This further emphasizes the significance of multimodal fine-grained ID information. To fully show the advantages of our ConsistentID, more visualized comparisons are provided in the supplementary materials, including com-
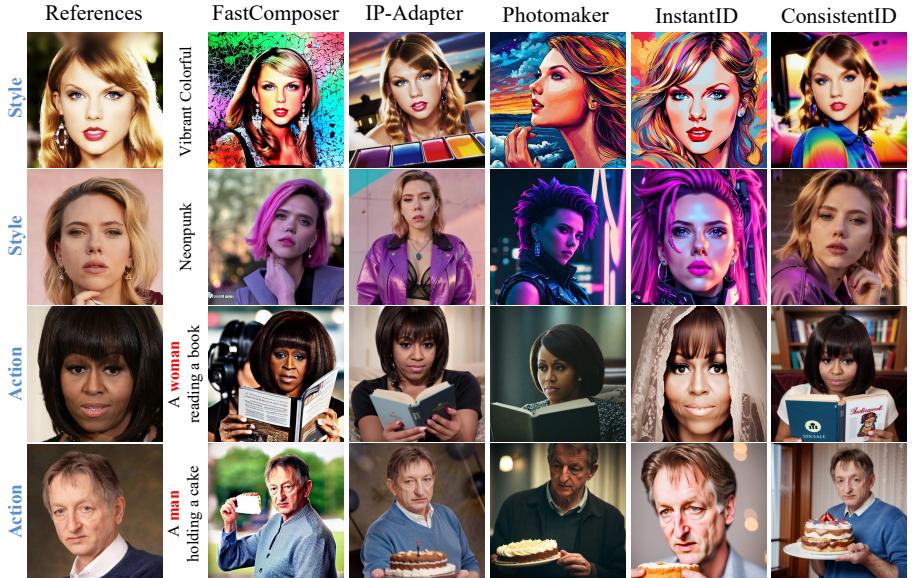
**Fig. 6:** Qualitative comparison of our model with other models on two special tasks: stylization and action instruction.



**Fig. 7:** Comparison of ConsistentID with IP-Adapter and its face version variants conditioned on different styles.

parative experiments with fine-tuning-based models Dreambooth [40], Textual Inversion [8], and CustomDiffusion [20].

Moreover, we compare ConsistentID with models specifically designed using IP-Adapter as the base model in Figure 7. From the figure, it is evident that this series of works currently falls short in achieving highly detailed ID preservation in facial areas without fine-grained textual and visual prompts. In contrast, ConsistentID exhibits a robust capability to preserve the integrity of facial ID and seamlessly blend it into various styles, utilizing multimodal fine-grained prompts. This comparison emphasizes the superiority of ConsistentID in retaining identity while simultaneously maintaining flexibility and control over style.

## 4.2   Human Study

We also investigate user preferences regarding image fidelity, fine-grained ID fidelity, and overall ID fidelity through surveys. In Figure 8, we present a visual-

ization of the proportion of total votes received by each method. Across all three metric dimensions, ConsistentID holds the most significant share.
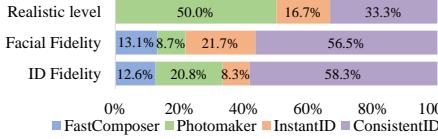


**Fig. 8:** User preferences across image fidelity, fine-grained ID fidelity, overall ID fidelity for different methods
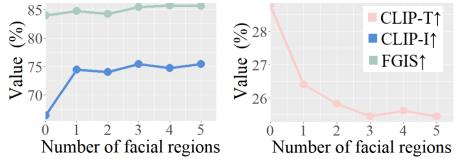


**Fig. 9:** Ablation study of the number of facial regions.
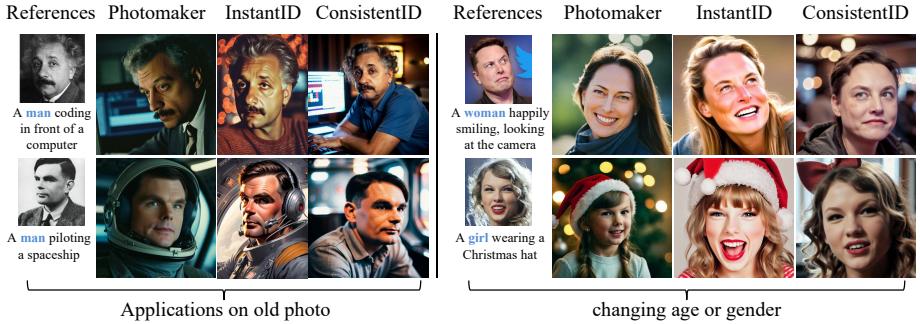
### 4.3   Ablation Study

**Facial ID Types:** We conducted an ablation study on facial ID, considering three variations: using solely overall facial ID, incorporating our designed fine-grained ID, and concurrently leveraging both facial ID and our designed fine-grained ID. Results in Table 2 show that using fine-grained ID features alone enhances ID consistency in facial regions, indicated by an increased DINO value. However, overall ID preservation, assessed by FaceSim, cannot be adequately maintained. Utilizing only overall facial ID enhances facial similarity, but facial feature consistency diminishes, performing significantly worse in CLIP-I and DINO metrics compared to using only facial features. Combining both overall facial ID and fine-grained ID leads to a balanced improvement in the model's overall ID fidelity compared to using only overall facial ID and fine-grained ID.

| $\mathcal{L}_{noise}$ | $\mathcal{L}_{loc}$ | CLIP-I ↑ | DINO ↑ | FGIS ↑ | LLaVA1.5 | CLIP-I ↑ | DINO ↑ | FGIS ↑ |
|---|---|---|---|---|---|---|---|---|
| ✓ | - | 66.4 | 77.8 | 82.9 | - | **75.5** | 86.1 | 85.6 |
| ✓ | ✓ | **75.5** | **86.1** | **85.6** | ✓ | 71.0 | **86.5** | 85.7 |

| Facial Feature | CLIP-I ↑ | DINO ↑ | FGIS ↑ | ImageProjection | CLIP-I ↑ | DINO ↑ | FGIS ↑ |
|---|---|---|---|---|---|---|---|
| Overall Facial Feature | 72.9 | 80.7 | 84.2 | - | 61.0 | 75.6 | 82.9 |
| Fine-grained Feature | 75.2 | **86.6** | 85.4 | ✓ | **75.5** | **86.1** | **85.6** |
| Overall Facial & Fine-grained Feature | **75.5** | 86.1 | **85.6** | | | | |

**Table 2:** Ablation study on ID features, loss functions, ImageProjection module, and the usage of LLaVA1.5 in inference.

**Facial attention localization strategy:** We investigated the effectiveness of facial attention localization strategies during training. The first strategy involves $\mathcal{L}_{noise}$, while the second strategy adds attention loss $\mathcal{L}_{loc}$. From Table 2, we observe that ConsistentID experiences a clear improvement in metrics related to facial feature consistency and fine-grained ID preservation when $\mathcal{L}_{loc}$ is considered. This confirms the effectiveness of maintaining ID consistency between facial regions and the entire face during the training process.

**Image projection module:** Additionally, we compared two training strategies. The first involves frozen weights of the image projection model and only training our designed FacialEncoder. The second strategy involves training both simultaneously. The results from Table 2 indicate that concurrently training ImageProjection brings the maximum benefits to the model. This is attributed to

References  Photomaker  InstantID  ConsistentID    References  Photomaker  InstantID  ConsistentID



A **man** coding in front of a computer

A **man** piloting a spaceship

A **woman** happily smiling, looking at the camera

A **girl** wearing a Christmas hat

Applications on old photo          changing age or gender

**Fig. 10:** The comparisons of two downstream applications.

our model being an ID preservation method of cooperative training with multi-modal text and image information.

**The LLaVA1.5 usage:** In Table 2 (bottom), we present comparative results conducted with textual descriptions using and not using LLaVA 1.5. When LLaVA1.5 is not utilized, we use the prompt 'The person has one nose, two eyes, two ears, and a mouth' to replace detailed descriptions of facial regions. We observe an improvement in the generated quality with the introduction of LLaVA1.5 in most metrics, with only CLIP-I showing degradation. This degradation is attributed to the local attention of CLIP, which tends to overlook finer facial details.

**Different Facial Areas' Number:** To explore the influence of different numbers of facial areas, we adhere to the sequence 'face, nose, eyes, ears, and mouth' and incrementally introduce the selected facial areas, as depicted in Figure 9. We note a progressive enhancement in image quality as the number of facial regions increases, attributed to the richer multi-modal prompts. However, with regard to the CLIP-T metric, detailed textual descriptions encompass a greater variety of objects, potentially leading to oversight by the CLIP model.

### 4.4 Applications

Following the approach proposed in Photomaker [22], we verify facial high-fidelity and naturalness through two downstream applications presented in Figure 10. These applications involve 'Change Age & Gender' and 'Bringing a person in an old photo into reality'. Compared to Photomaker and InstantID, our ConsistentID demonstrates robust capabilities in maintaining ID consistency, facial high-fidelity, and natural expressions. This superiority is attributed to the effective fine-grained ID preservation of facial features.

## 5    Conclusion and Limitation

In this work, we introduce ConsistentID, an innovative method designed to maintain identity consistency and capture diverse facial details. We have developed two novel modules: a multimodal facial prompt generator and an identity preservation network. The former is dedicated to generating multimodal facial prompts

by incorporating both visual and textual descriptions at the facial region level. The latter aims to ensure ID consistency in each facial area through a facial attention localization strategy, preventing the blending of ID information from different facial regions. By leveraging multimodal fine-grained prompts, our approach achieves remarkable identity consistency and facial realism using only a single facial image. Additionally, we present the FGID dataset, a comprehensive dataset containing fine-grained identity information and detailed facial descriptions essential for training the ConsistentID model. Experimental results demonstrate outstanding accuracy and diversity in personalized facial generation, surpassing existing methods on the MyStyle dataset.

**Limitations:** The utilization of MLLM in our approach may introduce limitations that could affect specific facets of model performance. The constraints posed by limited pose and expression may restrict the diversity of our method, impacting its capability to handle facial variations. These limitations underscore the necessity for in-depth discussion and exploration, specifically in addressing challenges related to pose, expression, and the integration of GPT-4V.

# References

1. Balaji, Y., Nah, S., Huang, X., Vahdat, A., Song, J., Kreis, K., Aittala, M., Aila, T., Laine, S., Catanzaro, B., et al.: ediffi: Text-to-image diffusion models with an ensemble of expert denoisers. arXiv preprint arXiv:2211.01324 (2022) 2

2. Beniaguev, D.: Synthetic faces high quality (sfhq) dataset (2022). https://doi.org/10.34740/kaggle/dsv/4737549, https://github.com/SelfishGene/SFHQ-dataset 8, 21

3. Cao, Q., Shen, L., Xie, W., Parkhi, O.M., Zisserman, A.: Vggface2: A dataset for recognising faces across pose and age. In: 2018 13th IEEE international conference on automatic face & gesture recognition (FG 2018). pp. 67–74. IEEE (2018) 3

4. Chefer, H., Alaluf, Y., Vinker, Y., Wolf, L., Cohen-Or, D.: Attend-and-excite: Attention-based semantic guidance for text-to-image diffusion models. ACM Transactions on Graphics (TOG) **42**(4), 1–10 (2023) 6

5. Chen, W., Hu, H., Li, Y., Ruiz, N., Jia, X., Chang, M.W., Cohen, W.W.: Subject-driven text-to-image generation via apprenticeship learning. Advances in Neural Information Processing Systems **36** (2024) 8

6. Cong, W., Zhang, J., Niu, L., Liu, L., Ling, Z., Li, W., Zhang, L.: Dovenet: Deep image harmonization via domain verification. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 8394–8403 (2020) 9

7. Deng, J., Guo, J., Xue, N., Zafeiriou, S.: Arcface: Additive angular margin loss for deep face recognition. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 4690–4699 (2019) 9

8. Gal, R., Alaluf, Y., Atzmon, Y., Patashnik, O., Bermano, A.H., Chechik, G., Cohen-Or, D.: An image is worth one word: Personalizing text-to-image generation using textual inversion. arXiv preprint arXiv:2208.01618 (2022) 2, 4, 5, 6, 9, 12

9. Gu, J., Wang, Y., Zhao, N., Fu, T.J., Xiong, W., Liu, Q., Zhang, Z., Zhang, H., Zhang, J., Jung, H., et al.: Photoswap: Personalized subject swapping in images. Advances in Neural Information Processing Systems **36** (2024) 5

10. Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., Hochreiter, S.: Gans trained by a two time-scale update rule converge to a local nash equilibrium. Advances in neural information processing systems **30** (2017) 9

11. Ho, J., Jain, A., Abbeel, P.: Denoising diffusion probabilistic models. Advances in neural information processing systems **33**, 6840–6851 (2020) 2

12. Hu, E.J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., Chen, W.: Lora: Low-rank adaptation of large language models. arXiv preprint arXiv:2106.09685 (2021) 4, 20

13. Huang, L., Chen, D., Liu, Y., Shen, Y., Zhao, D., Zhou, J.: Composer: Creative and controllable image synthesis with composable conditions. arXiv preprint arXiv:2302.09778 (2023) 4

14. Huang, Z., Chan, K.C., Jiang, Y., Liu, Z.: Collaborative diffusion for multi-modal face generation and editing. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 6080–6090 (2023) 2

15. Huang, Z., Li, Y.: Interpretable and accurate fine-grained recognition via region grouping. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 8662–8672 (2020) 6

16. Ju, X., Zeng, A., Zhao, C., Wang, J., Zhang, L., Xu, Q.: Humansd: A native skeleton-guided diffusion model for human image generation. arXiv preprint arXiv:2304.04269 (2023) 2

17. Karras, T., Laine, S., Aila, T.: A style-based generator architecture for generative adversarial networks. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 4401–4410 (2019) 5, 8, 9, 21

18. Kim, K., Kim, Y., Cho, S., Seo, J., Nam, J., Lee, K., Kim, S., Lee, K.: Diffface: Diffusion-based face swapping with facial guidance. arXiv preprint arXiv:2212.13344 (2022) 6

19. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980 (2014) 9

20. Kumari, N., Zhang, B., Zhang, R., Shechtman, E., Zhu, J.Y.: Multi-concept customization of text-to-image diffusion. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 1931–1941 (2023) 2, 4, 12, 20

21. Li, L., Zhang, T., Kang, Z., Jiang, X.: Mask-fpan: Semi-supervised face parsing in the wild with de-occlusion and uv gan. Computers & Graphics **116**, 185–193 (2023) 6

22. Li, Z., Cao, M., Wang, X., Qi, Z., Cheng, M.M., Shan, Y.: Photomaker: Customizing realistic human photos via stacked id embedding. arXiv preprint arXiv:2312.04461 (2023) 2, 4, 5, 6, 10, 11, 14

23. Liu, H., Li, C., Li, Y., Lee, Y.J.: Improved baselines with visual instruction tuning. arXiv preprint arXiv:2310.03744 (2023) 3, 5

24. Liu, R., Ma, B., Zhang, W., Hu, Z., Fan, C., Lv, T., Ding, Y., Cheng, X.: Towards a simultaneous and granular identity-expression control in personalized face generation. arXiv preprint arXiv:2401.01207 (2024) 5

25. Liu, X., Ren, J., Siarohin, A., Skorokhodov, I., Li, Y., Lin, D., Liu, X., Liu, Z., Tulyakov, S.: Hyperhuman: Hyper-realistic human generation with latent structural diffusion. arXiv preprint arXiv:2310.08579 (2023) 2, 5

26. Liu, Z., Li, M., Zhang, Y., Wang, C., Zhang, Q., Wang, J., Nie, Y.: Fine-grained face swapping via regional gan inversion. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 8578–8587 (2023) 5

27. Liu, Z., Luo, P., Wang, X., Tang, X.: Deep learning face attributes in the wild. In: Proceedings of the IEEE international conference on computer vision. pp. 3730–3738 (2015) 3, 8, 21

28. Meng, C., Song, Y., Song, J., Wu, J., Zhu, J.Y., Ermon, S.: Sdedit: Image synthesis and editing with stochastic differential equations. arXiv preprint arXiv:2108.01073 (2021) 5

29. Nasir, O.R., Jha, S.K., Grover, M.S., Yu, Y., Kumar, A., Shah, R.R.: Text2facegan: Face generation from fine grained textual descriptions. In: 2019 IEEE Fifth International Conference on Multimedia Big Data (BigMM). pp. 58–67. IEEE (2019) 5

30. Nichol, A., Dhariwal, P., Ramesh, A., Shyam, P., Mishkin, P., McGrew, B., Sutskever, I., Chen, M.: Glide: Towards photorealistic image generation and editing with text-guided diffusion models. arXiv preprint arXiv:2112.10741 (2021) 2, 4

31. Nitzan, Y., Aberman, K., He, Q., Liba, O., Yarom, M., Gandelsman, Y., Mosseri, I., Pritch, Y., Cohen-Or, D.: Mystyle: A personalized generative prior. ACM Transactions on Graphics (TOG) **41**(6), 1–10 (2022) 3, 10

32. Peebles, W., Xie, S.: Scalable diffusion models with transformers. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 4195–4205 (2023) 4

33. Podell, D., English, Z., Lacey, K., Blattmann, A., Dockhorn, T., Müller, J., Penna, J., Rombach, R.: Sdxl: Improving latent diffusion models for high-resolution image synthesis. arXiv preprint arXiv:2307.01952 (2023) 4

34. Qiu, Z., Liu, W., Feng, H., Xue, Y., Feng, Y., Liu, Z., Zhang, D., Weller, A., Schölkopf, B.: Controlling text-to-image diffusion by orthogonal finetuning. Advances in Neural Information Processing Systems **36** (2024) 6

35. Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al.: Learning transferable visual models from natural language supervision. In: International conference on machine learning. pp. 8748–8763. PMLR (2021) 4, 6, 9

36. Ramesh, A., Dhariwal, P., Nichol, A., Chu, C., Chen, M.: Hierarchical text-conditional image generation with clip latents. arXiv preprint arXiv:2204.06125 **1**(2), 3 (2022) 2

37. Rombach, R., Blattmann, A., Lorenz, D., Esser, P., Ommer, B.: High-resolution image synthesis with latent diffusion models. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 10684–10695 (2022) 2, 4, 20

38. Ronneberger, O., Fischer, P., Brox, T.: U-net: Convolutional networks for biomedical image segmentation. In: Medical Image Computing and Computer-Assisted Intervention–MICCAI 2015: 18th International Conference, Munich, Germany, October 5-9, 2015, Proceedings, Part III 18. pp. 234–241. Springer (2015) 4

39. Rosberg, F., Aksoy, E.E., Alonso-Fernandez, F., Englund, C.: Facedancer: pose- and occlusion-aware high fidelity face swapping. In: Proceedings of the IEEE/CVF winter conference on applications of computer vision. pp. 3454–3463 (2023) 5

40. Ruiz, N., Li, Y., Jampani, V., Pritch, Y., Rubinstein, M., Aberman, K.: Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 22500–22510 (2023) 4, 6, 9, 12, 20

41. Ruiz, N., Li, Y., Jampani, V., Wei, W., Hou, T., Pritch, Y., Wadhwa, N., Rubinstein, M., Aberman, K.: Hyperdreambooth: Hypernetworks for fast personalization of text-to-image models. arXiv preprint arXiv:2307.06949 (2023) 2

42. Saharia, C., Chan, W., Saxena, S., Li, L., Whang, J., Denton, E.L., Ghasemipour, K., Gontijo Lopes, R., Karagol Ayan, B., Salimans, T., et al.: Photorealistic text-to-image diffusion models with deep language understanding. Advances in Neural Information Processing Systems **35**, 36479–36494 (2022) 2

43. Schroff, F., Kalenichenko, D., Philbin, J.: Facenet: A unified embedding for face recognition and clustering. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 815–823 (2015) 9

44. Schuhmann, C., Beaumont, R., Vencu, R., Gordon, C., Wightman, R., Cherti, M., Coombes, T., Katta, A., Mullis, C., Wortsman, M., et al.: Laion-5b: An open large-scale dataset for training next generation image-text models. Advances in Neural Information Processing Systems **35**, 25278–25294 (2022) 9

45. Song, J., Meng, C., Ermon, S.: Denoising diffusion implicit models. arXiv preprint arXiv:2010.02502 (2020) 2, 9

46. Stypułkowski, M., Vougioukas, K., He, S., Zięba, M., Petridis, S., Pantic, M.: Diffused heads: Diffusion models beat gans on talking-face generation. In: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision. pp. 5091–5100 (2024) 2

47. Tumanyan, N., Geyer, M., Bagon, S., Dekel, T.: Plug-and-play diffusion features for text-driven image-to-image translation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 1921–1930 (2023) 6

48. Umirzakova, S., Whangbo, T.K.: Detailed feature extraction network-based fine-grained face segmentation. Knowledge-Based Systems **250**, 109036 (2022) 6

49. Valevski, D., Lumen, D., Matias, Y., Leviathan, Y.: Face0: Instantaneously conditioning a text-to-image model on a face. In: SIGGRAPH Asia 2023 Conference Papers. pp. 1–10 (2023) 2

50. Wan, L., Wan, J., Jin, Y., Tan, Z., Li, S.Z.: Fine-grained multi-attribute adversarial learning for face generation of age, gender and ethnicity. In: 2018 International Conference on Biometrics (ICB). pp. 98–103. IEEE (2018) 5

51. Wang, B., Zheng, H., Liang, X., Chen, Y., Lin, L., Yang, M.: Toward characteristic-preserving image-based virtual try-on network. In: Proceedings of the European conference on computer vision (ECCV). pp. 589–604 (2018) 2

52. Wang, K., Wu, Q., Song, L., Yang, Z., Wu, W., Qian, C., He, R., Qiao, Y., Loy, C.C.: Mead: A large-scale audio-visual dataset for emotional talking-face generation. In: European Conference on Computer Vision. pp. 700–717. Springer (2020) 3

53. Wang, Q., Bai, X., Wang, H., Qin, Z., Chen, A.: Instantid: Zero-shot identity-preserving generation in seconds. arXiv preprint arXiv:2401.07519 (2024) 2, 4, 5, 10, 11

54. Xiao, G., Yin, T., Freeman, W.T., Durand, F., Han, S.: Fastcomposer: Tuning-free multi-subject image generation with localized attention. arXiv preprint arXiv:2305.10431 (2023) 2, 4, 5, 6, 7, 9, 11

55. Xue, Z., Song, G., Guo, Q., Liu, B., Zong, Z., Liu, Y., Luo, P.: Raphael: Text-to-image generation via large mixture of diffusion paths. Advances in Neural Information Processing Systems **36** (2024) 2, 4

56. Ye, H., Zhang, J., Liu, S., Han, X., Yang, W.: Ip-adapter: Text compatible image prompt adapter for text-to-image diffusion models. arXiv preprint arXiv:2308.06721 (2023) 2, 4, 6, 10, 11

57. Yu, C., Wang, J., Peng, C., Gao, C., Yu, G., Sang, N.: Bisenet: Bilateral segmentation network for real-time semantic segmentation. In: Proceedings of the European conference on computer vision (ECCV). pp. 325–341 (2018) 6, 9

58. Yu, C., Lu, G., Zeng, Y., Sun, J., Liang, X., Li, H., Xu, Z., Xu, S., Zhang, W., Xu, H.: Towards high-fidelity text-guided 3d face generation and manipulation using only images. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 15326–15337 (2023) 6
59. Yu, J., Wang, Y., Zhao, C., Ghanem, B., Zhang, J.: Freedom: Training-free energy-guided conditional diffusion model. arXiv preprint arXiv:2303.09833 (2023) 6
60. Zhang, L., Rao, A., Agrawala, M.: Adding conditional control to text-to-image diffusion models. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 3836–3847 (2023) 2, 4, 6
61. Zheng, Y., Yang, H., Zhang, T., Bao, J., Chen, D., Huang, Y., Yuan, L., Chen, D., Zeng, M., Wen, F.: General facial representation learning in a visual-linguistic manner. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 18697–18709 (2022) 3, 8

# Appendix Materials of ConsistentID

Jiehui Huang[1], Xiao Dong[1], Wenhui Song[1], Hanhui Li[1], Jun Zhou[2], Yuhao Cheng[3], Shutao Liao[1], Long Chen[3], Yiqiang Yan[3], Shengcai Liao[4], and Xiaodan Liang[1]

[1]Shenzhen Campus of Sun Yat-sen University, [2]Zhuhai Campus of Sun Yat-sen University, [3]Lenovo Research, [4]Inception Institute of Artificial Intelligence

## 1 Compare with more fine-tuning-based models

We individually fine-tune CustomDiffusion [20], Dreambooth [40], and LoRA [12] using 192 face images of Obama and 158 of Taylor Swift. These fine-tuned models are then utilized for text-to-image inference. To ensure a fair comparison, all models are based on Stable Diffusion v1.5 [37]. In Figure 1, we compare the proposed ConsistentID with these fine-tuned models. It is noteworthy that our training of ConsistentID relies solely on a single reference image, yet it achieves comparable quality in synthesized images to methods trained with hundreds of reference images.

## 2 More ablation

**Attention loss $\mathcal{L}_{loc}$:** To further validate the effectiveness of our $\mathcal{L}_{loc}$, we present several comparisons using two different identities in Figure 2. From the figure, we draw the following conclusions: 1) The details of key facial areas, such as Li Feifei's eye shape, are well preserved. 2) The appearances of facial regions are highly consistent with the reference images, like eyes, ears, and mouths.

**Delay control:** In Figure 3, we provide visualized results to evaluate the impact of delay control during inference. The term 'merge step' denotes the first time step in which we incorporate fine-grained facial image features. It helps to control the balance between text prompts and face images. In general, as the 'merge step' increases, the influence of fine-grained facial image features gradually diminishes. For example, if the 'merge step' is set to 0, it indicates that fine-grained facial image features are dominant in the generation process and might result in semantic inconsistencies. On the contrary, setting the 'merge step' to 0 will maximize the guidance of text prompts, yet might harm the identity consistency.

To visually illustrate the impact of the 'merge step', we display the variation curves for the CLIP-I, CLIP-T, and FGIS metrics as the 'merge step' increases in Figure 4. From the figure, we observe a consistent trend where the textual control gradually strengthens with each increment of the 'merge step'.

**Fig. 1:** The comparisons with more fine-tuning-based models.

| | BiSeNet | InsightFace | FacialEncoder | UNet | Inference |
|---|---|---|---|---|---|
| Time (s) | 1 | 3 | 3 | 5 | 4 |

**Table 1:** Inference time of each module.

**Infer time of each module:** In Table 1, we present the inference time of each module, which adds up to 16 seconds for processing one image. It quantitatively demonstrates the competitive performance in terms of inference efficiency [1].

## 3    Dataset details

**Data source:** Our facial images are from three public datasets, including FFHQ [17], CelebA [27], and SFHQ [2]. We select 70,000, 30,000, and 424,258 images from these datasets, respectively. In a total of 524,258 images, 107,048 images have recognizable IDs. Figure 6 shows some examples of these images and their corresponding captions.

**Dataset pipeline:** Below is the detailed our dataset processing pipeline: Each image is initially fed into BiSeNet and LLaVA1.5 to obtain a fine-grained mask and facial descriptions. Let $I_{mask}$ denote the fine-grained image mask, which is a binary mask indicating our defined facial components. We then utilize $I_{mask}$ to

---

[1] During the inference phase, the LLaVa or ChatGPT$-4$v module is not necessary input. Effective outcomes can be achieved by solely utilizing predefined descriptors of the facial description 'face, nose, eyes, ears, and mouth'.

**Fig. 2:** Visualized results with or without using attention loss.



**Fig. 3:** Performance variations of CLIP-I, CLIP-T, and FGIS metrics with increasing 'merge step'.

segment out the corresponding regions from the original image and denote them as $I_{face}$. Meanwhile, we utilize the InsightFace model to extract facial identity features.

**Dataset characteristics:** In our dataset, we include 15 types of identities, which are listed in Table 2. Subsequently, in Figure 7, we display the distributions of gender and age across all identities. The figure illustrates a relative balance of both properties in our dataset.

**Dataset scenarios:** In table 3, we further provide all prompts used in 45 scenarios. The scenarios are divided into 4 categories based on the different applications, including Clothing&Accessory, Action, Background and Style.

| Evaluation IDs | | |
|---|---|---|
| ①Andrew Ng | ⑥Scarlett Johansson | ⑪Joe Biden |
| ②Barack Obama | ⑦Taylor Swift | ⑫ Kamala Harris |
| ③Dwayne Johnson | ⑧Albert Einstein | ⑬Kaming He |
| ④Fei-Fei Li | ⑨Elon Mask | ⑭Lecun Yann |
| ⑤Michelle Obama | ⑩Geoffrey Hinton | ⑮Sam Altman |

**Table 2:** ID names used for evaluation.

| Reference | Merge Step = 0 | Merge Step = 10 | Merge Step = 20 | Merge Step = 30 | Merge Step = 40 | Merge Step = 50 |

a person with a mountain in the background

a street art stencil of a person

a person on the beach

a person with a blue house in the background

**Fig. 4:** Visualized results under different 'merge steps'. 'Merge Step' indicates when to start adding facial image features to the text prompt.



**Fig. 5:** The statistical characteristics of age and gender distribution in the FGID training dataset.

# 4 More training details.

1. We define the keywords of $T_{\text{face}}$ as 'face', 'ears', 'eyes', 'nose', and 'mouth'. These keywords are then used to locate their positions in $I_{\text{face}}$, and a trigger word '<facial>' is inserted to replace the keywords at the matched positions. Next, the descriptions corresponding to facial regions are rearranged according to the order of facial feature keywords, while any descriptions unrelated to facial regions are eliminated. To resolve the problem of incomplete alignment between facial region descriptions and actual facial regions, we address two scenarios. If $T_{\text{face}}$ comprises fewer than 5 region descriptions, $I_{\text{face}}$ retains only the corresponding regions, with unidentified regions substituted by zero matrices. Conversely, if $T_{\text{face}}$ includes complete descriptions for 5 regions but $I_{\text{face}}$ lacks descriptions for all 5, the absent region is replaced with a zero matrix.

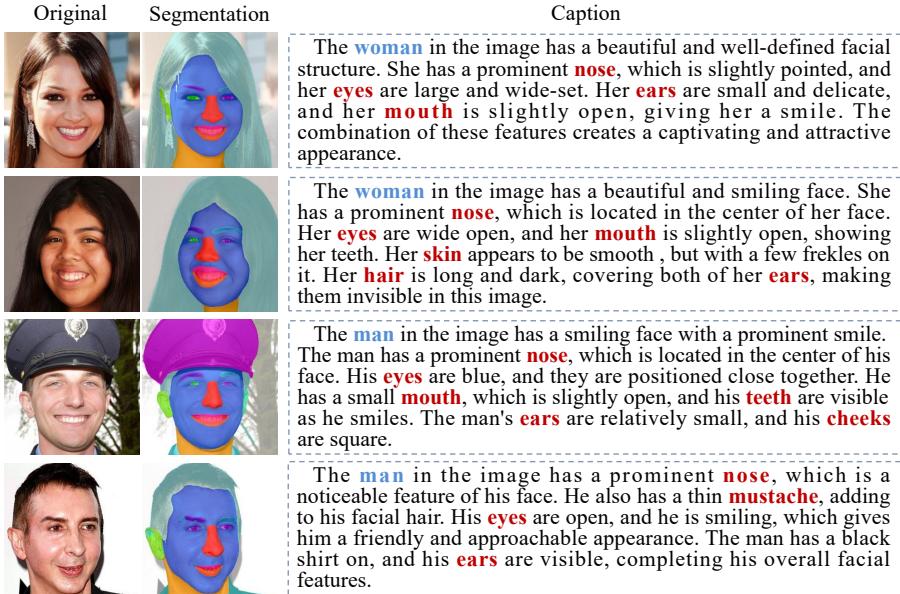| Original | Segmentation | Caption |
|---|---|---|
|  |  | The **woman** in the image has a beautiful and well-defined facial structure. She has a prominent **nose**, which is slightly pointed, and her **eyes** are large and wide-set. Her **ears** are small and delicate, and her **mouth** is slightly open, giving her a smile. The combination of these features creates a captivating and attractive appearance. |
|  |  | The **woman** in the image has a beautiful and smiling face. She has a prominent **nose**, which is located in the center of her face. Her **eyes** are wide open, and her **mouth** is slightly open, showing her teeth. Her **skin** appears to be smooth , but with a few frekles on it. Her **hair** is long and dark, covering both of her **ears**, making them invisible in this image. |
|  |  | The **man** in the image has a smiling face with a prominent smile. The man has a prominent **nose**, which is located in the center of his face. His **eyes** are blue, and they are positioned close together. He has a small **mouth**, which is slightly open, and his **teeth** are visible as he smiles. The man's **ears** are relatively small, and his **cheeks** are square. |
|  |  | The **man** in the image has a prominent **nose**, which is a noticeable feature of his face. He also has a thin **mustache**, adding to his facial hair. His **eyes** are open, and he is smiling, which gives him a friendly and approachable appearance. The man has a black shirt on, and his **ears** are visible, completing his overall facial features. |

**Fig. 6:** Several training data demos from our FGID dataset.

2. $T_{\text{face}}$ with trigger words undergoes ID encoding via a tokenizer, resulting in an encoded vector *input_ids* of length 77. The positions of trigger word representations in *input_ids* are recorded for FacialEncoder model localization during trigger word replacement.

3. To extract prior information for overall face ID, we referred to the IP-Adapter model and employed the InsightFace model to extract facial ID features for all images.

# 5    More visualization

**Downstream applications:** In Figures 7, 8, 9, and 10, we present additional visualized results to demonstrate the capabilities of our model in high-fidelity and flexible editing across various recontextualization scenarios, including age modification, identity mixing, and gender transformation.

# 6    Discussion of ethical principles:

In the work, we introduce ConsistentID, a method for generating high-quality facial images while preserving identity fidelity. Our approach emphasizes efficiency, diversity, and controllability in facial generation tasks, serving as a robust baseline for academic research. However, the widespread adoption of such technology raises ethical concerns regarding privacy, misinformation, and potential misuse. We advocate for the responsible development and use of these tools, emphasizing the importance of ethical guidelines to ensure their safe and ethical application in computer vision.

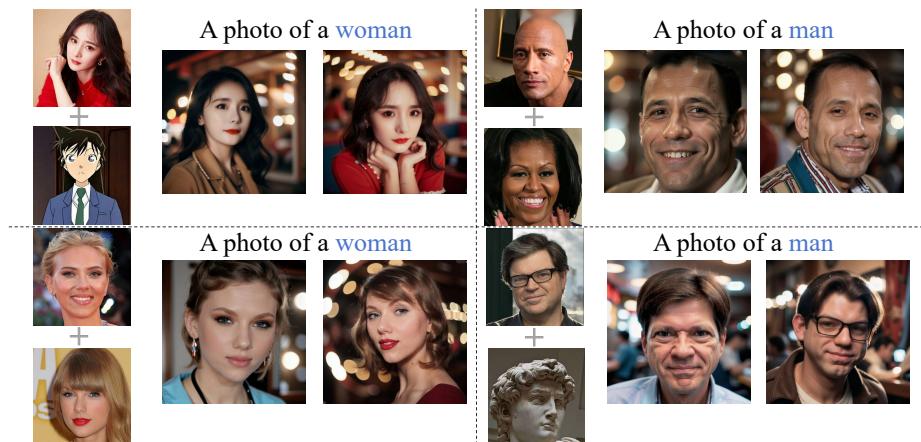**Fig. 7:** Additional application cases of ConsistentID for altering the age attribute of a character.



**Fig. 8:** Additional application cases of ConsistentID for identity confusion. Utilizing the overall facial ID feature of one character (top). Leveraging the fine-grained multimodal features from another character (bottom).
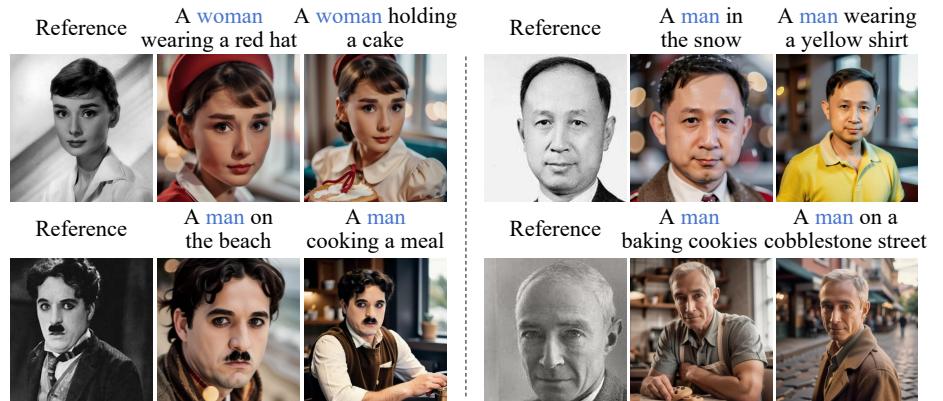


**Fig. 9:** Additional application cases of ConsistentID for bringing old photos back to life.

|  Reference | IP-Adapter | Photomaker | InstantID | ConsistentID |

a person img in a chef outfit

a person img working out at the gym

a person img in the jungle

a person img wearing a yellow shirt

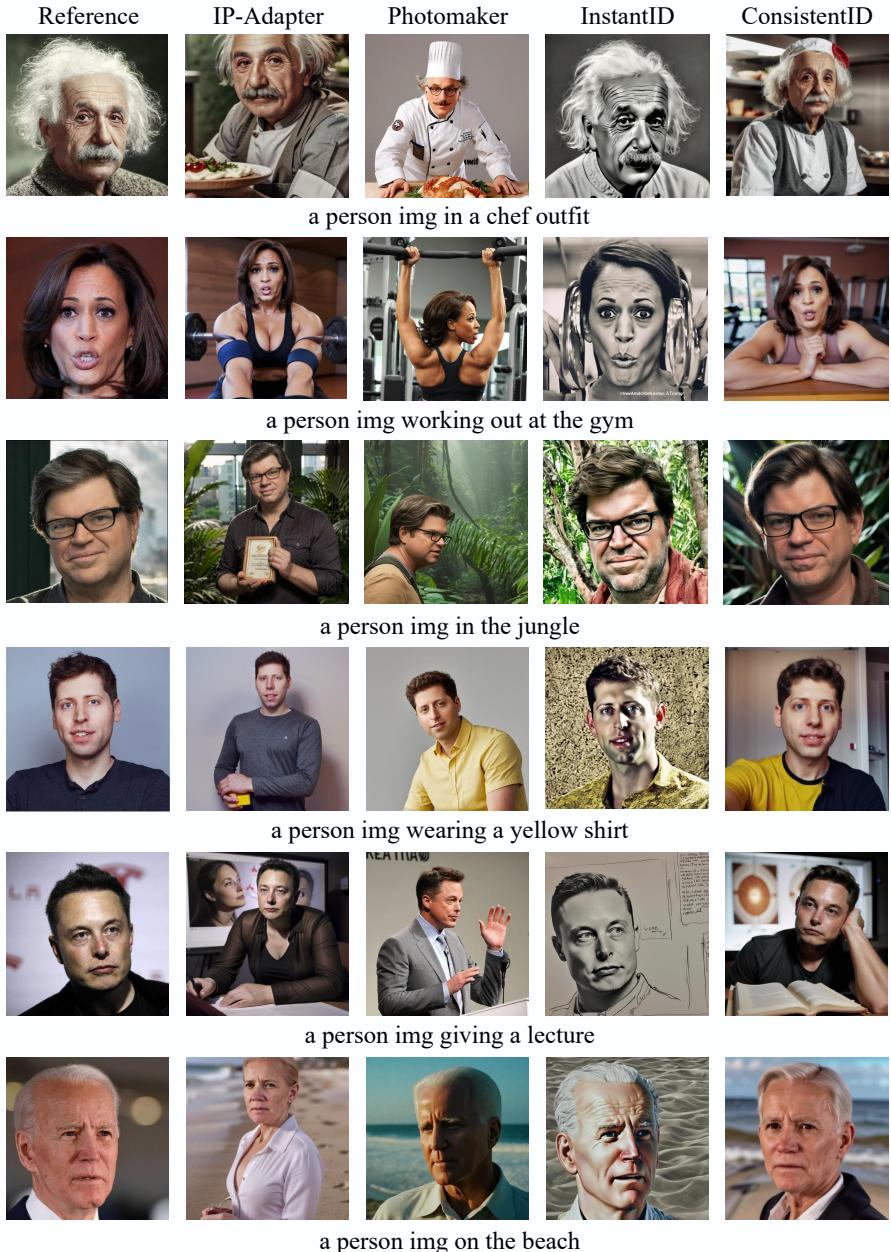a person img giving a lecture

a person img on the beach

**Fig. 10:** Extended visualization in re-contextualization settings. These examples demonstrate the high-identity fidelity and text editing capability of ConsistentID.

| Category | Prompt |
|---|---|
| Clothing& Accessory | a \<class word\> wearing a red hat |
| | a \<class word\> wearing a santa hat |
| | a \<class word\> wearing a rainbow scarf |
| | a \<class word\> wearing a black top hat and a monocle |
| | a \<class word\> in a chef outfit |
| | a \<class word\> in a firefighter outfit |
| | a \<class word\> in a police outfit |
| | a \<class word\> wearing pink glasses |
| | a \<class word\> wearing a yellow shirt |
| | a \<class word\> in a purple wizard outfit |
| Background | a \<class word\> in the jungle |
| | a \<class word\> in the snow |
| | a \<class word\> on the beach |
| | a \<class word\> on a cobblestone street |
| | a \<class word\> on top of pink fabric |
| | a \<class word\> on top of a wooden floor |
| | a \<class word\> with a city in the background |
| | a \<class word\> with a mountain in the background |
| | a \<class word\> with a blue house in the background |
| | a \<class word\> on top of a purple rug in a forest |
| Action | a \<class word\> holding a glass of wine |
| | a \<class word\> holding a piece of cake |
| | a \<class word\> giving a lecture |
| | a \<class word\> reading a book |
| | a \<class word\> gardening in the backyard |
| | a \<class word\> cooking a meal |
| | a \<class word\> working out at the gym |
| | a \<class word\> walking the dog |
| | a \<class word\> baking cookies |
| | a \<class word\> wearing a doctoral cap |
| | a \<class word\> wearing a spacesuit |
| | a \<class word\> wearing sunglasses and necklace |
| | a \<class word\> coding in front of a computer |
| | a \<class word\> in a helmet and vest riding a motorcycle |
| Style | a painting of a \<class word\> in the style of Banksy |
| | a painting of a \<class word\> in the style of Vincent Van Gogh |
| | a colorful graffiti painting of a \<class word\> |
| | a watercolor painting of a \<class word\> |
| | a Greek marble sculpture of a \<class word\> |
| | a street art mural of a \<class word\> |
| | a black and white photograph of a \<class word\> |
| | a pointillism painting of a \<class word\> |
| | a Japanese woodblock print of a \<class word\> |
| | a street art stencil of a \<class word\> |

**Table 3:** Evaluation text prompts are categorized by Clothing&Accessories, Background, Action, and Style. During inference, the term 'class' will be substituted with 'man', 'woman', 'girl', etc. For each identity and prompt, we generated 1,000 images randomly for evaluation.