# From Parts to Whole: A Unified Reference Framework for Controllable Human Image Generation

Zehuan Huang*    Hongxing Fan*    Lipeng Wang*    Lu Sheng†

Beihang University

{huangzehuan, fanhongxing, wanglipeng, lsheng}@buaa.edu.cn

https://huanngzh.github.io/Parts2Whole/

(a) Generated human images conditioned on reference images from **different humans**. Each pair includes 4 human parts input (1st column) and generated image (2nd column).

(b) Generated human images from **varying numbers** of reference images. Each pair includes 1 or 2 human parts input (1st row) and generated image (2nd row).

Figure 1. We propose **Parts2Whole**, which can generate realistic and high-quality human figures in various postures from referential human part images of any quantity and different origins. Our method maintains the high alignment with the corresponding conditional semantic regions, while ensuring diversity and harmony among the whole body.

## Abstract

*Recent advancements in controllable human image generation have led to zero-shot generation using structural signals (e.g., pose, depth) or facial appearance. Yet, generating human images conditioned on multiple parts of human appearance remains challenging. Addressing this, we introduce Parts2Whole, a novel framework designed for generating customized portraits from multiple reference images, including pose images and various aspects of human appearance. To achieve this, we first develop a semantic-aware appearance encoder to retain details of different human parts, which processes each image based on its textual label to a series of multi-scale feature maps rather than one image token, preserving the image dimension. Second, our framework supports multi-image conditioned generation through a shared self-attention mechanism that operates across reference and target features during the diffusion process. We enhance the vanilla attention mechanism by incorporating mask information from the reference human images, allowing for the precise selection of any part. Extensive experiments demonstrate the superiority of our approach over existing alternatives, offering advanced capabilities for multi-part controllable human image customization.*

---

*Equal Contribution.

†Corresponding author.

# 1. Introduction

Controllable human image generation aims to synthesize human images that align with specific textual descriptions, structural signals or more precise appearance conditions. It emerges as a significant technology within the realm of digital content creation, providing users with a portrait customization solution. However, due to the complexity of the control conditions, this task presents significant challenges, especially when it comes to multi-type condition input and control of various aspects of human appearance.

As diffusion models [7, 12, 30, 35, 36, 39] have brought great success in image generation, the task of controllable human image generation has experienced rapid development. Several works [18] utilize languages as condition, generating human images by providing attributes about the textures of clothes. Due to the rough control of texts, it struggles to accurately guide the generation of human appearance. Another group of works [24, 29, 56, 58] focuses on introducing structural signals to control human posture. Although these methods have achieved impressive results, they do not consider appearance as a condition, which is crucial for portrait customization.

Recently, several works [3, 9, 13, 21, 23, 26, 38, 40, 45, 52, 57] have emerged that use appearance conditions to guide human image generation. They learn human representation from reference images and generate images aligning with the specific face identity. One prominent approach involves test-time fine-tuning [9, 13, 21, 26, 38]. It requires substantial computational resources to learn each new individual, which costs about half an hour to achieve satisfactory results. Another approach [3, 23, 40, 45, 52, 57] investigates the zero-shot setting to bypass the fine-tuning cost. It encodes the reference image into one or several tokens and injects them into the generation process along with text tokens. These zero-shot methods make human image customization practical with faster speed. However, due to the loss of spatial representations when encoding the reference images into one or a few tokens, they struggle to preserve appearance details. And they lack the design to obtain specified information from the images, but instead utilize all the information, resulting in ambiguous subject representation.

In this paper, we target generating human images from multi-part images of human appearance, along with specific pose maps or optionally text descriptions. The above-mentioned generation methods conditioned on structural signals [29, 56] or face identity [9, 13, 38, 45, 52, 57] have their limitations on this task (results shown in Fig. 4 and Fig. 5). It is attributed to the spatial misalignment of the input multi-body parts with the target image, and the lack of specific design in existing methods to **address the variation in spatial positions during feature injection**. Methods like IP-Adapter [52] and SSR-Encoder [57] incorporate features into the denoising U-Net through cross-

attention mechanisms. They encode reference images into other modal features (e.g., semantic features) and utilize the cross-attention keys and values from them rather than from **image dimensional** feature maps. As a result, the spatial relationship in the original image dimensions between the reference images and the target image is lost, resulting in a mixture of attributes from different subjects. Although methods like ControlNet [56] encode the reference images into image-dimensional features, they add the features to the feature maps in the U-Net decoder. It is suitable for tasks where the condition maps and the target map have the same structure, such as guiding generation using line drawings. However, in the case of the spatial misalignment of the conditional images with the target image, it is difficult to model the correlation of spatial information by **directly adding or concat features** on the channel dimension.

To address the above issues, we present Parts2Whole, a unified reference framework for portrait customization from multiple reference images, including various parts of human appearance (e.g., hair, face, clothes, shoes, etc.) and pose maps. Inspired by the effective reference mechanism used in image-to-video tasks [14, 49], we develop a semantic-aware appearance encoder based on the Reference U-Net architecture. It encodes each image with its textual label into **a series of multi-scale feature maps in image dimension**, preserving appearance details and spatial information of multiple reference images. The additional semantic condition represents a category instruction, which helps retain richer shapes and detailed attributes of each aspect. Furthermore, to preserve the positional relationship when injecting reference features into the image generation process, we employ a **shared self-attention** operation across reference and target features during the diffusion process. We also build a tiny convolution network to extract the pose features and inject them into the generation. To precisely select the specified part from each reference image, we enhance the vanilla self-attention mechanism by **incorporating masks** of the subjects in the reference images.

Equipped with these techniques, Parts2Whole demonstrates superior quality and controllability for human image generation. Our contributions are summarized as follows:

- We construct a novel framework, Parts2Whole, which supports the controllable generation of human images conditioned on texts, pose signals, and multiple aspects of human appearance.
- We propose an advanced multi-reference mechanism consisting of a semantic-aware image encoder and the shared attention operation, which retains details of the specific key elements and achieves precise subject selection with the help of our proposed mask-guided approach.
- Experiments show that our Parts2Whole generates high-quality human images from multiple conditions and maintains high consistency with the given conditions.

## 2. Related Work

**Text-to-Image Generation.** In recent years, text-to-image generation has made remarkable progress, particularly with the development of diffusion models [7, 12, 16, 30, 32, 35, 36, 39] and auto-regressive models [2, 43, 53], which have propelled text-to-image generation to large-scale commercialization. Since DALLE2 [35], Stable Diffusion [36] and Imagen [39] employ diffusion models as generative models and train the models on large datasets, text-to-image synthesis ability has been significantly enhanced. More recently, Stable Diffusion XL [32], a two-stage cascade diffusion model, has greatly improved the generation of high-frequency details and overall image color, taking aesthetic appeal to a higher level. However, these existing methods are limited to generating images solely from text prompts, and they do not meet the demand for producing customized images with the preservation of appearance.

**Controllable Image Generation.** Given the robust generative capabilities of image diffusion models, a series of research [3, 13, 29, 33, 38, 52, 56, 57] attempts to explore the controllability of image generation, enabling image synthesis guided by multi-modal conditions. Some work [15, 19, 29, 33, 56, 58] focuses on introducing structural signals such as edges, depth maps, and segmentation maps, to control the spatial structure of generated images. Another group of work [3, 9, 13, 38, 52] uses appearance conditions to guide image generation, aiming to generate images aligning with specific concepts like identity and style, known as subject-driven image generation. The methods generally fall into two categories: those requiring test-time fine-tuning and those that do not. Test-time fine-tuning methods [9, 13, 21, 26, 38] often optimizes additional text embedding, parameter residuals or direct fine-tune the whole model to fit the specified subject. Although these methods have achieved impressive results, they cost about half an hour to achieve satisfactory results. Fine-tuning-free methods [3, 10, 27, 40, 47, 52, 57] typically train an additional encoding network to encode the reference image into embeddings or image prompts. However, due to the loss of spatial representations when encoding the reference images into one or a few tokens, they struggle to preserve appearance details.

**Controllable Human Image Generation.** In this paper, we mainly focus on controllable human image generation and aim to synthesize human images aligning with specific text prompts, pose signals, and various parts of human appearance. Text2Human [18] generates full-body human images using detailed descriptions about the textures of clothes, but is limited by the coarse-grained textual condition. Test-time fine-tuning methods [13, 21, 38] produce satisfactory results, but when it comes to customizing portraits using multiple parts of human appearance, they take much more

time to fit each aspect. Recently, methods like IP-Adapter-FaceID [52], FastComposer [48], PhotoMaker [23], and InstantID [45] show promising results on zero-shot human image personalization. They encode the reference face to one or several tokens as conditions to generate customized images. With the addition of adaptable structural control networks [29, 56], these methods can generate portraits aligned with specified poses and human identities. However, they usually fail to maintain the details of human identities and utilize all the information from a single image, resulting in ambiguous subject representation. These make it difficult to apply these schemes to precisely generation conditioned on multiple parts of the human appearance. In contrast, our Parts2Whole is both generalizable and efficient, and precisely retains details in multiple parts of human appearance.

## 3. Method

We target controllable human image generation guided by multiple reference images. Given $N$ images that capture distinct parts of human appearance $x^{1:N}$ and a pose map $p$, and optionally text inputs, our objective is to synthesize a human image $\hat{x}$ aligning with the specified appearances and posture. To achieve this goal, we propose Parts2Whole, a specialized framework designed to interpolate various reference images and generate high-quality portraits.

In general, Parts2Whole is built on text-to-image diffusion models [36]. In the following sections, we start with an overview of T2I diffusion models, and in particular, the self-attention mechanism in Sec. 3.1. We continue by presenting our unified reference framework in Sec. 3.2, which consists of a semantic-aware appearance encoder, a shared self-attention that queries referential features within the self–attention layers, and the enhanced mask-guided subject selection. These methods enable Parts2Whole to accurately obtain the specific subject information from multiple reference images while preserving appearance details.

### 3.1. Preliminaries

**Text-to-Image Diffusion Models.** Diffusion models [12] exhibit promising capabilities in image generation. In this study, we select the widely adopted Stable Diffusion [36] as our foundational model, which is also known as Latent Diffusion Models (LDM). The model operates the denoising process in the latent space of an autoencoder [20], namely $\mathcal{E}(\cdot)$ and $\mathcal{D}(\cdot)$. During the training phase, an input image $x_0$ is initially mapped to the latent space using a frozen encoder, yielding $z_0 = \mathcal{E}(x_0)$, then perturbed by a pre-defined Markov process:

$$q(z_t|z_{t-1}) = \mathcal{N}(z_t; \sqrt{1 - \beta_t}z_{t-1}, \beta_t I) \quad (1)$$

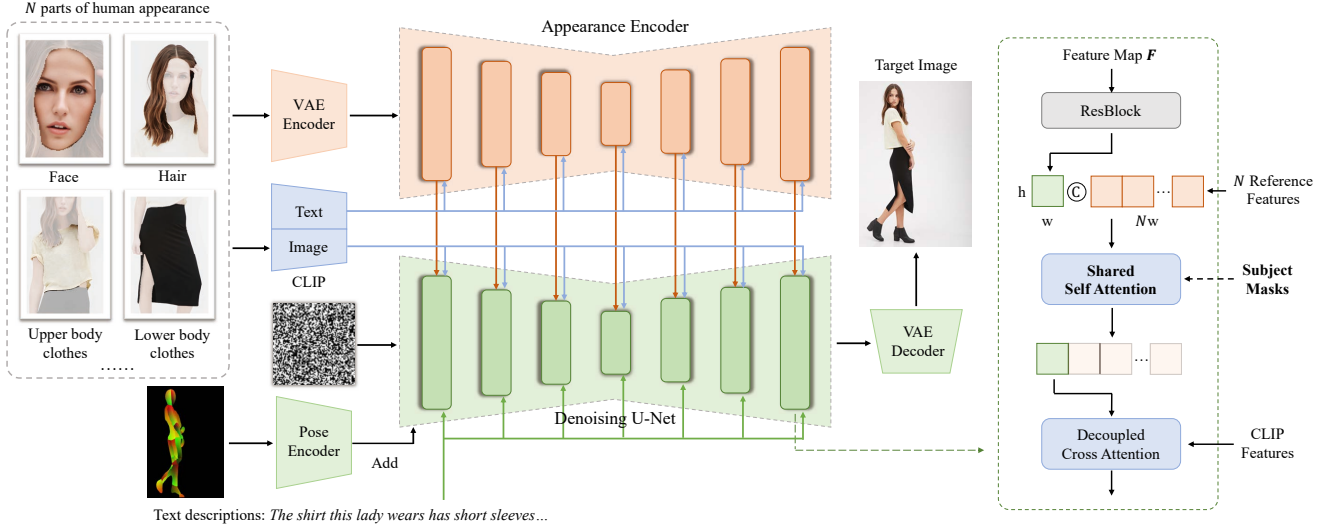For $t = 1, \cdots, T$, where $T$ represents the number of steps

Figure 2. Overview of Parts2Whole. Based on the text-to-image diffusion model, our method designs an appearance encoder for encoding various parts of human appearance into multi-scale feature maps. We build this encoder by copying the network structure and pre-trained weights from denoising U-Net. Features obtained from reference images with their textual labels are injected into the generation process by shared attention mechanism layer by layer. To precisely select the specified parts from reference images, we enhance the vanilla self-attention mechanism by incorporating subject masks in the reference images. An illustration of one block in U-Net is shown on the right part.

in the forward diffusion process. The sequence of hyperparameters $\beta_t$ determines the noise strength at each step. The denoising UNet $\epsilon_\theta$ is trained to approximate the reverse process $q(z_{t-1}|z_t)$. The training objective is expressed as:

$$\mathcal{L} = \mathbb{E}_{\mathcal{E}(x_0), \epsilon \sim \mathcal{N}(0,I), c, t}[\|\epsilon - \epsilon_\theta(z_t, c, t)\|_2^2] \quad (2)$$

Here, $c$ denotes the conditioning texts. At the inference stage, Stable Diffusion effectively reconstructs an image from Gaussian noise step by step, predicting the noise added at each stage. The denoised results are then fed into a latent decoder $\mathcal{D}(\cdot)$ to regenerate colored images from the latent representations, denoted as $\hat{x}_0 = \mathcal{D}(\hat{z}_0)$.

**Self-Attention in T2I Models.** Stable Diffusion [36] employs a U-Net architecture [37] that consists of convolution layers and transformer attention blocks [44]. In these attention mechanisms, self-attention layers are used to aggregate the spatial features of the image itself and cross-attention layers are designed to query information from text embedding. The main difference is that the cross-attention layer uses text features as keys and values, while in self-attention layers, image features with spatial dimensions serve as query, key, and value by themselves, preserving more freedom to represent spatially varying visual elements. The self-attention layer takes a feature map $F$ of the image as input and computes the attention of the feature in location $s$ with the entire feature map:

$$\tilde{F}_s = \text{SoftMax}\left(\frac{Q(F_s) \cdot K(F)^\top}{\sqrt{d}}\right) \cdot V(F) \quad (3)$$

where $Q, K, V$ are linear projection layers, $F \in \mathbb{R}^{(hw) \times d}$ is a flattened feature map obtained from the denoiser $\epsilon_\theta$, where $d$ is the feature dimension, and $h, w$ are intermediate spatial dimensions. $F_s, \tilde{F}_s$ is the input and output feature for location $s$ respectively.

Several works extend the self-attention layer to inject the reference image features [14, 49], or generate style-aligned or subject-consistent images [11, 17, 42], and demonstrate the effectiveness of this mechanism. Inspired by them, we extend the keys and values of the self-attention layer to multiple reference images and preserve the details of the referential appearance successfully.

## 3.2. Unified Reference Framework

As demonstrated in Fig. 2, our Parts2Whole consists of two branches: the reference branch used to encode multiple parts of human appearance, and the denoising branch, to gradually denoise the randomly sampled noise to finally obtain the image. The two branches utilize the same network architecture U-Net, initialized with the pretrained weights of Stable Diffusion [36]. In detail, our framework mainly consists of three crucial components: 1) Semantic-Aware Appearance Encoder, encoding the multi-scale features of various human parts from reference images; 2) Shared Self-Attention, which obtains detailed information and spatial information by sharing keys and values in self-attention layers between denoising U-Net and appearance encoder, and supports pose control by utiliz-
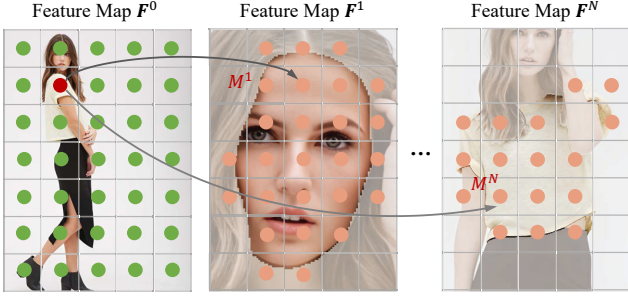
Feature Map $\boldsymbol{F}^0$     Feature Map $\boldsymbol{F}^1$     Feature Map $\boldsymbol{F}^N$

Figure 3. **Illustration of our Mask-Guided Attention.** For each patch $s$ (red point) on the feature map $\boldsymbol{F}^0$, given subject masks $M^{1:N}$ on the $N$ reference images, we only attend patch $s$ to features in these masks along with the patches on itself.

ing a lightweight pose encoder; 3) Enhanced Mask-Guided Subject Selection, achieving precisely subject selection by explicitly introducing subject masks into the self-attention mechanism.

**Semantic-Aware Appearance Encoder.** In image conditioned generation tasks, previous work [50, 52, 57] employs CLIP image encoder [34], or combined with some simple linear layers to encode reference images, thereby replacing the original text encoder in Stable Diffusion [36]. However, such methods struggle to preserve appearance details due to the loss of spatial representations when encoding the reference images into semantic-level features.

Inspired by recent works [1, 14, 49, 55] on dense reference image conditioning, we propose a semantic-aware appearance encoder with improved identity and details preservation. Specifically, we adopt a framework identical to the denoising U-Net for the appearance encoder. Unlike the denoising branch, we do not add any noise to the reference images. Given $N$ images capturing various parts of human appearance, we first compress them into latent features and then input them into the copied trainable U-Net. Instead of simply piecing multiple reference images together, we pass the latent features of different parts through the appearance encoder one by one and provide a textual class label for each part. These text labels, such as face, hair, upper body clothes, etc., are converted into feature representations by CLIP text encoder [34] and then injected into the appearance encoder through cross-attention. This simple yet effective external condition provides a classifier-like guidance, which enables the encoder to have semantic awareness of different parts of the human appearance rather than simply performing operations such as image downsampling and upsampling. This helps produce results that are not only rich in detail, but also flexible and realistic.

In the encoding process, we set the timestep to 0 and only perform one processing instead of iterating successively, so it will not cause time burden at the inference stage. We

cache the features before each self-attention layer for the next multi-image conditioned generation.

**Shared Self-Attention.** After obtaining the multi-layer feature maps of $N$ reference images, we do not directly add them to the features in denoising U-Net, but use shared keys and values in self-attention to achieve feature injection. This is because our reference and target images are not structurally aligned.

Take one certain self-attention layer as an example. Given the features of $N$ reference images $\boldsymbol{F}^{1:N}$ and the feature maps $\boldsymbol{F}^0$ in the denoising U-Net, we concatenate the feature maps of them side-by-side as input to the self-attention layer, denoted as $[\boldsymbol{F}^0|\boldsymbol{F}^1|\cdots|\boldsymbol{F}^N]$. This allows each location $s$ on $\boldsymbol{F}^0$ to attend to all locations on itself and reference feature maps, calculated as:

$$
\begin{aligned}
\tilde{\boldsymbol{F}}_s^0 = \text{SoftMax} & \left( \frac{Q(\boldsymbol{F}_s^0) \cdot K([\boldsymbol{F}^0|\boldsymbol{F}^1|\cdots|\boldsymbol{F}^N])^\top}{\sqrt{d}} \right) \\
& \cdot V([\boldsymbol{F}^0|\boldsymbol{F}^1|\cdots|\boldsymbol{F}^N])
\end{aligned}
\tag{4}
$$

We retain the cross-attention layers in Stable Diffusion [36] for injecting CLIP features of reference images and optional text input. We use the decoupled cross-attention proposed by IP-Adapter [52] to support both images and text input. Specifically, feature maps $\tilde{\boldsymbol{F}}^0$ obtained from shared self-attention serve as the origin of the query, and the reference image features and text features are each used as the key and value of the two cross-attention. The final feature maps are the sum of the two cross-attention outputs.

To further enhance the controllability of the human image generation, we add the pose map as an additional control. We construct a tiny convolution network, which is similar to the condition embedding network in ControlNet [56], to extract the features of the pose map. The features are then added to the initial feature maps in the denoising U-Net.

**Enhanced Mask-Guided Subject Selection.** We find that the vanilla shared self-attention leads to interference from irrelevant subjects in the reference images (shown in the 6th column in Fig. 5), resulting in an unnatural appearance and background. To synthesize human images conditioned on specified parts from each reference image, we enhance the vanilla self-attention mechanism by incorporating subject masks in the reference images. Fig. 3 presents this mechanism. Starting with a patch $s$ on a feature map $\boldsymbol{F}^0$ in the denoising U-Net, and subject masks $M^{1:N}$ on the $N$ reference images. When computing the attention map between the one in the denoising U-Net and those from the appearance encoder, patches that do not lie in these masks are ignored. Hence, the target patch $s$ only has access to features of the subjects specified by masks in the reference images, thereby avoiding interference from other elements such as the background. The final formulation of the mask-guided

attention is defined as follows:

$$\tilde{\boldsymbol{F}}_s^0 = \text{SoftMax} \left( \frac{Q(\boldsymbol{F}_s^0) \cdot K([\boldsymbol{F}^0|\cdots|\boldsymbol{F}_{M^N}^N])^\top}{\sqrt{d}} \right)$$
$$\cdot V([\boldsymbol{F}^0|\cdots|\boldsymbol{F}_{M^N}^N]) \quad (5)$$

In practice, to avoid misalignment between masks and original images caused by downsampling, a full-ones convolutional kernel is applied to the mask before each attention layer, ensuring that the mask preserves critical regions. Overall, the mask-guided attention enhances the ability of Parts2Whole to precisely extract the appearance of specified subjects in reference images.

# 4. Experiments

## 4.1. Implementation Details

**Dataset.** To train the Parts2Whole model, we build a multi-modal dataset comprising about 41,500 reference-target pairs from the open-source DeepFashion-MultiModal dataset [18, 25]. Each pair in this newly constructed dataset includes multiple reference images, which encompass human pose images (e.g., OpenPose, Human Parsing, Dense-Pose), various aspects of human appearance (e.g., hair, face, clothes, shoes) with their short textual labels, and a target image featuring the same individual (ID) in the same outfit but in a different pose, along with textual captions.

The DeepFashion-MultiModal dataset exhibits noise in its ID data. For example, different images are tagged with the same ID but depict different individuals. To address this issue, we first cleanse the IDs by extracting facial ID features from images tagged with the same ID using InsightFace[5, 6]. Cosine similarity is then used to evaluate the similarity between image ID feature pairs to distinguish between different ID images within the same ID group. Subsequently, we utilize DWPose[51] to generate pose images corresponding to each image. Guided by human parsing files, we crop human images into various parts. Due to the low resolution of the cropped parts, we apply Real-ESRGAN[46] to enhance the image resolution, thus obtaining clearer reference images. Textual descriptions of the original dataset are used as captions. For constructing pairs, we select images with cleaned IDs that feature the same clothes and individual but in different poses. Specifically, a pair contains multiple parts from one human image as reference images, and an image of the person in another pose as the target. Finally, we build a total of about 41,500 pairs, of which the training set is about 40,000 and the test set is about 1,500 pairs.

**Detailed Configurations.** In this work, the denoising U-Net and the appearance encoder both leverage the pre-trained weights from Stable Diffusion-1.5 [36]. We use CLIP Vision Model with projection layers as our image encoder, initialized with Stable Diffusion Image Variations [22]. During training, we set the initial learning rate 1e-5 with a batch size of 64. The model is trained using 8 A800 GPUs, for a total of 30000 iterations. To maintain the capability of image generation, we randomly drop all of the reference image features and the pose condition with a probability of 0.2. At the same time, to improve the flexibility of generation, we randomly drop each appearance condition with a probability of 0.2, so that the human images can be generating from indefinite reference images. At the inference stage, we adopt DDIM sampler [41] with 50 steps, and set the guidance scale to 7.5.

## 4.2. Comparison with Existing Alternatives

Our Parts2Whole targets at controllable human image generation conditioned on multiple parts of human appearance. To evaluate the performance of our proposed framework, we compare our Parts2Whole with existing subject-driven solutions. For fairness, we make some improvements to the methods, to make them more suitable for generating human images from multiple conditions.

**Test-time Fine-tuning Methods.** Among the tuning-based methods, we adopt DreamBooth LoRA [13, 38] and Custom Diffusion [21] as baseline methods for comparison, as these methods are relatively robust and effective. DreamBooth LoRA inserts a smaller number of new weights into Stable Diffusion [36] and only trains these parameters on just a few images of a subject or style, thereby associating a special word in the prompt with the example images. Custom Diffusion fine-tunes only key and value projection matrices in the cross-attention layers to customize text-to-image models. Given as input several aspects of human appearance, we use these two methods to fine-tune Stable Diffusion, such that it learns to bind identifiers with specific human parts. As shown in Fig. 4, when it comes to multi-aspect composition, the attributes of different parts in images generated by these tuning-base methods mix together, resulting in unrealistic human images. In contrast, Parts2Whole generates high-fidelity results without the need for parameter tuning.

**Reference-based Methods.** Among the tuning-free methods, we adopt IP-Adapter [52] and SSR-Encoder [57] for comparison. IP-Adapter is an image prompt adapter that can be plugged into diffusion models to enable image prompting, and can be combined with other adapters like ControlNet [56]. We firstly use IP-Adapter FaceID and ControlNet to generate human images from facial appearance and pose maps. Then we repaint hair, clothes, shoes and other areas using the specific image step by step, thereby achieving multi-image conditioned generation of portraits in a multi-step way. SSR-Encoder is an effective encoder designed for selectively capturing any subject from single or multiple reference images by the text query or mask query. For fairness, we fine-tune it in our human

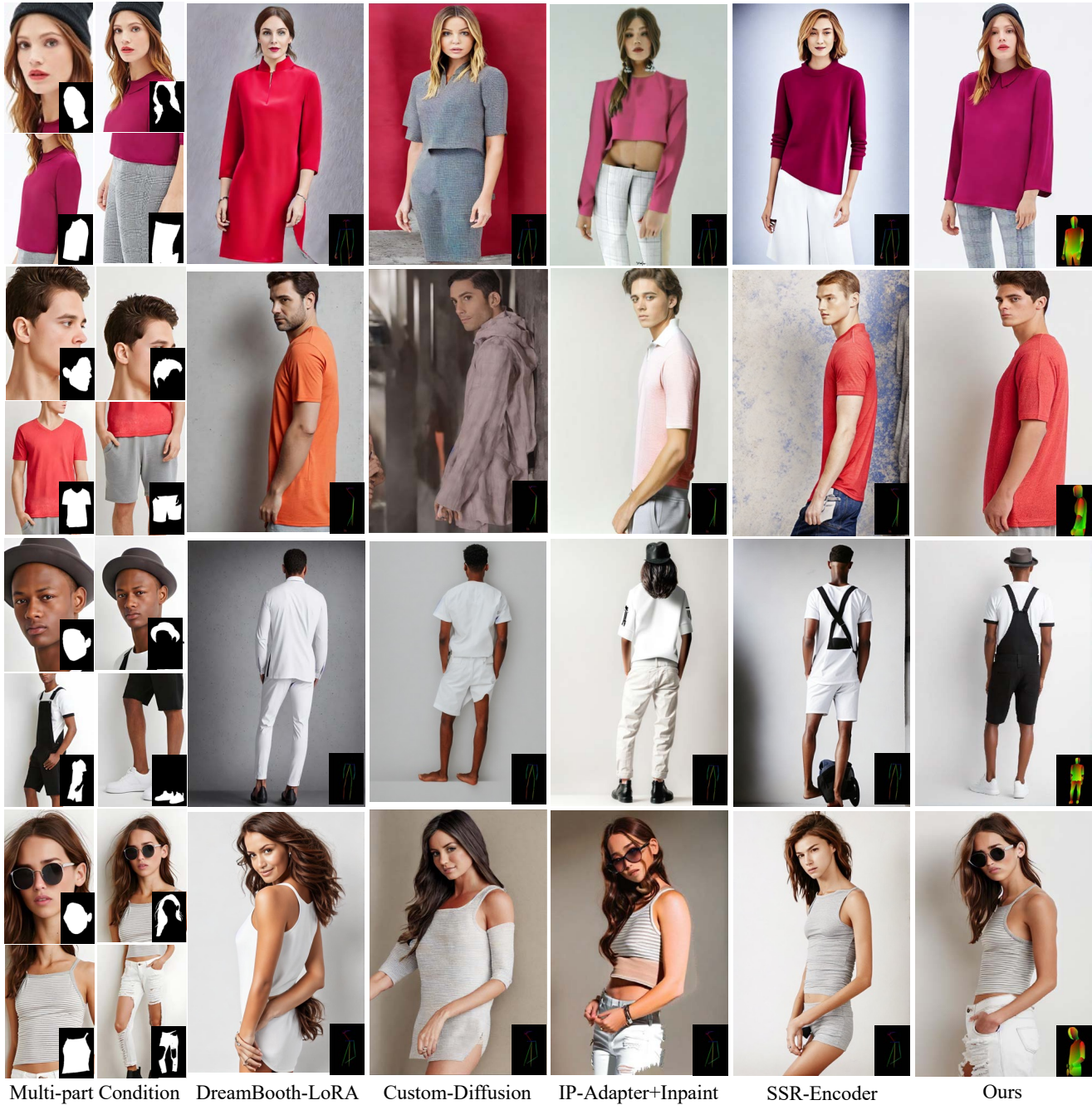| Multi-part Condition | DreamBooth-LoRA | Custom-Diffusion | IP-Adapter+Inpaint | SSR-Encoder | Ours |

Figure 4. Qualitative results generated by Parts2Whole and existing alternatives on our partitioned test set. We do not show the text condition in the figure, but notably, when we input the reference images to our proposed appearance encoder, we will pass in short labels such as face, hair or headwear, upper body clothes, lower body clothes, whole body clothes, shoes, etc.

dataset to enhance its ability for human images.

We compare our Parts2Whole with the above two reference-based alternatives in the test set. For quantitative comparison, we compute the commonly used CLIP score and DINO score to evaluate the similarity between the generated image and the specified human parts. For further alignment evaluation, we use DreamSim [8], a new metric for perceptual image similarity that bridges the gap between "low-level" metrics (e.g., LPIPS, PSNR, SSIM) and "high-level" measures (e.g., CLIP). Since the generated image is conditioned on multiple parts of human appearance, it is difficult to evaluate the degree of alignment by calculating

Table 1. Quantitative comparison between our Parts2Whole and existing reference-based alternatives.

| Method | CLIP↑ | DINO↑ | DreamSim[8]↓ |
|---|---|---|---|
| IP-Adapter [52] + Inpaint | 80.1 | 69.8 | 0.445 |
| SSR-Encoder [57] | 86.9 | 75.1 | 0.346 |
| Parts2Whole (Ours) | **91.2** | **93.7** | **0.221** |

Table 2. User study on the comparison with existing reference-based alternatives. "Quality" and "Similarity" measures synthesis quality and appearance preservation. Each metric is rated from 1 (worst) to 5 (best).

| Method | Quality↑ | Similarity↑ |
|---|---|---|
| IP-Adapter [52] + Inpaint | 3.78 | 3.58 |
| SSR-Encoder [57] | 3.64 | 3.14 |
| Parts2Whole (Ours) | **4.52** | **4.55** |

the average metric with these multiple images. Therefore, we calculate the above three indicators between the output image and the original reference portrait from which these different parts come. We present the quantitative results in Tab. 1 and the qualitative results in Fig. 4. Both IP-Adapter and SSR-Encoder fail to maintain alignment with the specified appearance images and often produce unrealistic results when multi-part combinations are involved. In comparison, our method achieves the best results in terms of image quality and appearance alignment.

**User Study.** We conduct a user study to further evaluate the reference-based methods IP-Adapter [52], SSR-Encoder [57] and our Parts2Whole. We randomly select 20 pairs of reference-target pairs from the test set. For each pair, we provide multiple referential appearance images, pose images, textual captions, and the generated human images. We evaluate the performance from two main aspects: first, the **quality** of the generated images, which primarily refers to the realism, rationality, and clarity of the images; and second, the **similarity** between the generated images and the reference images. The similarity assessment includes consistency in ID, pose, texture, and color between the generated images and the reference images. We involve 20 users in the user study, who are required to score the three methods based on these two evaluative aspects. The final experimental results are shown in Tab. 2, from which we observe that our model owns obvious superiorities for alignment with given appearance conditions.

### 4.3. Ablation Studies

**Appearance Encoder for Multiple Images.** As described in Sec. 3.2, to extract detailed features from multiple reference images, our Parts2Whole designs an appearance en-

Table 3. Quantitative analysis of using semantic-aware encoder and mask-guided subject selection.

| Method | CLIP↑ | DINO↑ | DreamSim[8]↓ | FID↓ |
|---|---|---|---|---|
| w/o text labels | 90.1 | 91.9 | 0.248 | 23.95 |
| w/o mask | 90.8 | 91.6 | 0.243 | 19.79 |
| Parts2Whole | **91.2** | **93.7** | **0.221** | **17.29** |

coder by copying the network structure and pre-trained weights from the denoising U-Net. Here, we compare it to the baseline with other image encoders. Specifically, we leverage the CLIP image encoder [34], DINOv2 [31], and ControlNet [56] as feature extractors and apply the same training settings for fair comparison. The qualitative results of generated human images are presented in Fig. 5. From the second and third columns of the figure, we observe that these semantic-level feature extractors cannot preserve the appearance details of multiple reference images and only extract color and rough texture. ControlNet directly adds different image features with misaligned structures to the feature maps, resulting in unstable image quality. In contrast, our proposed appearance encoder provides fine-grained details of multiple aspects of human appearance.

**Semantic-Aware Encoder.** In the process of encoding multiple reference image features, we provide a textual class label for each aspect of human appearance, thus providing a classifier-like guidance. To assess the effectiveness of the additional external condition, we compare it with directly concatenating multiple reference images in the dimension of width as input to the appearance encoder. As shown in the 5th column in Fig. 5, simply piecing reference images produces images relatively aligned with the given images, but leads to stiff-looking and unrealistic results. This is because modeling only the image itself makes the model lack awareness of different types of appearances. Conversely, after injecting different semantic labels for each reference image, the model has an awareness of various parts of the human appearance, producing realistic and flexible portraits.

**Mask-Guided Subject Selection.** To precisely select subjects from multiple reference images, we introduce the subject masks into the shared self-attention mechanism. To evaluate the effectiveness of our proposed mask-guided attention, we compare it with that without masks. When not using subject masks, the image is generated with reference to all patches of the conditional images, including the unexpected background or other parts. As shown in Fig. 5, due to the generation being interfered with by irrelevant subjects, the model produces homogeneous colors or appears with unexpected backgrounds or subjects. On the contrary, with the support of mask-guided attention, Parts2Whole accurately refers to the appearance of the specified parts to generate real human images.

Figure 5. Qualitative analysis of using different backbones for the appearance encoder, and our proposed methods.
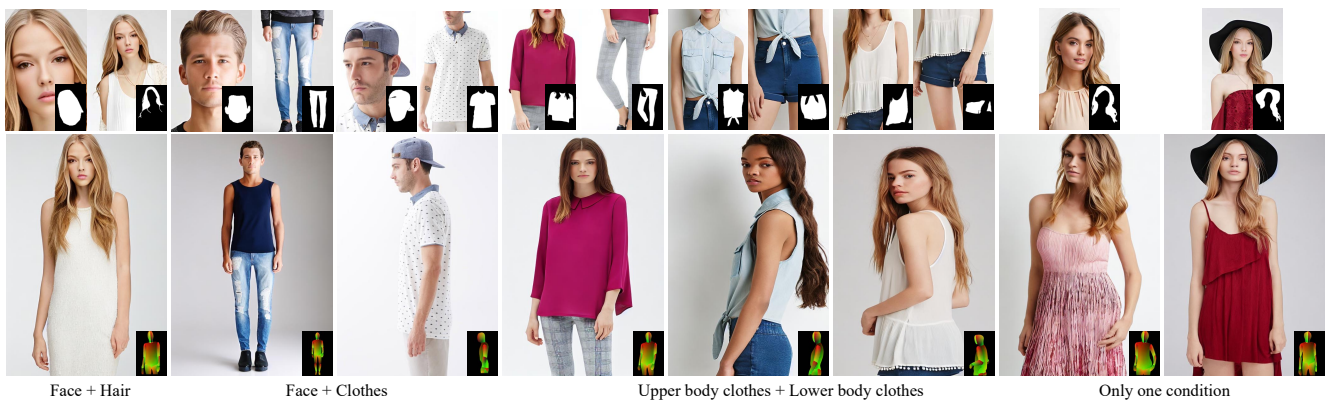
Multi-part Condition    CLIP    DINO    ControlNet    w/o text labels    w/o mask    Parts2Whole



Face + Hair    Face + Clothes    Upper body clothes + Lower body clothes    Only one condition

Figure 6. The generated results from combinations of a different number of conditions.

## 4.4. More Results

**Body Parts of Any Quantity.** Our Parts2Whole is able to generate human images from varying numbers of condition images, such as single hair or face input, or arbitrary combinations like "Face + Hair", "Face + Clothes", and "Upper body clothes + Lower body clothes". The experimental results are presented in Fig. 6. The generated results under different control condition combinations still maintain high quality and realism. This flexibility enables our method to have broader application.

**Multiple Parts from Different Humans.** We select various parts from different human images to serve as conditional images. For example, the face from person A, the hair or headwear from person B, the upper clothes from person

C, and the lower clothes from person D. These parts are collectively used as control conditions for generation. The experimental results, as shown in Fig. 7, demonstrate that our method not only accurately maps different parts of the reference image to the corresponding regions in the target image but also effectively preserves the details of the conditions, producing realistic images.

## 5. Conclusion

In this work, we propose Parts2Whole, a novel framework for controllable human image generation conditioned on multiple reference images, including various aspects of human appearance (e.g., hair, face, clothes, shoes, etc.) and pose maps. Based on a dual U-Net design, we develop a semantic-aware appearance encoder to process each condi-

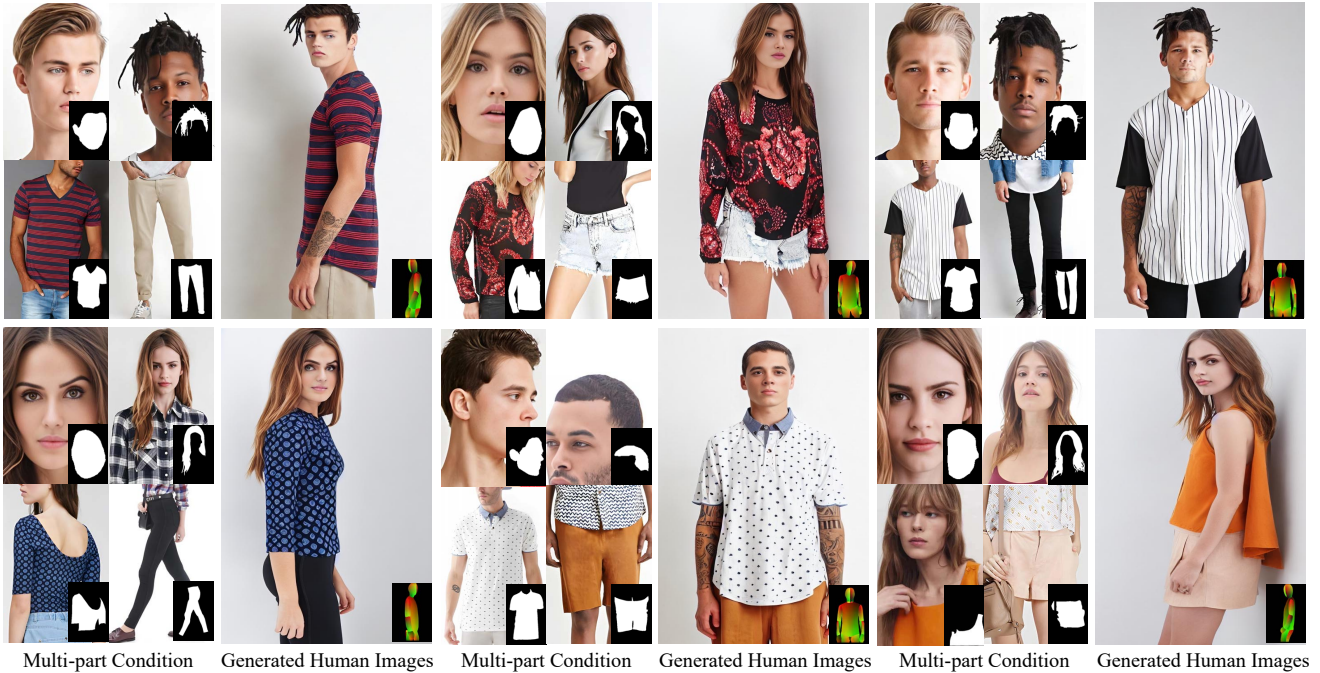| Multi-part Condition | Generated Human Images | Multi-part Condition | Generated Human Images | Multi-part Condition | Generated Human Images |

Figure 7. Results of image generation using selected parts from different individuals as control conditions.

tion image with its label into multi-scale feature maps and inject those detail-rich reference features into the generation via a shared self-attention mechanism. This design retains details from multiple references and looks very good. We also enhance vanilla self-attention by incorporating subject masks, enabling Parts2Whole to synthesize human images from specified parts from condition images. Extensive experiments demonstrate that our Parts2Whole performs well in terms of image quality and condition alignment.

**Future Works.** Our Parts2Whole is currently trained at the resolution of 512, which may cause artifacts in some generated results. This could be improved by using higher resolutions and larger diffusion models like SD-XL [32] as our backbone. Furthermore, it will be valuable to achieve the try-on of layer-wise clothing based on our Parts2Whole.

# References

[1] Mingdeng Cao, Xintao Wang, Zhongang Qi, Ying Shan, Xiaohu Qie, and Yinqiang Zheng. Masactrl: Tuning-free mutual self-attention control for consistent image synthesis and editing. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 22560–22570, 2023. 5

[2] Huiwen Chang, Han Zhang, Jarred Barber, AJ Maschinot, Jose Lezama, Lu Jiang, Ming-Hsuan Yang, Kevin Murphy, William T Freeman, Michael Rubinstein, et al. Muse: Text-to-image generation via masked generative transformers. *arXiv preprint arXiv:2301.00704*, 2023. 3

[3] Xi Chen, Lianghua Huang, Yu Liu, Yujun Shen, Deli Zhao, and Hengshuang Zhao. Anydoor: Zero-shot object-level image customization. *arXiv preprint arXiv:2307.09481*, 2023. 2, 3

[4] Seunghwan Choi, Sunghyun Park, Minsoo Lee, and Jaegul Choo. Viton-hd: High-resolution virtual try-on via misalignment-aware normalization. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 14131–14140, 2021. 1

[5] Jiankang Deng, Jia Guo, et al. InsightFace: 2D and 3D Face Analysis Project. https://github.com/deepinsight/insightface, 2018. Accessed: 2024-04-11. 6, 1

[6] Jiankang Deng, Jia Guo, Xue Niannan, and Stefanos Zafeiriou. Arcface: Additive angular margin loss for deep face recognition. In *CVPR*, 2019. 6, 1

[7] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021. 2, 3

[8] Stephanie Fu, Netanel Tamir, Shobhita Sundaram, Lucy Chai, Richard Zhang, Tali Dekel, and Phillip Isola. Dreamsim: Learning new dimensions of human visual similarity using synthetic data, 2023. 7, 8

[9] Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit H. Bermano, Gal Chechik, and Daniel Cohen-Or. An image is worth one word: Personalizing text-to-image generation using textual inversion, 2022. 2, 3

[10] Rinon Gal, Moab Arar, Yuval Atzmon, Amit H Bermano, Gal Chechik, and Daniel Cohen-Or. Encoder-based domain tuning for fast personalization of text-to-image models. *ACM Transactions on Graphics (TOG)*, 42(4):1–13, 2023. 3

[11] Amir Hertz, Andrey Voynov, Shlomi Fruchter, and Daniel

Cohen-Or. Style aligned image generation via shared attention. *arXiv preprint arXiv:2312.02133*, 2023. 4

[12] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. 2, 3

[13] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021. 2, 3, 6, 1

[14] Li Hu, Xin Gao, Peng Zhang, Ke Sun, Bang Zhang, and Liefeng Bo. Animate anyone: Consistent and controllable image-to-video synthesis for character animation. *arXiv preprint arXiv:2311.17117*, 2023. 2, 4, 5

[15] Minghui Hu, Jianbin Zheng, Daqing Liu, Chuanxia Zheng, Chaoyue Wang, Dacheng Tao, and Tat-Jen Cham. Cocktail: Mixing multi-modality control for text-conditional image generation. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. 3

[16] Zehuan Huang, Hao Wen, Junting Dong, Yaohui Wang, Yangguang Li, Xinyuan Chen, Yan-Pei Cao, Ding Liang, Yu Qiao, Bo Dai, et al. Epidiff: Enhancing multi-view synthesis via localized epipolar-constrained diffusion. *arXiv preprint arXiv:2312.06725*, 2023. 3

[17] Jaeseok Jeong, Junho Kim, Yunjey Choi, Gayoung Lee, and Youngjung Uh. Visual style prompting with swapping self-attention. *arXiv preprint arXiv:2402.12974*, 2024. 4

[18] Yuming Jiang, Shuai Yang, Haonan Qiu, Wayne Wu, Chen Change Loy, and Ziwei Liu. Text2human: Text-driven controllable human image generation. *ACM Transactions on Graphics (TOG)*, 41(4):1–11, 2022. 2, 3, 6, 1

[19] Zeyinzi Jiang, Chaojie Mao, Yulin Pan, Zhen Han, and Jingfeng Zhang. Scedit: Efficient and controllable image diffusion generation via skip connection editing. *arXiv preprint arXiv:2312.11392*, 2023. 3

[20] Diederik P Kingma and Max Welling. Auto-encoding variational bayes, 2022. 3

[21] Nupur Kumari, Bingliang Zhang, Richard Zhang, Eli Shechtman, and Jun-Yan Zhu. Multi-concept customization of text-to-image diffusion, 2023. 2, 3, 6, 1

[22] Lambda Labs. Stable diffusion image variations, 2022. 6

[23] Zhen Li, Mingdeng Cao, Xintao Wang, Zhongang Qi, Ming-Ming Cheng, and Ying Shan. Photomaker: Customizing realistic human photos via stacked id embedding. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024. 2, 3

[24] Xian Liu, Jian Ren, Aliaksandr Siarohin, Ivan Skorokhodov, Yanyu Li, Dahua Lin, Xihui Liu, Ziwei Liu, and Sergey Tulyakov. Hyperhuman: Hyper-realistic human generation with latent structural diffusion. *arXiv preprint arXiv:2310.08579*, 2023. 2

[25] Ziwei Liu, Ping Luo, Shi Qiu, Xiaogang Wang, and Xiaoou Tang. Deepfashion: Powering robust clothes recognition and retrieval with rich annotations. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 6, 1

[26] Zhiheng Liu, Ruili Feng, Kai Zhu, Yifei Zhang, Kecheng Zheng, Yu Liu, Deli Zhao, Jingren Zhou, and Yang Cao.

[27] Jian Ma, Junhao Liang, Chen Chen, and Haonan Lu. Subject-diffusion: Open domain personalized text-to-image generation without test-time fine-tuning. *arXiv preprint arXiv:2307.11410*, 2023. 3

[28] Davide Morelli, Matteo Fincato, Marcella Cornia, Federico Landi, Fabio Cesari, and Rita Cucchiara. Dress code: high-resolution multi-category virtual try-on. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2231–2235, 2022. 1

[29] Chong Mou, Xintao Wang, Liangbin Xie, Yanze Wu, Jian Zhang, Zhongang Qi, Ying Shan, and Xiaohu Qie. T2i-adapter: Learning adapters to dig out more controllable ability for text-to-image diffusion models. *arXiv preprint arXiv:2302.08453*, 2023. 2, 3

[30] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models, 2022. 2, 3

[31] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Mahmoud Assran, Nicolas Ballas, Wojciech Galuba, Russell Howes, Po-Yao Huang, Shang-Wen Li, Ishan Misra, Michael Rabbat, Vasu Sharma, Gabriel Synnaeve, Hu Xu, Hervé Jegou, Julien Mairal, Patrick Labatut, Armand Joulin, and Piotr Bojanowski. Dinov2: Learning robust visual features without supervision, 2024. 8

[32] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis. *arXiv preprint arXiv:2307.01952*, 2023. 3, 10

[33] Can Qin, Shu Zhang, Ning Yu, Yihao Feng, Xinyi Yang, Yingbo Zhou, Huan Wang, Juan Carlos Niebles, Caiming Xiong, Silvio Savarese, et al. Unicontrol: A unified diffusion model for controllable visual generation in the wild. *arXiv preprint arXiv:2305.11147*, 2023. 3

[34] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision, 2021. 5, 8

[35] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1 (2):3, 2022. 2, 3

[36] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 2, 3, 4, 5, 6

[37] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. 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*, pages 234–241. Springer, 2015. 4

[38] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. 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*, pages 22500–22510, 2023. 2, 3, 6, 1

[39] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in neural information processing systems*, 35:36479–36494, 2022. 2, 3

[40] Jing Shi, Wei Xiong, Zhe Lin, and Hyun Joon Jung. Instantbooth: Personalized text-to-image generation without test-time finetuning, 2023. 2, 3

[41] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020. 6

[42] Yoad Tewel, Omri Kaduri, Rinon Gal, Yoni Kasten, Lior Wolf, Gal Chechik, and Yuval Atzmon. Training-free consistent text-to-image generation. *arXiv preprint arXiv:2402.03286*, 2024. 4

[43] Keyu Tian, Yi Jiang, Zehuan Yuan, Bingyue Peng, and Liwei Wang. Visual autoregressive modeling: Scalable image generation via next-scale prediction. 2024. 3

[44] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. 4

[45] Qixun Wang, Xu Bai, Haofan Wang, Zekui Qin, and Anthony Chen. Instantid: Zero-shot identity-preserving generation in seconds. *arXiv preprint arXiv:2401.07519*, 2024. 2, 3

[46] Xintao Wang, Liangbin Xie, Chao Dong, and Ying Shan. Real-esrgan: Training real-world blind super-resolution with pure synthetic data. In *International Conference on Computer Vision Workshops (ICCVW)*. 6, 1

[47] Yuxiang Wei, Yabo Zhang, Zhilong Ji, Jinfeng Bai, Lei Zhang, and Wangmeng Zuo. Elite: Encoding visual concepts into textual embeddings for customized text-to-image generation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15943–15953, 2023. 3

[48] Guangxuan Xiao, Tianwei Yin, William T Freeman, Frédo Durand, and Song Han. Fastcomposer: Tuning-free multi-subject image generation with localized attention. *arXiv preprint arXiv:2305.10431*, 2023. 3

[49] Zhongcong Xu, Jianfeng Zhang, Jun Hao Liew, Hanshu Yan, Jia-Wei Liu, Chenxu Zhang, Jiashi Feng, and Mike Zheng Shou. Magicanimate: Temporally consistent human image animation using diffusion model. 2024. 2, 4, 5

[50] Binxin Yang, Shuyang Gu, Bo Zhang, Ting Zhang, Xuejin Chen, Xiaoyan Sun, Dong Chen, and Fang Wen. Paint by example: Exemplar-based image editing with diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18381–18391, 2023. 5

[51] Zhendong Yang, Ailing Zeng, Chun Yuan, and Yu Li. Effective whole-body pose estimation with two-stages distillation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4210–4220, 2023. 6, 1

[52] Hu Ye, Jun Zhang, Sibo Liu, Xiao Han, and Wei Yang. Ip-adapter: Text compatible image prompt adapter for text-to-image diffusion models. 2023. 2, 3, 5, 6, 8, 1

[53] Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gunjan Baid, Zirui Wang, Vijay Vasudevan, Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, Ben Hutchinson, Wei Han, Zarana Parekh, Xin Li, Han Zhang, Jason Baldridge, and Yonghui Wu. Scaling autoregressive models for content-rich text-to-image generation, 2022. 3

[54] Polina Zablotskaia, Aliaksandr Siarohin, Bo Zhao, and Leonid Sigal. Dwnet: Dense warp-based network for pose-guided human video generation. *arXiv preprint arXiv:1910.09139*, 2019. 1

[55] Lvmin Zhang. Reference-only controlnet, 2023. 5

[56] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models, 2023. 2, 3, 5, 6, 8

[57] Yuxuan Zhang, Yiren Song, Jiaming Liu, Rui Wang, Jinpeng Yu, Hao Tang, Huaxia Li, Xu Tang, Yao Hu, Han Pan, and Zhongliang Jing. Ssr-encoder: Encoding selective subject representation for subject-driven generation, 2024. 2, 3, 5, 6, 8, 1

[58] Shihao Zhao, Dongdong Chen, Yen-Chun Chen, Jianmin Bao, Shaozhe Hao, Lu Yuan, and Kwan-Yee K Wong. Uni-controlnet: All-in-one control to text-to-image diffusion models. *Advances in Neural Information Processing Systems*, 36, 2024. 2, 3

# From Parts to Whole: A Unified Reference Framework for Controllable Human Image Generation

## Supplementary Material

## A. Dataset

Generating the human body from multiple conditional parts is a significant undertaking, but lacking a directly available dataset. The datasets related to this task, such as those for Virtual Try-on[4, 28, 54], primarily suffer from a lack of multiple controllable conditions and are often limited to single control conditions (clothing). Issues with these datasets include a limited variety of clothing types, absence of facial data, low resolution, and lack of textual captions. The DeepFashion-MultiModal dataset [18, 25] aligns more closely with our task as it includes a vast array of human body images, the same person and same clothes in different poses, and precise human parsing labels. However, this dataset cannot be used directly and requires data cleansing and further post-processing.

**ID Cleansing.** In the DeepFashion-MultiModal dataset, there is some confusion with IDs where images of different individuals are mistakenly labeled under the same ID. We start by cleansing these IDs, extracting facial ID features from images tagged with the same ID using InsightFace[5, 6]. Cosine similarity is then used to evaluate the similarity between image ID feature pairs, allowing us to reclassify IDs within the same ID group. After cleansing, images from the same ID and the same clothes are selected if there are two or more images available.

**Building Reference-Target Pair.** We use images with human parsing labels from the dataset as reference images. Target images are then selected from the same ID and clothing, creating pairs with the reference image.

**Obtaining Reference Human Part.** We crop the images according to the provided human parsing labels. Specifically, we divide the human image into six parts: upper body clothes, lower body clothes, whole body clothes, hair or headwear, face, and footwear. Each part is cropped according to the human parsing labels to obtain the crop image and corresponding mask image. Due to the low resolution of the cropped parts, we apply Real-ESRGAN[46] to enhance the image resolution, thus obtaining clearer reference images.

**Obtaining Target Description.** Based on the reference human parts, we need to generate images that resemble the target image, requiring a description of the target image. The description is divided into two parts: one for the human body's pose and another for the target image's textual description. For pose information, we utilize DWPose[51] to generate pose images corresponding to each image, and for DensePose, we use the provided DensePose files from the dataset. The textual description for each image is taken directly from the dataset's accompanying text description.

**Introduction to the Final Dataset.** Finally, we have constructed a multimodal dataset with approximately 41,500 reference-target pairs derived from the open-source DeepFashion-MultiModal dataset [18, 25]. The controllable conditions for each pair are categorized into two main types. The first type is the appearance reference image, which is subdivided into six parts: upper body clothes, lower body clothes, whole body clothes, hair or headwear, face, and footwear. Each image is accompanied by a corresponding mask and has undergone super-resolution processing. These data elements are sourced from the original reference image. The second type is the target description, primarily consisting of pose and text description. The pose is further divided into OpenPose and DensePose, all of which are derived from the target image. A sample of reference-target image pair in our dataset is shown in Fig. 8.

## B. Different Types of Pose Maps

To show the ability of Parts2Whole to generate human image conditions on different types of pose maps, we train a new Parts2Whole model but with **OpenPose** as a condition. As shown in Fig. 9, the generated images strictly maintain consistency with the target pose, and each body part retains the appearance information from the reference images.

## C. Discussion about Existing Methods

In our main text, we conduct experiments with both tuning-based methods DreamBooth LoRA [13, 38] and Custom Diffusion [21] and tuning-free methods such as IP-Adapter [52] and SSR-Encoder [57]. The results of these experiments are further illustrated in Fig 10, which showcases the performance of these methods conditioned on both a single image and multiple images.

The results demonstrate that while these methods generally perform well under a single control condition, they exhibit significant issues when multiple conditions are applied. For example, DreamBooth LoRA encountered cases where pants were omitted, and both Custom Diffusion and SSR-Encoder show alterations in facial ID. The phenomenon observed is attributed to the spatial misalignment of the input multi-body parts with the target image, and the lack of specific design in existing methods to address the variation in spatial positions during feature injection. For instance, methods like SSR-Encoder, Custom Diffu-
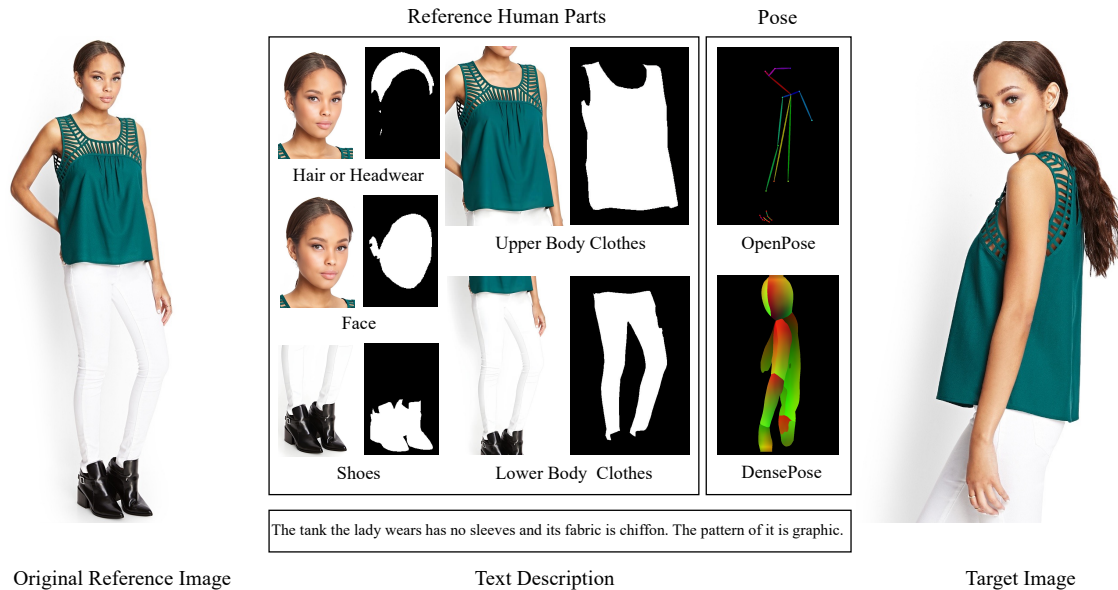
Figure 8. A sample of reference-target image pair in our dataset. Reference Human Part images are obtained from the reference image based on the human parsing image, while the pose and text description are descriptions of the target image.



Figure 9. Generated human images with OpenPose as pose condition by Parts2Whole. The upper row displays multiple reference human parts, while the lower row shows the results generated by our Parts2Whole method under the control of reference images and OpenPose image.

sion, and IP-Adapter incorporate features into the denoising UNet through cross-attention mechanisms. They encode reference images into other modal features (e.g. semantic features) and utilize the cross-attention keys ($K$) and values ($V$) from them rather than from **image dimensional** feature maps. In this process, the correlation between the reference images and the target image **loses the spatial re-** **lationship** of the original image dimensions. It is difficult for these methods to effectively model the attention from various conditional feature maps at different locations in the target image, resulting in a mixture of attributes from different subjects.

Conversely, our Parts2Whole model employs shared self-attention between the reference features and the feature

Figure 10. Results of single-condition and multi-condition generation. The top three rows labeled "face", "upper body clothes", "lower body clothes" represent separate conditions. The fourth row demonstrates joint control under multi-part conditions.

maps in the Denoising U-Net, executed on the image dimension. This allows our model to establish a more precise correlation between different condition images and distinct positions within the feature maps, thereby generating results that are consistent with the detailed attributes of multi-condition images.