

# DMID: DYNAMIC MASK ATTENTION FOR HIGH-FIDELITY IDENTITY PRESERVATION UNDER LIMITED DATA

**Anonymous authors**

Paper under double-blind review



Figure 1: DMID ensures high-fidelity identity preservation while preserving textual semantics. On the left are examples of some famous figures, and on the right are examples of age editing.

## ABSTRACT

We present Dynamic Mask Attention for High-Fidelity Identity Preservation under Limited Data (DMID), which aims to precisely reconstruct fine-grained identity features under scarce data conditions while alleviating conflicts between textual and conditional semantics. At its core, DMID employs a Variational Autoencoder (VAE) for meticulous identity encoding and introduces a **Dynamic Attention Mask mechanism**, coupled with **Distribution Consistency Loss** and **Identity Mask Loss**, ensuring identity fidelity while mitigating semantic conflicts. To further reduce annotation and training costs, we have designed an efficient data construction pipeline. Furthermore, our method enables the dynamic adjustment of the **AttnMask strength factor** during inference, ensuring precise modifications and fine-grained control over identity features and semantics across various scenarios. The training process is divided into three stages: (1) identity embedding stage, (2) dynamic attention mask learning stage, and (3) Diffusion-DPO post-training stage. Evaluated on our newly constructed ID Benchmark, DMID achieves state-of-the-art performance in both identity consistency and textual semantics, demonstrating its strong competitiveness in data-limited scenarios. Among them, the parameter count of AttnMaskNet is only approximately 1% of that of Flux.1-dev.

## 1 INTRODUCTION

In text-to-image (T2I) generation Esser et al. (2024); Rombach et al. (2022); Peebles & Xie (2023); Chen et al. (2023); Betker et al. (2023), Transformer-based diffusion models such as Flux Labs (2024) can produce high-fidelity images with complex semantics. However, maintaining consistent human identity (ID) across diverse generated scenes remains challenging, as existing techniques struggle to balance identity fidelity and textual semantics.

Early personalization methods, such as LoRA Devalal & Karthikeyan (2018) and DreamBooth Ruiz et al. (2023), rely on tens of target identity images for fine-tuning, which poses challenges to the robustness of identity consistency tasks. Recent approaches Li et al. (2024); Peng et al. (2024); Xiao

054 et al. (2025); Yuan et al. (2023); Han et al. (2024) including IP-Adapter Ye et al. (2023) and InstantID  
055 Wang et al. (2024b) significantly improve ID consistency by extracting identity features through  
056 auxiliary facial encoders and injecting them into generative models. While enabling personalized  
057 generation with single/few reference images, they suffer from inherent information compression  
058 during facial encoding that causes identity detail loss. Moreover, the introduced conditioning tokens  
059 frequently conflict with textual semantic tokens, reducing prompt responsiveness.

060 To address semantic degradation, Pulid-Flux Guo et al. (2024) introduces semantic alignment loss  
061 and layout alignment loss, while InfiniteYou Jiang et al. (2025) designs InfuseNet for improved  
062 fusion of facial encodings with foundation models. Despite these advances, effective alignment  
063 between external facial features and the foundation model’s latent space demands 100K to 1M high-  
064 quality training samples. The substantial data and computational costs severely limit reproducibility  
065 and practical deployment, highlighting the need for robust low-cost solutions.

066 OmniControl Tan et al. (2024) adopts a distinct approach: it directly encodes an additional condition  
067 image into the feature space of the foundation model using a pre-trained VAE encoder Kingma et al.  
068 (2013), thus unifying the encodings of condition image and noise image in a shared space. Similar  
069 structures were also used in earlier works on ReferenceNet Hu (2024); Xu et al. (2024b); Tian et al.  
070 (2024); Chang et al. (2023); Xu et al. (2024a); Choi et al. (2024); Wang et al. (2024c); Huang  
071 et al. (2024); Gu et al. (2024); Zhang et al. (2023). This strategy, which ensures high consistency  
072 between the encodings of conditional input and noise, significantly enhances the training efficiency  
073 and consistency of the details. Although the method demonstrates great potential for various tasks,  
074 such high consistency leads to the issue of reference image replication. Additionally, the number of  
075 conditional tokens is much greater than in earlier work such as IP-Adapter Ye et al. (2023), which  
076 exacerbates conflicts between textual semantics and image semantics.

077 To preserve high facial identity consistency while minimizing semantic interference, we propose a  
078 high-fidelity, identity-consistent diffusion model driven by dynamic attention masks. Our frame-  
079 work employs a multistage learning approach incorporating a dynamic attention masking mecha-  
080 nism guided by two losses: (1) **Distribution Consistency Loss**, which aligns conditional and tex-  
081 tual feature distributions; and (2) **Identity Mask Loss**, which directly optimizes facial similarity  
082 between generated and target ID image.

083 This mechanism compresses identity information into a facial-specific subspace and dynamically  
084 masks conflict regions between tokens, achieving robust identity preservation with minimal training  
085 data (approximately 40,000 high-quality image pairs). We also develop a specialized data acquisi-  
086 tion pipeline to resolve semantic conflicts and source image replication. Our principal contributions:

087 **Consistency-driven dynamic attention masking:** A novel masking scheme utilizing distribution  
088 consistency loss and identity mask loss, significantly enhancing face similarity while maintaining  
089 textual semantic consistency. We introduce the AttnMask strength factor, which dynamically ad-  
090 justs attention levels during inference, enabling a precise balance between identity preservation and  
091 semantic consistency, thus allowing fine-grained control and preventing semantic distortion due to  
092 over-attention.

093 **Efficient multi-stage training framework:** A three-phase approach leveraging unified VAE encod-  
094 ing:

095 *Stage 1 (Identity Embedding):* Learns identity embeddings from condition image.

096 *Stage 2 (Dynamic Mask Optimization):* Jointly trains AttnMaskNet using both consistency losses.

097 *Stage 3 (Preference Alignment):* We adopt Diffusion-DPO Wallace et al. (2024) training to enhance  
098 the stability of model generation.

099 **ID consistency benchmark:** We establish and open-source a standardized evaluation benchmark  
100 spanning diverse demographics (race, gender, age) with standardized protocols.

101  
102  
103  
104  
105  
106  
107

2 RELATED WORK

2.1 ID CONSISTENT GENERATION

To maintain specified identity (ID) in text-to-image generation, mainstream approaches focus on feature injection and model optimization. Feature injection adapters (e.g., IP-Adapter Ye et al. (2023), InstantID Wang et al. (2024b)) utilize pre-trained encoders to extract facial features, injecting them into foundation models (e.g., SDXL Podell et al. (2023)) via lightweight plugins (e.g., cross-attention). However, feature compression causes detail loss, while additional conditioning tokens frequently conflict with text tokens. Loss optimization methods (e.g., PuLID Guo et al. (2024)) explicitly incorporate ID loss based on face recognition models (e.g., ArcFace Deng et al. (2019)), directly optimizing generated face similarity to achieve quantifiable metric improvements. Architectural innovations explore new pathways: FlashFace Zhang et al. (2024) trains end-to-end isomorphic UNet encoders for fine-grained fusion, while ACE++ Mao et al. (2025) unifies conditional and noise image in DiT architecture by mapping to VAE latent space Kingma et al. (2013). Although enhancing detail consistency, this tends to cause reference image over-replication and exacerbates semantic conflicts due to excessive conditional tokens.

2.2 HUMAN PREFERENCE ALIGNMENT IN DIFFUSION MODELS

Inspired by Reinforcement Learning from Human Feedback (RLHF) paradigms in large language models (LLMs), alignment research for diffusion models has expanded rapidly. DRAFT Clark et al. (2023) and AlignProp Prabhudesai et al. (2023) integrate reward model gradients directly into training. Diffusion-DPO Wallace et al. (2024) adopts an offline approach using preference pairs, eliminating reward models as a simpler RLHF alternative. Concurrently, Flow-GRPO Liu et al. (2025) and DanceGRPO Xue et al. (2025), MixGRPOLi et al. (2025) embed online reinforcement learning within flow-matching frameworks, demonstrating significant gains in human preference tasks.

3 METHODS

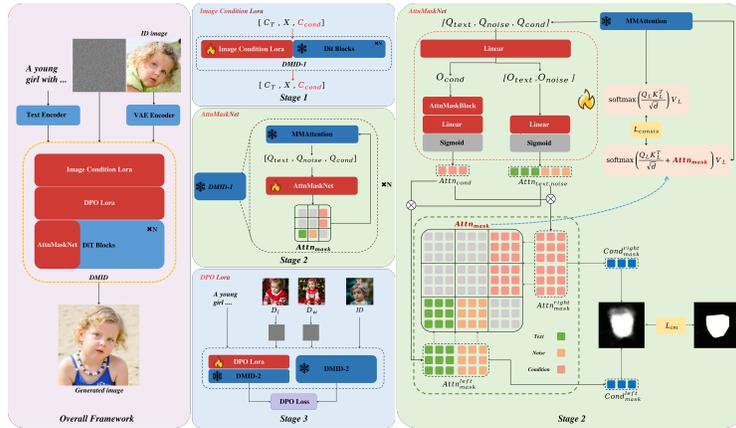


Figure 2: **Overall training framework of DMID**, consisting of three core stages: 1) The identity embedding stage; 2) The dynamic attention mask learning stage; 3) The Diffusion-DPO post-training stage. The AttnMaskNetBlock primarily comprises residual blocks and convolutional layers.

3.1 ENSEMBLE-BASED CONDITIONAL CONTROL LORA

To achieve ID conditioned image embedding, we adopt the framework of Ominicontrol Tan et al. (2024) for the first-stage training. Specifically, we employ Rectified Flows Liu et al. (2022) as the forward sampling process and use Conditional Flow Matching Lipman et al. (2022) as the optimiza-

tion objective, with the loss function defined as:

$$\mathcal{L}_{\text{CFM}} = \mathbb{E}_{t, p_t(z|\epsilon), p(\epsilon)} \|v_\theta(z, t) - u_t(z|\epsilon)\|_2^2, \quad (1)$$

where  $\epsilon \sim \mathcal{N}(0, I)$  and  $v_\theta$  is parameterized by neural network weights. The detailed network structure is illustrated in Stage 1 of Figure 2. First, the condition image is encoded via a VAE to obtain condition image tokens  $C_{\text{cond}} \in \mathbb{R}^{K \times d}$ . These tokens are then integrated with text tokens  $C_T \in \mathbb{R}^{M \times d}$  and image tokens  $X \in \mathbb{R}^{N \times d}$  into a unified space  $L = [C_T, X, C_{\text{cond}}]$  to allow interactions between multiple conditions. To further enhance flexibility, Ominicontrol proposes a method for manually controlling the strength of condition image. It constructs an attention mask matrix  $Attn(\alpha)$  using a given strength factor  $\alpha$  to adjust attention weights between condition image tokens and other tokens:

$$\text{MMAttention}(L) = \sigma \left( \frac{Q_L K_L^\top}{\sqrt{d}} + Attn(\alpha) \right) V_L, \quad (2)$$

where  $\sigma(\cdot) \triangleq \text{softmax}(\cdot)$  denotes the softmax function, and  $Attn(\alpha)$  is a structured masking matrix that selectively modifies attention weights for the condition image tokens. Specifically,  $Attn(\alpha)$  is defined as:

$$Attn(\alpha) = \begin{bmatrix} 0_{M \times M} & 0_{M \times N} & \alpha_{M \times K} \\ 0_{N \times M} & 0_{N \times N} & \alpha_{N \times K} \\ \alpha_{K \times M} & \alpha_{K \times N} & 0_{K \times K} \end{bmatrix}. \quad (3)$$

While the strength factor setting effectively controls the condition image effect, it suffers from several drawbacks. First, the strength factor operates globally across the entire condition image, leading to re-weighting of irrelevant regions that should be suppressed. In the face identity preservation task, the current method assigns equal weight to the face and the background. In fact, attention should be focused on the face area to accurately maintain identity consistency. Second, when the strength factor is set too large, it will seriously damage the original attention distribution and lead to the loss of text semantic information.

To address these issues, we design a dynamic Attention Mask learning scheme. The model learns an adaptive Attention Mask to enhance attention to the main subject while ensuring no significant loss of textual semantics.

### 3.2 DYNAMIC ATTENTION MASK LEARNING

As shown in Stage 2 of Figure 2, we take the query matrix from the attention module as input to AttnMaskNet and finally obtain an attention mask that acts on the condition image region. Its function is to improve face similarity while keeping the text and image as aligned as possible. As mentioned earlier, over-enhancement of some regions will disrupt the overall distribution, leading to text-image misalignment. Therefore, our main idea is to help the model suppress redundant attention regions and enhance key regions, which can maximize the attention enhancement while mitigating the loss of semantic information. We have designed corresponding components from the aspects of network structure and loss functions.

#### 3.2.1 ATTNMASKNET

As shown in Figure 2, the detailed workflow of AttnMaskNet is illustrated. Although our focus is on modifying the attention of condition image, we need to consider the original attention relationships between text, noise, and condition image. To preserve the inherent attention relationships among different conditional tokens in the model, we concatenate the query matrices of text, noise, and condition image from the attention module as the input to the network  $[Q_{\text{text}}, Q_{\text{noise}}, Q_{\text{cond}}] \in \mathbb{R}^{L_1 \times d_1}$  where  $L_1 = M + N + K$ .

The projected output is divided into the text-noise part  $Output_1 = [O_{\text{text}}, O_{\text{noise}}] \in \mathbb{R}^{L_2 \times d_2}$  (where  $L_2 = M + N$ ) and the condition image part  $Output_2 = O_{\text{cond}} \in \mathbb{R}^{K \times d_2}$ . We simply project the text-noise part further to retain their original attention information, ultimately obtaining  $Attn_{\text{text, noise}} \in \{\mathbf{x} \in \mathbb{R}^{L_2} \mid 0 \leq x_i \leq 1, \forall i = 1, \dots, L_2\}$ , where each element represents the strength of attention of each token in text and noise towards the condition image.

For the condition image part, we aim to enable it to focus precisely and effectively on key regions while masking non-critical ones. Attention modules can effectively capture global information, but

are relatively weak in local spatial information. To obtain more refined spatial features,  $Output_2$  is reshaped into  $\mathbb{R}^{d_2 \times \sqrt{K} \times \sqrt{K}}$  and fed into a residual network. Finally, the output of the residual network is reshaped and projected to obtain  $Attn_{cond} \in \{\mathbf{x} \in \mathbb{R}^K \mid 0 \leq x_i \leq 1, \forall i = 1, \dots, K\}$ .

Meanwhile, to ensure the attention mask can effectively enhance or suppress attention weights, we map the  $Attn_{cond}$  of each condition to the interval  $[-\beta, \beta]$ :

$$Attn'_{cond} = 2\beta \cdot Attn_{cond} - \beta, \quad (4)$$

where  $Attn_{cond}$  denotes the original matrix with elements in  $[0, 1]$ , and  $\beta$  is the strength factor controlling the range of the transformed values ( $\beta$  is set to 1 by default unless explicitly stated otherwise.). The same transformation yields  $Attn_{text,noise} \in \mathbb{R}^{L_2}$ .

As shown in Figure 2, we compute the regions of  $Attn_{mask}$  corresponding to the condition image, i.e., the upper-right region  $Attn_{mask}^{right} = Attn_{text,noise}^\top Attn'_{cond}$ ,  $Attn_{mask}^{right} \in \mathbb{R}^{L_2 \times K}$  and the lower-left region  $Attn_{mask}^{left} = Attn_{cond}^\top Attn'_{text,noise}$ ,  $Attn_{mask}^{left} \in \mathbb{R}^{K \times L_2}$ . The resulting attention mask is:

$$Attn_{mask} = \begin{bmatrix} 0_{(M+N) \times (M+N)} & Attn_{mask}^{right} \\ Attn_{mask}^{left} & 0_{K \times K} \end{bmatrix}, \quad (5)$$

The obtained  $Attn_{mask}$  acts on the attention map as a weight matrix, as shown in Equation 5. Here,  $Attn_{mask}^{right}$  and  $Attn_{mask}^{left}$  preserve the original attention relationships between text, noise, and the condition image while learning the information that needs to be enhanced or suppressed. This avoids disrupting the original self-attention distribution of the model and helps correct attention effectively. As a branch network, AttnMaskNet is incorporated into the attention module at every layer of the model. To guide AttnMaskNet toward our desired learning direction, we further design an identity mask Loss and a distribution consistency loss.

### 3.2.2 IDENTITY MASK LOSS

To help the model effectively learn key and non-key regions, we introduce a face mask as prior information to constrain the learning process of  $Attn_{mask}$ . As shown in Stage 2 of Figure 2, we first average  $Attn_{mask}^{right}$  and  $Attn_{mask}^{left}$  over the text and noise dimensions, yielding the condition image vectors  $Cond_{mask}^{right}$  and  $Cond_{mask}^{left}$  with dimension  $\mathbb{R}^K$ . These vectors are then resized to match the spatial size of the face mask, and the pixel-wise error between them is computed using binary cross-entropy Wang et al. (2024d):

$$\mathcal{L}_{im} = -\frac{1}{K} \sum_{i=1}^{\sqrt{K}} \sum_{j=1}^{\sqrt{K}} \hat{f}_{i,j} \log(\hat{a}_{i,j}) + (1 - \hat{f}_{i,j}) \log(1 - \hat{a}_{i,j}), \quad (6)$$

where  $\hat{f}_{i,j}$  denotes the pixel value of the face mask,  $\hat{a}_{i,j}$  denotes the pixel value corresponding to the reshaped condition image vector, and  $\sqrt{K}$  is the spatial dimension size.

### 3.2.3 DISTRIBUTION CONSISTENCY LOSS

ID consistency and text-image alignment are inherently conflicting. Essentially, when the model over-focuses on the condition image, the original attention distribution is disrupted, leading to loss of semantic information. Experimental observations show that increasing the scale of the condition image improves face similarity but weakens textual semantics. Expanding the softmax operation in Equation 2, the probability value for the text region can be expressed as:

$$M_{ik} = \frac{e^{x_{ik}}}{\sum_{j=1}^M e^{x_{ij}} + \sum_{j=M+1}^{M+N} e^{x_{ij}} + \sum_{j=M+N+1}^{L_1} e^{x_{ij}+a_{ij}}}, \quad (7)$$

where  $0 < i < L_1, 0 < k < M$ . Let  $S_c = \sum_{j=M+N+1}^{L_1} e^{x_{ij}+a_{ij}}$ , and  $a_{ij}$  denotes an element in  $Attn_{mask}$ . When  $a_{ij} \rightarrow -\infty$ , the expression degenerates to the standard text-to-image formulation.

It can be seen from the above formula that an increase in  $S_c$  leads to a decrease in  $M_{ik}$  ( $0 < k < M$ ). To reduce the weakening of textual semantics by adjusting  $a_{ij}$ , an effective method is to keep  $S_c$  unchanged. Therefore, we introduce a distribution consistency loss:

$$\mathcal{L}_{\text{consis}} = \frac{1}{L_1} \sum_{i=1}^{L_1} \left\| \sum_{k=M+N+1}^{L_1} M_{ik} - \sum_{k=M+N+1}^{L_1} M_{ik}^* \right\|, \quad (8)$$

where  $M_{ik}^*$  denotes the probability value when  $a_{ij} = 0$ , which is also the target value of the loss function. This regularizes the sum of probability values in the condition image region, indirectly preserving the probability of the text region to alleviate semantic conflicts.

The final optimization objective:

$$\mathcal{L}_{\text{stage2}} = \mathcal{L}_{\text{CFM}} + \lambda (\mathcal{L}_{\text{im}} + \mathcal{L}_{\text{consis}}), \quad (9)$$

where  $\lambda$  is a hyperparameter ( $0 < \lambda < 1$ ) that balances the weights of identity mask loss and distribution consistency loss.

### 3.3 DIFFUSION-DPO

To address the instability issue caused by the absence of perceptual loss (e.g. PuLID) in VAE-based image encoding during training, where the similarity of generated results fluctuates significantly under the same image conditions, this paper employs Diffusion-DPO for post-training refinement. As shown in Figure 2, Diffusion-DPO requires offline construction of paired preference data. To this end, we design a comprehensive scoring model:

$$R = R_{\text{CosSim}} + \mu R_{\text{CLIP}}, \quad (10)$$

where  $R_{\text{CosSim}}$  denotes the facial similarity metric,  $R_{\text{CLIP}}$  represents the semantic consistency metric from CLIP model Radford et al. (2021), and  $\mu$  balances their contributions. During the offline phase, we generate multiple outputs for the same input using different random seeds, compute their scores, and select the highest-scoring output as the winning sample  $D_w$  and the lowest-scoring one as the losing sample  $D_l$  to form paired preference data for Diffusion-DPO optimization.

## 4 EXPERIMENT

### 4.1 EXPERIMENTAL SETUP

#### 4.1.1 IMPLEMENTATION DETAILS

We implement DMID using PyTorch and HuggingFace Diffusers. The DiT base model is FLUX.1-dev. All experiments were conducted on 2xNVIDIA H100 GPUs. The training batch size is 1 with gradient accumulation steps of 2, and target image size is 768x768. For Stage 1 training, we use the Prodigy optimizer Mishchenko & Defazio (2024) with safe warmup and bias correction, setting weight decay to 0.01 for 30,000 iterations. In Stage 2, we train for 2,000 iterations: the first 1,500 steps use weight coefficient  $\lambda = 1$ , and the last 500 steps use  $\lambda = 0.1$ . For Stage 3 DPO training, we fix  $W = 5000$ ,  $\mu = 4$  and run 2,500 iterations with AdamW optimizer Loshchilov & Hutter (2017) (learning rate=0.0001,  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ ).

#### 4.1.2 TRAINING DATASET

Through our data construction pipeline, we obtain around 20,000 image pairs—totaling 40,000 individual images—covering 12,000 unique IDs. For video generation, we use Wan2.1 Wan et al. (2025) for batch image-to-video conversion, sampling frame pairs at 20fps. For image augmentation, ACE++ generates diverse pairs through different mask inputs. Low-similarity data undergoes face swapping Zhou et al. (2022) to enhance ID consistency. Based on the first-stage training, we develop a local face- and head-swapping ID-inpainting model and, by leveraging an external aesthetic LoRA, produce high-quality data pairs in batches. The data filtering pipeline includes: (1) resolution filtering (768px), (2) aesthetic scoring  $> 5.2$  using aesthetic-predictor-v2-5, (3) face quality filtering (removing occluded, multi-person, or extreme-pose image via face detection and landmark models), and (4) similarity filtering with ArcFace threshold  $> 0.8$ . All images are annotated using Qwen2-VL-7B Wang et al. (2024a) for person and background descriptions.

324  
325  
326  
327  
328  
329  
330  
331  
332  
333  
334  
335  
336  
337  
338  
339  
340  
341  
342  
343  
344  
345  
346  
347  
348  
349  
350  
351  
352  
353  
354  
355  
356  
357  
358  
359  
360  
361  
362  
363  
364  
365  
366  
367  
368  
369  
370  
371  
372  
373  
374  
375  
376  
377



Figure 3: Qualitative comparison. A qualitative comparison of DMID with the latest baseline methods including PuLID-FLUX, InfiniteYou, and ACE++. To demonstrate that DMID can maintain good image-text consistency, we also use FLUX.1-dev to generate examples without an ID image. More example demonstrations can be found in the supplementary materials.

### 4.1.3 BASELINES AND EVALUATION

Since DMID is trained on Flux.1-dev with VAE-encoded image conditions, we compare against three state-of-the-art DiT-based methods: PuLID-FLUX, InfiniteYou, and ACE++ (Portrait LoRA). For evaluation, existing works use private small-scale ID datasets (e.g., 170 IDs in PuLID, 15 IDs in InfiniteYou). To enable fair comparison, we created a new ID consistency benchmark consisting of 504 unique IDs covering three age groups (elderly, young, children) and four races (Caucasian, African, Asian, South Asian) from LAION-human Ju et al. (2023); Schuhmann et al. (2022). Each test case includes a pair of images with different IDs, and all images are annotated using Qwen2-VL-7B Wang et al. (2024a).

## 4.2 COMPARISON WITH PRIOR METHODS

### 4.2.1 QUALITATIVE COMPARISON

As shown in Figure 3, our qualitative comparison reveals distinct performance characteristics. Methods using facial encoders (PuLID and InfiniteYou) maintain reasonable ID consistency but fail to capture fine-grained facial features like wrinkles and spots. While preserving basic semantic structures, their outputs exhibit noticeable artificial generation artifacts. In contrast, VAE-encoded approaches (ACE++ and our DMID) produce more photorealistic results with accurate reproduction of subtle facial details, as demonstrated by the faithful rendering of elderly subjects’ spots and wrinkles in the fourth example column. Crucially, DMID achieves superior image-text alignment, maintaining not only semantic structures but also precise pose preservation and textual element fidelity across all examples. These visual comparisons confirm DMID’s state-of-the-art performance in terms of identity consistency and textual semantics.

Table 1: Quantitative comparison results

Method	CosSim $\uparrow$	CLIPScore $\uparrow$
Ace++	0.5095	30.4532
InfiniteYou	0.6724	31.4062
PuLID	0.7051	32.0730
<b>DMID<math>_{\beta=0.5}</math></b>	<b>0.7672</b>	<b>32.4694</b>

### 4.2.2 QUANTITATIVE COMPARISON

Quantitative results are presented in Table 1. We employ CLIPScore to measure image-text semantic alignment, which evaluates the model’s ability to follow input prompts. Cosine similarity (CosSim) is used to assess facial identity consistency. The results demonstrate that DMID achieves state-of-the-art performance in both identity consistency and textual semantics.

### 4.3 ABLATION STUDIES

#### 4.3.1 IMPACT OF STRENGTH FACTOR $\alpha$

We conduct ablation experiments to validate the effectiveness of different components in our proposed method and the multi-stage training strategy. First, we present a preliminary experiment to illustrate the motivation behind our approach. We set the strength factor  $\alpha$  using the methods described in Equations 2 and 3, and evaluate performance on our proposed ID benchmark.



Figure 4: Demonstration of strength factor  $\alpha$  testing.

these visual observations. These findings motivate our dynamic mask learning scheme.

#### 4.3.2 ABLATION STUDY ON DYNAMIC ATTENTION MASK

To validate the effectiveness of the proposed dynamic masking scheme, we conduct an ablation study on Stage 2 components. This stage incorporates two loss functions ( $\mathcal{L}_{\text{consis}}$  and  $\mathcal{L}_{\text{im}}$ ), whose individual contributions are examined in Figure 5 and Exp 2 of Table 2. Figure 5 visualizes generated samples with corresponding identity-focused attention heatmaps, revealing three key observations:

**DMID-2 w/o  $\mathcal{L}_{\text{im}}$ :** As shown in Table 2, the model fails to effectively adjust its attention distribution, resulting in negligible visual changes compared to the baseline (DMID-1).  
**DMID-2 w/o  $\mathcal{L}_{\text{consis}}$ :** Excessive attention on facial regions induces copy-paste artifacts and mouth distortions. Quantitative metrics align with visual observations: facial similarity increases while CLIP scores significantly drop.

**DMID-2:** When  $\mathcal{L}_{\text{consis}}$  and  $\mathcal{L}_{\text{im}}$  are optimized together, the model effectively focuses on key facial regions, improving facial similarity while reducing the loss of semantic information. Notably, DMID-2 alone already achieves performance on par with the current state-of-the-art PuLID-Flux.

#### 4.3.3 MULTI-STAGE TRAINING STRATEGY

Table 2: Performance Comparison of Models

Experiment	Model	CosSim $\uparrow$	CLIPScore $\uparrow$
Exp 1	DMID-1	0.4906	32.6607
	DMID-1 $\alpha=1.2$	0.6288	31.8678
	DMID-1 $\alpha=1.5$	0.7455	30.3416
Exp 2	DMID-2 w/o $\mathcal{L}_{\text{im}}$	0.4890	32.7084
	DMID-2 w/o $\mathcal{L}_{\text{consis}}$	0.7595	31.5468
	DMID-2	0.7021	32.2621
Exp 3	DMID	0.7803	32.1851
	DMID $\beta=0.75$	0.7712	32.3485
	DMID $\beta=0.6$	0.7722	32.4371
	DMID $\beta=0.5$	0.7672	32.4694
	DMID $\beta=0.4$	0.7611	32.5147

Figure 4 shows selected test cases. It can be observed that DMID-1 preserves the original semantic structure, including complete text information. As the strength factor  $\alpha$  increases, the similarity between the generated face and the ID image face improves significantly; however, the original background, clothing, and text undergo noticeable changes. The quantitative results in Exp 1 of Table 2 are consistent with



Figure 5: Ablation study on dynamic attention mask learning.

432 Table 2 and Figure 6 demonstrate the effective-  
 433 ness of our multi-stage training strategy. As  
 434 shown in Table 2, facial similarity metrics sig-  
 435 nificantly improve with minimal semantic in-  
 436 formation loss. This phenomenon is also re-  
 437 flected in visual results: Figure 6 shows that as  
 438 training progresses, the model captures increas-  
 439 ingly fine-grained facial features such as wrin-  
 440 kles and spots. Therefore, we conclude that  
 441 the proposed multi-stage training strategy ef-  
 442 fectively enhances face similarity while reduc-  
 443 ing semantic loss.

444  
 445 4.3.4 ATTNMASK STRENGTH FACTOR  $\beta$

446 During the experiments, we observed a key  
 447 phenomenon: there exists a saturation thresh-  
 448 old for the model’s effective attention to iden-  
 449 tity features. Beyond this threshold, increasing  
 450 the attention weights not only fails to alter the  
 451 visual appearance of facial regions but also ex-  
 452 acerbates the loss of semantic information. Ex-  
 453 tensive test results demonstrate that the facial  
 454 similarity (CosSim) of most samples reaches a  
 455 saturation threshold between 0.7 and 0.8; fur-  
 456 ther improving this metric only yields redund-  
 457 ant gains, while the semantic consistency met-  
 458 ric (CLIPScore) continues to decline. There-  
 459 fore, we introduce the design of an attention  
 460 mask scaling channel.

461 Based on this phenomenon, we achieve a re-  
 462 fined trade-off between the two consistencies  
 463 by adjusting the scaling factor  $\beta$ . The core  
 464 of this method lies in strategically sacrificing  
 465 redundant identity attention to specifically en-  
 466 hance semantic control capability. As illus-  
 467 trated in Figure 7, regulating  $\beta$  enables: 1)  
 468 Fine-grained semantic control (decorative at-  
 469 tributes such as glasses, hats, and expressions); 2)  
 470 Complete preservation of identity details (includ-  
 471 ing biological features like tiny spots); 3) Avoidance of semantic distortion caused by over-attention.

472 The quantitative results in Exp 3 of Table 2 confirm that this method can further improve semantic  
 473 alignment by sacrificing redundant similarity while maintaining the integrity of identity details. This  
 474 indicates that the  $\beta$ -based attention scaling channel effectively resolves the consistency conflict.

475 5 CONCLUSION

476 This paper presents DMID, a high ID-consistency approach for limited-data scenarios. By integrat-  
 477 ing dynamic attention masking with a joint loss enforcing both identity and distribution consistency,  
 478 and employing a three-stage progressive training strategy, DMID significantly improves identity  
 479 fidelity and detail reconstruction while mitigating the conflict between textual semantics and con-  
 480 dition image. The AttnMask strength factor design allows more refined editing. On the data side,  
 481 we develop an automated, high-quality pipeline that constructs training pairs with preserved identity  
 482 consistency under small-scale data and release the ID Consistency Benchmark, a dataset covering  
 483 504 distinct identities. Extensive experiments show DMID outperforms existing methods in both  
 484 identity consistency and facial detail recovery. Future work can adopt DMID’s training paradigm  
 485 for identity consistency in specific scenarios with minimal data and reduced cost.



Figure 6: Ablation study on multi-stage training strategy.



Figure 7: Strength factor  $\beta$  ablation.

## 6 REPRODUCIBILITY STATEMENT

The specific implementation code of the proposed method will be open-sourced, and the anonymous link will be provided in the comments. The data construction pipeline can be referred to in the appendix section.

## REFERENCES

- James Betker, Gabriel Goh, Li Jing, Tim Brooks, Jianfeng Wang, Linjie Li, Long Ouyang, Juntang Zhuang, Joyce Lee, Yufei Guo, et al. Improving image generation with better captions. *Computer Science*. <https://cdn.openai.com/papers/dall-e-3.pdf>, 2(3):8, 2023.
- Di Chang, Yichun Shi, Quankai Gao, Jessica Fu, Hongyi Xu, Guoxian Song, Qing Yan, Yizhe Zhu, Xiao Yang, and Mohammad Soleymani. Magicpose: Realistic human poses and facial expressions retargeting with identity-aware diffusion. *arXiv preprint arXiv:2311.12052*, 2023.
- Junsong Chen, Jincheng Yu, Chongjian Ge, Lewei Yao, Enze Xie, Yue Wu, Zhongdao Wang, James Kwok, Ping Luo, Huchuan Lu, and Zhenguo Li. Pixart- $\alpha$ : Fast training of diffusion transformer for photorealistic text-to-image synthesis, 2023. URL <https://arxiv.org/abs/2310.00426>.
- Yisol Choi, Sangkyung Kwak, Kyungmin Lee, Hyungwon Choi, and Jinwoo Shin. Improving diffusion models for authentic virtual try-on in the wild. In *European Conference on Computer Vision*, pp. 206–235. Springer, 2024.
- Kevin Clark, Paul Vicol, Kevin Swersky, and David J Fleet. Directly fine-tuning diffusion models on differentiable rewards. *arXiv preprint arXiv:2309.17400*, 2023.
- Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos Zafeiriou. 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.
- Shilpa Devalal and A Karthikeyan. Lora technology-an overview. In *2018 second international conference on electronics, communication and aerospace technology (ICECA)*, pp. 284–290. IEEE, 2018.
- Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis. In *Forty-first international conference on machine learning*, 2024.
- Yuming Gu, Hongyi Xu, You Xie, Guoxian Song, Yichun Shi, Di Chang, Jing Yang, and Linjie Luo. Diffportrait3d: Controllable diffusion for zero-shot portrait view synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10456–10465, 2024.
- Zinan Guo, Yanze Wu, Chen Zhuowei, Peng Zhang, Qian He, et al. Pulid: Pure and lightning id customization via contrastive alignment. *Advances in neural information processing systems*, 37: 36777–36804, 2024.
- Yue Han, Junwei Zhu, Keke He, Xu Chen, Yanhao Ge, Wei Li, Xiangtai Li, Jiangning Zhang, Chengjie Wang, and Yong Liu. Face-adapter for pre-trained diffusion models with fine-grained id and attribute control. In *European Conference on Computer Vision*, pp. 20–36. Springer, 2024.
- Li Hu. Animate anyone: Consistent and controllable image-to-video synthesis for character animation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8153–8163, 2024.
- Zehuan Huang, Hongxing Fan, Lipeng Wang, and Lu Sheng. From parts to whole: A unified reference framework for controllable human image generation. *arXiv preprint arXiv:2404.15267*, 2024.
- Liming Jiang, Qing Yan, Yumin Jia, Zichuan Liu, Hao Kang, and Xin Lu. Infinitelyyou: Flexible photo recrafting while preserving your identity. *arXiv preprint arXiv:2503.16418*, 2025.

- 540 Xuan Ju, Ailing Zeng, Chenchen Zhao, Jianan Wang, Lei Zhang, and Qiang Xu. Humansd: A native  
541 skeleton-guided diffusion model for human image generation, 2023. URL <https://arxiv.org/abs/2304.04269>.  
542  
543
- 544 Diederik P Kingma, Max Welling, et al. Auto-encoding variational bayes, 2013.  
545  
546 Black Forest Labs. Flux. <https://github.com/black-forest-labs/flux>, 2024.
- 547 Junzhe Li, Yutao Cui, Tao Huang, Yinping Ma, Chun Fan, Miles Yang, and Zhao Zhong. Mix-  
548 grpo: Unlocking flow-based grpo efficiency with mixed ode-sde, 2025. URL <https://arxiv.org/abs/2507.21802>.  
549
- 550 Zhen Li, Mingdeng Cao, Xintao Wang, Zhongang Qi, Ming-Ming Cheng, and Ying Shan. Pho-  
551 tomaker: Customizing realistic human photos via stacked id embedding. In *Proceedings of the*  
552 *IEEE/CVF conference on computer vision and pattern recognition*, pp. 8640–8650, 2024.  
553
- 554 Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching  
555 for generative modeling. *arXiv preprint arXiv:2210.02747*, 2022.
- 556 Jie Liu, Gongye Liu, Jiajun Liang, Yangguang Li, Jiaheng Liu, Xintao Wang, Pengfei Wan,  
557 Di Zhang, and Wanli Ouyang. Flow-grpo: Training flow matching models via online rl. *arXiv*  
558 *preprint arXiv:2505.05470*, 2025.
- 559 Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate and  
560 transfer data with rectified flow. *arXiv preprint arXiv:2209.03003*, 2022.  
561
- 562 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint*  
563 *arXiv:1711.05101*, 2017.
- 564 Chaojie Mao, Jingfeng Zhang, Yulin Pan, Zeyinzi Jiang, Zhen Han, Yu Liu, and Jingren Zhou.  
565 Ace++: Instruction-based image creation and editing via context-aware content filling. *arXiv*  
566 *preprint arXiv:2501.02487*, 2025.
- 567 Konstantin Mishchenko and Aaron Defazio. Prodigy: An expeditiously adaptive parameter-free  
568 learner, 2024. URL <https://arxiv.org/abs/2306.06101>.  
569
- 570 William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of*  
571 *the IEEE/CVF international conference on computer vision*, pp. 4195–4205, 2023.  
572
- 573 Xu Peng, Junwei Zhu, Boyuan Jiang, Ying Tai, Donghao Luo, Jiangning Zhang, Wei Lin, Taisong  
574 Jin, Chengjie Wang, and Rongrong Ji. Portraitbooth: A versatile portrait model for fast identity-  
575 preserved personalization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*  
576 *Pattern Recognition*, pp. 27080–27090, 2024.
- 577 Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe  
578 Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image  
579 synthesis. *arXiv preprint arXiv:2307.01952*, 2023.
- 580 Mihir Prabhudesai, Anirudh Goyal, Deepak Pathak, and Katerina Fragkiadaki. Aligning text-to-  
581 image diffusion models with reward backpropagation. 2023.  
582
- 583 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agar-  
584 wal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya  
585 Sutskever. Learning transferable visual models from natural language supervision, 2021. URL  
586 <https://arxiv.org/abs/2103.00020>.
- 587 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-  
588 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF confer-*  
589 *ence on computer vision and pattern recognition*, pp. 10684–10695, 2022.  
590
- 591 Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman.  
592 Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *Pro-*  
593 *ceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 22500–  
22510, 2023.

- 594 Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi  
595 Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, Patrick Schramowski,  
596 Srivatsa Kundurthy, Katherine Crowson, Ludwig Schmidt, Robert Kaczmarczyk, and Jenia Jitsev.  
597 Laion-5b: An open large-scale dataset for training next generation image-text models, 2022. URL  
598 <https://arxiv.org/abs/2210.08402>.
- 599 Zhenxiong Tan, Songhua Liu, Xingyi Yang, Qiaochu Xue, and Xinchao Wang. Ominicontrol: Min-  
600 imal and universal control for diffusion transformer. *arXiv preprint arXiv:2411.15098*, 2024.
- 601  
602 Linrui Tian, Qi Wang, Bang Zhang, and Liefeng Bo. Emo: Emote portrait alive generating ex-  
603 pressive portrait videos with audio2video diffusion model under weak conditions. In *European*  
604 *Conference on Computer Vision*, pp. 244–260. Springer, 2024.
- 605  
606 Bram Wallace, Meihua Dang, Rafael Rafailov, Linqi Zhou, Aaron Lou, Senthil Purushwalkam,  
607 Stefano Ermon, Caiming Xiong, Shafiq Joty, and Nikhil Naik. Diffusion model alignment using  
608 direct preference optimization. In *Proceedings of the IEEE/CVF Conference on Computer Vision*  
609 *and Pattern Recognition*, pp. 8228–8238, 2024.
- 610  
611 Team Wan, Ang Wang, Baole Ai, Bin Wen, Chaojie Mao, Chen-Wei Xie, Di Chen, Feiwu Yu,  
612 Haiming Zhao, Jianxiao Yang, et al. Wan: Open and advanced large-scale video generative  
613 models. *arXiv preprint arXiv:2503.20314*, 2025.
- 614  
615 Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu,  
616 Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng  
617 Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. Qwen2-vl: Enhancing vision-language model’s  
618 perception of the world at any resolution, 2024a. URL <https://arxiv.org/abs/2409.12191>.
- 619  
620 Qixun Wang, Xu Bai, Haofan Wang, Zekui Qin, Anthony Chen, Huaxia Li, Xu Tang, and Yao Hu.  
621 Instantid: Zero-shot identity-preserving generation in seconds. *arXiv preprint arXiv:2401.07519*,  
622 2024b.
- 623  
624 Rui Wang, Hailong Guo, Jiaming Liu, Huaxia Li, Haibo Zhao, Xu Tang, Yao Hu, Hao Tang,  
625 and Peipei Li. Stablegarment: Garment-centric generation via stable diffusion. *arXiv preprint*  
626 *arXiv:2403.10783*, 2024c.
- 627  
628 Zirui Wang, Zhizhou Sha, Zheng Ding, Yilin Wang, and Zhuowen Tu. Tokencompose: Text-to-  
629 image diffusion with token-level supervision. In *Proceedings of the IEEE/CVF Conference on*  
630 *Computer Vision and Pattern Recognition*, pp. 8553–8564, 2024d.
- 631  
632 Guangxuan Xiao, Tianwei Yin, William T Freeman, Frédo Durand, and Song Han. Fastcomposer:  
633 Tuning-free multi-subject image generation with localized attention. *International Journal of*  
634 *Computer Vision*, 133(3):1175–1194, 2025.
- 635  
636 Mingwang Xu, Hui Li, Qingkun Su, Hanlin Shang, Liwei Zhang, Ce Liu, Jingdong Wang, Yao Yao,  
637 and Siyu Zhu. Hallo: Hierarchical audio-driven visual synthesis for portrait image animation.  
638 *arXiv preprint arXiv:2406.08801*, 2024a.
- 639  
640 Zhongcong Xu, Jianfeng Zhang, Jun Hao Liew, Hanshu Yan, Jia-Wei Liu, Chenxu Zhang, Jiashi  
641 Feng, and Mike Zheng Shou. Magicanimate: Temporally consistent human image animation  
642 using diffusion model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*  
643 *Pattern Recognition*, pp. 1481–1490, 2024b.
- 644  
645 Zeyue Xue, Jie Wu, Yu Gao, Fangyuan Kong, Lingting Zhu, Mengzhao Chen, Zhiheng Liu, Wei  
646 Liu, Qiushan Guo, Weilin Huang, et al. Dancegrp: Unleashing grp on visual generation. *arXiv*  
647 *preprint arXiv:2505.07818*, 2025.
- 648  
649 Hu Ye, Jun Zhang, Sibao Liu, Xiao Han, and Wei Yang. Ip-adapter: Text compatible image prompt  
650 adapter for text-to-image diffusion models. *arXiv preprint arXiv:2308.06721*, 2023.
- 651  
652 Ge Yuan, Xiaodong Cun, Yong Zhang, Maomao Li, Chenyang Qi, Xintao Wang, Ying Shan,  
653 and Huicheng Zheng. Inserting anybody in diffusion models via celeb basis. *arXiv preprint*  
654 *arXiv:2306.00926*, 2023.

648 Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image  
 649 diffusion models. In *Proceedings of the IEEE/CVF international conference on computer vision*,  
 650 pp. 3836–3847, 2023.

651 Shilong Zhang, Lianghua Huang, Xi Chen, Yifei Zhang, Zhi-Fan Wu, Yutong Feng, Wei Wang,  
 652 Yujun Shen, Yu Liu, and Ping Luo. Flashface: Human image personalization with high-fidelity  
 653 identity preservation. *arXiv preprint arXiv:2403.17008*, 2024.

654 Shangchen Zhou, Kelvin C. K. Chan, Chongyi Li, and Chen Change Loy. Towards robust blind  
 655 face restoration with codebook lookup transformer, 2022. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2206.11253)  
 656 [2206.11253](https://arxiv.org/abs/2206.11253).  
 657

## 658 A APPENDIX

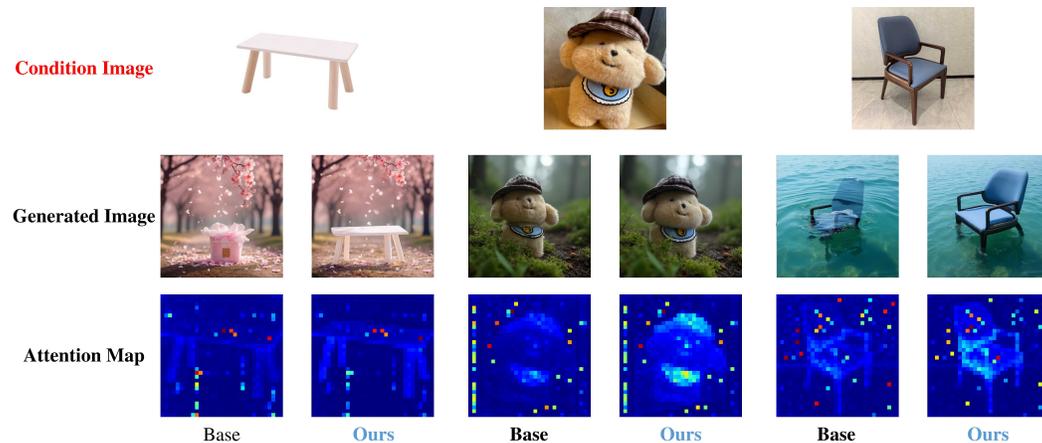
### 659 A.1 THE USE OF LARGE LANGUAGE MODELS

660 This paper only uses large language models for English translation and English grammar correction  
 661 and polishing.  
 662

### 663 A.2 EXTENDED EXPERIMENT

664 There is an inherent conflict between facial similarity and textual semantic alignment. Based on this  
 665 premise, we propose DMID. In the previous section, we demonstrated the strong identity-preserving  
 666 capability and editability of DMID in real-world scenarios. To further showcase the effectiveness of  
 667 our proposed method, we will provide additional experimental results.  
 668

### 669 A.3 THE USABILITY OF THE METHOD IN OTHER TASKS



691 Figure 8: The Usability of the Method in Other Tasks

692 To verify that the AttnMaskNet scheme is not restricted exclusively to identity consistency preser-  
 693 vation tasks, we adopt the open-source subject consistency preservation model *subject\_512* from  
 694 Ominicontrol as the base model—corresponding to the stage-1 model described in our methodol-  
 695 ogy. For the training process of AttnMaskNet, we utilize Subject200K, the official training dataset  
 696 associated with *subject\_512*. The experimental results are presented in the Figure 8, which includes  
 697 both the generated output images and self-attention maps of the base model (*subject\_512*) and our  
 698 AttnMaskNet-augmented model. As observed from these results, AttnMaskNet retains its effective-  
 699 ness when applied to subject consistency tasks.

#### 700 A.3.1 DMID-INPAINTING

701 We reproduce the first stage of DMID on our self-constructed dataset and further propose a region-  
 controllable, high-fidelity local identity-preserving model, DMID-Inpainting. DMID-Inpainting not

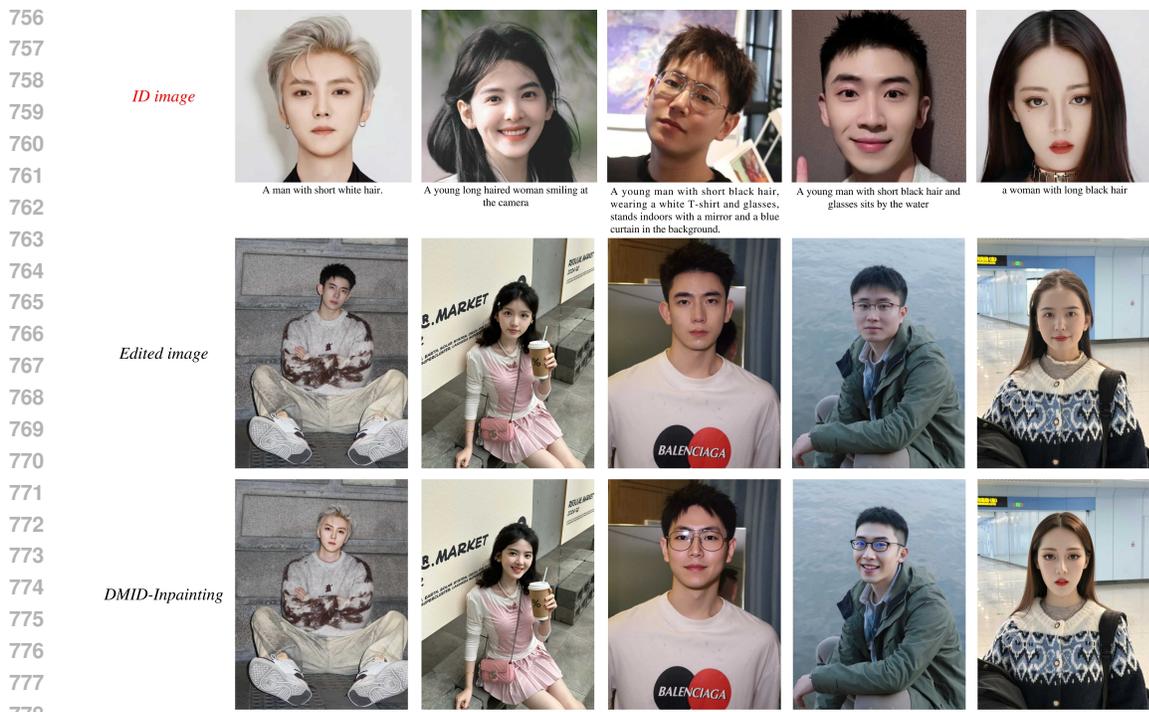


Figure 9: **Condition Image.** (a) Condition image of the normal inpainting model, where the black regions indicate the areas to be generated. (b) Condition image generated by hair and face partitioning, where the white regions represent the face, and the black regions correspond to the hair and background.

only verifies the effectiveness of ID preservation tasks in VAE encoding but also serves as a data construction pipeline. As shown in Figure 9, by explicitly introducing local masks, DMID-Inpainting divides the input space into facial regions and hairstyle/background regions, enabling the network to generate independently by region and optimize collaboratively. This significantly enhances the controllability and overall consistency of character generation. Additionally, DMID-Inpainting supports seamless integration with aesthetic LoRA and weight fusion with PuLID. It can flexibly construct diverse data pipelines by adjusting only a few hyperparameters. Experiments demonstrate that DMID-Inpainting exhibits strong generalization ability and ID consistency on the dataset constructed in this paper. As shown in Figure 10, DMID-Inpainting can still maintain a high degree of identity consistency in local editing tasks. Benefiting from the proposed region-wise control strategy, the disentangled representation of hairstyle and face significantly enhances the stability and controllability of each region. In practical applications, the accurate preservation of hairstyle further amplifies the benefits of identity preservation. It is worth noting that when the mask region only covers the face, DMID-Inpainting can serve as a stable face-swapping model. Additionally, using VAE as the encoder for image conditions brings consistency in details, making the edited regions more realistic. As shown in columns 3 and 5 of Figure 10, DMID-Inpainting can even accurately restore subtle facial features (such as facial moles), further verifying its advantages in detail consistency.

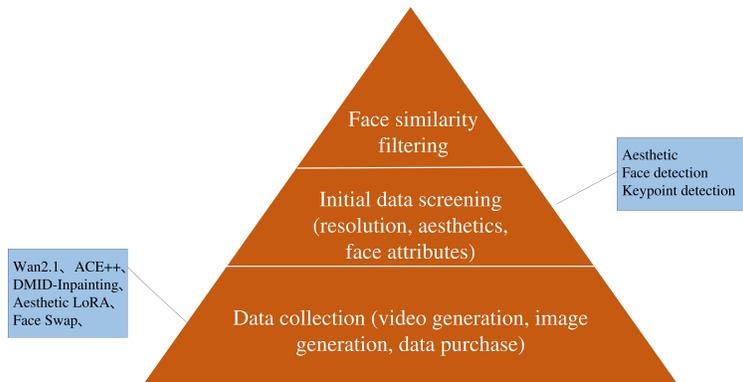
#### A.4 HIGH-QUALITY DATA CONSTRUCTION PIPELINE

Due to the use of VAE for conditional encoding, the model’s requirements for the quality of data pairs are significantly higher. For pairs with low similarity, the model fails to learn consistent ID representations. Unlike facial feature encoding, where the space is relatively small, the large output space of the VAE complicates training. On the other hand, for pairs with high similarity, the model is more prone to duplication issues, resulting in a lack of facial pose diversity. This underscores the critical role of the data construction process. Unlike existing approaches that rely on millions of samples, we do not adopt a large-scale data strategy. Instead, we base our approach on the following observation: noise images and condition images share the same VAE encoding path, which inherently provides high consistency in identity features. Additionally, the diversity of poses required for face generation tasks is relatively limited. Therefore, we argue that a small amount of carefully selected high-quality data can be used to train a competitive identity-preserving model, significantly reducing training costs. To achieve this, we have designed and fully implemented a systematic data construction pipeline, including steps from raw data collection, quality evaluation, similarity filtering, to final pairing, as illustrated in the figure 11.



779 **Figure 10: DMID-Inpainting.** This figure presents the results of DMID-Inpainting, from which it  
780 can be observed that this method can stably generate high-quality data.

781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809



800 **Figure 11: Overall Data Construction Pipeline.** The overall process of data construction comprises  
801 three stages: (1) data collection, (2) data filtering, and (3) face similarity filtering.

#### 802 A.4.1 DATA COLLECTION

803 We utilize Wan2.1 to batch-generate videos with high-quality portraits as the initial frames. Of-  
804 fline frameworks such as ACE++ are employed, integrating inpainting, face swapping, and aesthetic  
805 LoRA to output high-fidelity identity-consistent images. Additionally, DMID-Inpainting is used to  
806 generate high-quality facial images. Furthermore, high-resolution portraits are collected in compli-  
807 ance with relevant regulations to enhance diversity.

#### 808 A.4.2 DATA FILTERING

809 Initial filtering is conducted using four thresholds: resolution, clarity, facial attributes (including  
810 face angle, face quality, face size, etc.), and aesthetic score. Subsequently, intra-identity pairs are  
811 retained by applying a high threshold to the facial feature similarity matrix. The entire process

810 involves no manual intervention, ultimately obtaining 12,000 identities with approximately 40,000  
811 high-quality images. Experimental validation demonstrates that this dataset is sufficient to support  
812 excellent ID consistency.  
813

#### 814 A.4.3 CONSTRUCTION OF DPO TRAINING DATA

815 We collected 188 ID images covering diverse ethnicities, ages, and genders. Based on these data,  
816 we constructed 300 image pairs, where one image serves as the ID image and the other as the target  
817 image. We used Qwen2-VL-7B to generate descriptions of the target images as prompts, and input  
818 the ID images (as conditional images) and these prompts into DMID-2. According to Equation 10,  
819 we selected the samples with the highest scores and the lowest scores among the 8 generation results,  
820 which were used as our training data.  
821

822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863