

StyleGAN-Human: A Data-Centric Odyssey of Human Generation

Jianglin Fu^{1*} Shikai Li^{1*} Yuming Jiang² Kwan-Yee Lin¹ Chen Qian¹

Chen Change Loy² Wayne Wu^{1,3} Ziwei Liu²

¹ SenseTime Research ² S-Lab, Nanyang Technological University ³ Shanghai AI Laboratory

{fujianglin, lishikai, linjunyi, qianchen}@sensetime.com

{yuming002, ccloy, ziwei.liu}@ntu.edu.sg

wuwenyan0503@gmail.com



Figure 1. **A Data-Centric Odyssey of Human Generation.** With good “data engineering” practices, given a random latent code z , the StyleGAN-Human model could generate high-resolution photo-realistic human images as presented. Zoom in for the best view.

Abstract

Unconditional human image generation is an important task in vision and graphics, enabling various applications in the creative industry. Existing studies in this field mainly focus on “network engineering” such as designing new components and objective functions. This work takes a data-centric perspective and investigates multiple critical aspects in “data engineering”, which we believe would complement the current practice. To facilitate a comprehensive study, we collect and annotate a large-scale human image dataset with over 230K samples capturing di-

verse poses and textures. Equipped with this large dataset, we rigorously investigate three essential factors in data engineering for StyleGAN-based human generation, namely data size, data distribution, and data alignment. Extensive experiments reveal several valuable observations w.r.t. these aspects: 1) Large-scale data, more than 40K images, are needed to train a high-fidelity unconditional human generation model with a vanilla StyleGAN. 2) A balanced training set helps improve the generation quality with rare face poses compared to the long-tailed counterpart, whereas simply balancing the clothing texture distribution does not effectively bring an improvement. 3) Human GAN models that employ body centers for alignment outperform

*Equal contribution.

models trained using face centers or pelvis points as alignment anchors. In addition, a model zoo and human editing applications are demonstrated to facilitate future research in the community. Code and models are publicly available¹.

1. Introduction

Generating photo-realistic images of clothed humans unconditionally can provide great support for downstream tasks such as human motion transfer [9, 47], digital human animation [45], fashion recommendation [35, 43], and virtual try-on [16, 59, 89]. Traditional methods create dressed humans with classical graphics modeling and rendering processes [20, 34, 53, 64, 66, 74, 82, 97]. Although impressive results have been achieved, these prior works are easy to suffer from the limitation of robustness and generalizability in complex environments. Recent years, Generative Adversarial Networks (GANs) have demonstrated remarkable abilities in real-world scenarios, generating diverse and realistic images by learning from large-quantity and high-quality datasets. [27, 36, 39, 69].

Among the GAN family, StyleGAN2 [40] stands out in generating faces and simple objects with unprecedented image quality. A major driver behind recent advancements [2, 37, 40, 84, 94] on such StyleGAN architectures is the prosperous discovery of “network engineering” like designing new components [2, 37, 94] and loss functions [40, 84]. While these approaches show compelling results in generating diverse objects (*e.g.*, faces of humans and animals), applying them to the photo-realistic generation of articulated humans in natural clothing is still a challenging and open problem.

In this work, we focus on the task of Unconditional Human Generation, with a specific aim to train a good StyleGAN-based model for articulated humans from a *data-centric* perspective. First, to support the data-centric investigation, collecting a large-scale, high-quality, and diverse dataset of human bodies in clothing is necessary. We propose the Stylish-Humans-HQ Dataset (SHHQ), which contains 230K clean full-body images with a resolution of 1024×512 at least and up to 2240×1920 . The SHHQ dataset lays the foundation for extensive experiments on unconditional human generation. Second, based on the proposed SHHQ dataset, we investigate three fundamental and critical questions that were not thoroughly discussed in prior works and attempt to provide useful insights for future research on unconditional human generation.

To extract the questions that are indeed *important* for the community of Unconditional Human Generation, we make an extensive survey on recent literature in the field of general unconditional generation [5, 6, 23, 27, 36, 39, 54]. Based on the survey, three questions that are investigated actively can be concluded as below. **Question-1:** What is the re-

lationship between the *data size* and the generation quality? Several previous works [6, 31, 37, 85, 102] pointed out that the quantity of training data is the primary factor to determine the strategy for improving image quality in face and other object generation tasks. In this study, we want to examine the minimum quantity of training data required to generate human images of high quality without any extensive “network engineering” effort. **Question-2:** What is the relationship between the *data distribution* and the generation quality? This question has received extensive attention [15, 25, 55, 75, 96] and leads to a research topic dealing with data imbalance [49]. In this study, we aim to exploit data imbalance problem in the human generation task. **Question-3:** What is the relationship between the scheme of *data alignment* and the generation quality? Different alignment schemes applied to uncured faces [39, 41] and non-rigid objects [6, 14, 76] show success in enhancing training performance. In this study, we seek a better data alignment strategy for human generation.

Based on the proposed SHHQ dataset and observations from our experiments, we establish a Model Zoo with three widely-adopted unconditional generation models, *i.e.*, StyleGAN [39], StyleGAN2 [40], and alias-free StyleGAN [38], in both resolution of 1024×512 and 512×256 . Although hundreds of StyleGAN-based studies exist for *face* generation/editing tasks, a high-quality and public model zoo for *human* generation/editing with StyleGAN family is still missing. The provided model zoo is positioned to complement existing facial model zoo. We believe it has great potentials in many human-centric tasks, *e.g.*, human editing, neural rendering, and virtual try-on.

We further construct a human editing benchmark by adapting previous editing methods based on facial models to human body models (*i.e.*, PTI [71] for image inversion, InterFaceGAN [79], StyleSpace [94], and SeFa [80] for image manipulation). The impressive results in editing human clothes and attributes demonstrate the potential of the given model zoo in downstream tasks. In addition, a concurrent work, InsetGAN [18], is also evaluated with our baseline model, further showing the potential usage of our pre-trained human generative models.

Here is the summary of the main contributions of this paper: 1) We collect a large-scale, high-quality, and diverse dataset, Stylish-Humans-HQ (SHHQ), containing 230K human full-body images for unconditional human generation task. 2) We investigate three crucial questions that have aroused broad interest in the community and discuss our observation through comprehensive analysis. 3) We build a model zoo for unconditional human generation to facilitate future research. An editing benchmark is also established to demonstrate the potential of the proposed model zoo.

2. Related Work

2.1. Dataset For Human Generation

Large-scale and high-quality clothed human-centric training datasets are the critical fuel for the training of Style-

¹Project page: <https://stylegan-human.github.io/>
Code and models: <https://github.com/stylegan-human/StyleGAN-Human>

GAN models. A qualified dataset should conform to the following aspects: **1) Image quality:** high-resolution images with rich textures offer more raw detailed semantic information to the model. **2) Data volume:** the size of dataset should be sufficient to avoid generative overfitting [4, 101]. **3) Data coverage:** the dataset should cover multiple attribute dimensions to guarantee diversity of the model, for instance, gender, clothing type, clothing texture, and human pose. **4) Data content:** since this report only focuses on the generation of single full-body human, occlusion caused by other people or objects is not considered here, whereas self-occlusion is taken into account. That is, each image should contain only one complete human body.

Publicly available datasets built particularly for full human-body generation are rare, but there are several practices [51, 52, 81] cooperating with DeepFashion [48] and Market1501 [103]. DeepFashion dataset [48] with well-labeled attributes and diverse garment categories is satisfactory for image classification and attribute prediction, but not adequate for unconditional human generation since it emphasizes fashion items rather than human bodies. Thus the number of close-up shots of clothing is much higher than that of full-body images. Market1501 dataset [103] fails for human generation tasks due to its low resolution (128×64). There are some human-related datasets in other domains rather than GAN-based applications: datasets related to human parsing [22, 46] are limited by scalability and diversity; common datasets for virtual try-on tasks either contain only the upper body [28] or are not public [100]. A detailed comparison of the above datasets in terms of data scale, average resolution, attributes labeling, and proportion of full-body images across the whole dataset is listed in Table 1. In general, there is no high-quality and large-scale full human-body dataset publicly available for the generative purpose.

2.2. StyleGAN

In recent years, the research focus has gradually shifted to generating high-fidelity and high-resolution images through Generative Adversarial Networks [6, 36]. The StyleGAN generator [39] was introduced and became the state-of-the-art network of unconditional image generation. Compared to previous GAN-based architectures [5, 27, 58], StyleGAN injects a separate attribute factor (i.e., style) into the generator to influence the appearance of generated images. Then StyleGAN2 [40] redesigns the normalization, multi-scale scheme, and regularization method to rectify the artifacts in StyleGAN images. The latest update to StyleGAN [38] reveals the non-ideal case of detailed textures sticking to fixed pixel locations and proposes an alias-free network.

2.3. Human Generation

In human generation research, most of the existing applications focus on precise control of pose and appearance by leveraging conditional VAE and U-Net [17, 73] or StyleGAN-related architectures [3, 24, 44, 72]. Specifically,

the 3D method [24] renders StyleGAN-generated neural textures on the parametric human models, but the results are restricted by the quantity and quality of training data. The other works [3, 72, 73] preserve texture quality by spatial modulation using the extracted UV texture map, and perform pose transfer conditioned by extracted pose features. The limitation of these works is that paired data with satisfied volume is required for training. Moreover, studies [3, 17, 72, 73] trained with DeepFashion indicate that DeepFashion can produce decent results in human generation, at least for head-to-waist images. These works rely on additional network modifications and certain human priors. The above works can be summarized as “network engineering”; they require some architectural changes and certain priors. In contrast, this study probes unconditional human generation challenges from a data perspective.

2.4. Image Editing

Benefiting from StyleGAN, one of the significant downstream applications is image editing [1, 63, 79, 94, 99]. A standard image editing pipeline usually involves inversion from a real image to the latent space and manipulating the embedded latent code. Existing works for *image inversion* can be categorized into optimization-based [2, 83], encoder-based [71, 86, 91], and hybrid methods [8], which exploit encoders to embed images into latent space first and then refine with optimization. As for *image manipulation*, studies explore the capability of attribute disentanglement in the latent space with supervised [2, 33, 79] and unsupervised [29, 80, 87, 94] networks. In specific, Jiang *et al.* [33] proposes to use manually labeled fine-grained annotations to find non-linear manipulation directions in the StyleGAN latent space, while SeFa [80] search for semantic directions without supervision. StyleSpace [94] defines the style space S and proves that it is more disentangled than W and $W+$ space. In this report, we perform image editing on real images by the inversion method PTI [71] and various directions through the chosen methods [79, 80, 94] on our model to verify whether our model with human images could preserve the characteristics demonstrated on rigid objects.

3. Stylish-Humans-HQ Dataset

To investigate the key factors in unconditional human generation task from a data-centric perspective, we propose a large-scale, high-quality, and diverse dataset, Stylish-Humans-HQ (SHHQ). In this section, we first present the data collection and preprocessing (Section 3.1), in which we construct the SHHQ dataset. Then, we analyze the data statistic (Section 3.2) of Stylish-Humans-HQ dataset to demonstrate the superiority of SHHQ compared to other datasets from a statistical perspective.

3.1. Data Collection and Preprocessing

We first obtain over 500K raw data of human images from the Internet, covering a wide variety of races, ages and clothing styles. Some representative samples of raw data are

Table 1. **Comparison of SHHQ with other publicly available datasets.** The proposed Stylish-Humans-HQ (SHHQ) dataset covers the largest number of human images to date and it is the only dataset that satisfies all four data requirements mentioned in Section 2.1. For each dataset, we report the total number of images, resolutions, label properties, and the proportion of full-body images in the dataset. “Labeled Attributes” records whether the dataset provides human-relevant labels, while “Full-Body Ratio” indicates the proportion of full-body images each dataset contains.

Dataset	Total Image #	Mean Resolution	Labeled Attributes	Full-Body Ratio
ATR [46]	7,700	400×600	✓	76%
Market1501 [103]	32,668	128×64	✓	100%
DeepFashion [48]	146,680	1101×750	✓	6.8%
LIP [22]	50,462	196×345	✓	37%
VITON [28]	16,253	256×192	✗	0%
Stylish-Humans-HQ	231,176	1024×512 up to 2240×1920 in raw data	✓	100%



Figure 2. **Data Preprocessing.** The following types of images will be removed during our data preprocessing pipeline. (a) Images of low resolution. (b) Images of a person not placed in the center. (c) Images with missing body parts, *e.g.*, with presence of the upper body only. (d) Extreme posture, *e.g.*, handstand. (e) Images with multi-person.

shown in Appendix A. We preprocess the data with six factors taken into consideration (*i.e.*, resolution [48], body position [39], body-part occlusion, human pose [39,48], multi-person, and background), which are critical for the quality of a human dataset. After the data preprocessing procedure, we obtain a clean dataset of 231,176 images with high quality; see Figure 6 (a) for examples.

Resolution. We discard images with a resolution lower than 1024×512 (Figure 2 (a)).

Body Position. The position of the body varies widely in different images, as shown in Figure 2 (b). We design a procedure in which each person is appropriately cropped based on human segmentation [13], padded and resized to the same scale, and then placed in the image such that the body center is aligned. The body center is defined as the average coordinate of the entire body using segmentation.

Body-Part Occlusion. This work aims at generating full-body human images, images with any missing body parts are removed (*e.g.*, the half-body portrait shown in Figure 2 (c)).

Human Pose. We remove images with extreme poses (*e.g.*, lying postures, handstand in Figure 2 (d)) to ensure learnability of the data distribution. We exploit human pose estimation [7] to detect those extreme poses.

Multi-Person Images. Some raw images contain multiple persons, such as in Figure 2 (e). In this work, our goal is to generate single-person full-body images, so we keep unoccluded single-person full-body images, and remove those with occluded people.

Background. Some images contain complicated backgrounds, requiring additional representation ability. To focus on the generation of the human body itself and eliminate the influence of various backgrounds, we use a segmentation mask [13] to modify the image background to pure white. The edges of the mask are smoothed by Gaussian blur.

3.2. Data Statistics

Table 1 presents the comparison between SHHQ and other public datasets. **Dataset Scale.** As shown in the table, our proposed SHHQ is currently the largest dataset in scale compared to other datasets. Among them, the data volume of the SHHQ dataset is 1.6 times that of DeepFashion [48] dataset and is much larger than that of others (30 times to ATR [46], 7 times to Market1501 [103], 4.6 times to LIP [22], and 14 times to VITON [28]). **Resolution.** Images from ATR [46], Market1501 [103], LIP [22], and VITON [28] are lower in resolution, which is insufficient for our generation task, while the proposed SHHQ and DeepFashion provide high-definition images up to 2240×1920 . **Labels.** All datasets beside VITON provide various labeled attributes. Specifically, DeepFashion [48] and SHHQ label the clothing types and textures, which is useful for human generation/editing tasks. Full-body ratio denotes the proportion of full-body images in the dataset. **Full-Body Ratio.** Although DeepFashion [48] offers over 146K images with decent resolution, only 6.8% of them are full-body images, while SHHQ achieves a 100% full-body ratio. Appendix A shows a visual comparison among several representative datasets and the proposed SHHQ dataset.

In summary, SHHQ covers the largest number of human images with high-resolution, labeled clothing attributes, and 100% full-body ratio. It again confirms that our dataset

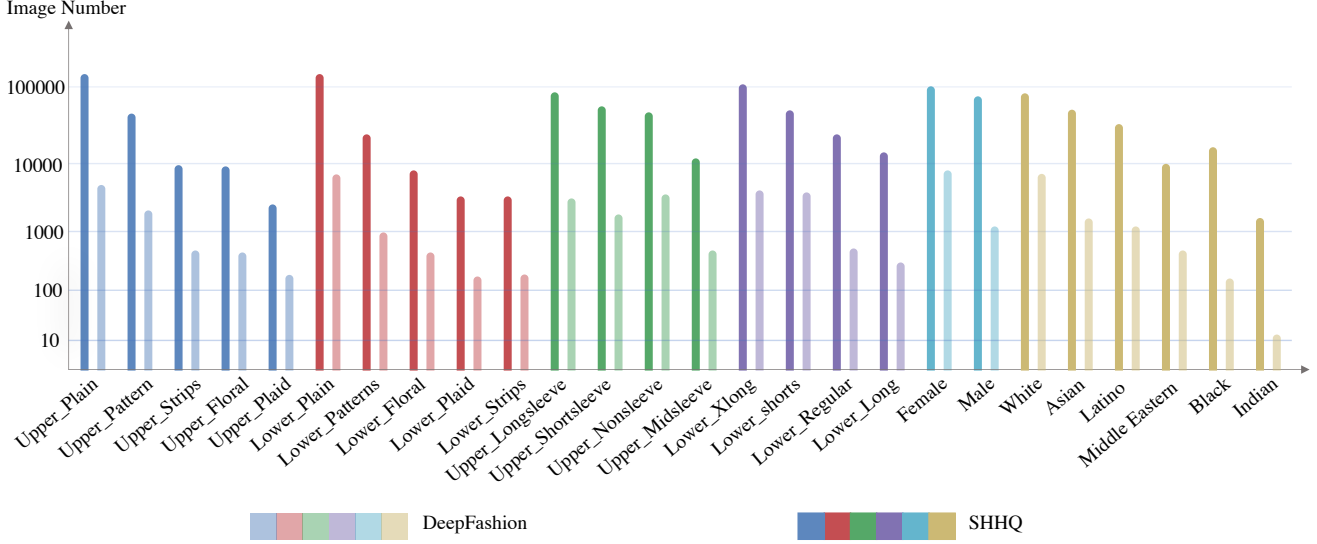


Figure 3. **Attribute Distribution.** Comparison of different attributes between the pruned DeepFashion and our Stylish-Humans-HQ dataset: Texture of the upper clothing, texture of the lower clothing, length of upper clothing, length of lower clothing, gender, and ethnicity. Note that the scale of y-axis uses *logarithmic* scaling, with a base of 10.

is more suitable for full-body human generation than other public datasets.

Of all the datasets compared above, DeepFashion [48] is the most relevant to our human generation task. In Figure 3, we further present the comparison of different attributes between filtered DeepFashion [48] (removing occluded body) and SHHQ in a more detailed view. The bar chart depicts the distributions along six dimensions: upper cloth texture, lower cloth texture, upper cloth length, lower cloth length, gender, and ethnicity. In particular, the number of females is approximately 4 times the number of males in filtered DeepFashion [48], while our dataset features a more balanced female-to-male ratio of 1.49. With the help of DeepFace API [78], it is shown that SHHQ is more diverse in terms of ethnicity. Advantages are also shown in the other five attributes. In terms of garment-related attributes, images with specific labels in filtered DeepFashion [48] are too scarce to be used as a training set. The Stylish-Humans-HQ dataset boosts the number of each category by an average of 24.4 times.

4. Systematic Investigation

Our investigations are built on the official StyleGAN2 codebase² and StyleGAN2 architecture. The detailed training settings can be found in Appendix C.

We conduct extensive experiments to study three factors concerning the quality of generated images: 1) data size (Section 4.1), 2) data distribution (Section 4.2), and 3) data alignment (Section 4.3).

4.1. Data Size

Motivation. Data size is an essential factor that determines the quality of generated images. Previous literature always takes different strategies to improve the generation performance according to different dataset sizes: regularization techniques [6] are employed to train a large dataset, while augmentation [37, 85, 102] and conditional feature transferring [55, 92] are proposed to tackle the limited data of faces and non-rigid objects. Here, we design sets of experiments to examine the relationship between training data size and the image quality of generated humans.

Experimental Settings. To determine the relationship between data size and image quality for the unconditional human GAN, we construct 6 sub-datasets and denoted these subsets as S_0 (10K), S_1 (20K), S_2 (40K), S_3 (80K), S_4 (160K) and S_5 (230K). Here, S_0 is the pruned DeepFashion dataset. We perform the training on two resolution settings for each set: 1024×512 and 512×256 . Considering the case of limited data, we also conduct additional training experiments with adaptive discriminator augmentation (ADA) [37] for small datasets S_0 , S_1 , and S_2 . Fréchet Inception Distance (FID) and Inception Score (IS) are the indicators for evaluating the model performance.

Results. As shown in Figure 4 (a), the FID scores (solid lines) decrease as the size of the training dataset increases for both resolution settings. The declining trend is gradually flattening and tends to converge. S_0 generates the least satisfactory results, with FID of 7.80 and 7.23 for low- and high-resolution, respectively, while S_1 achieves corresponding improvements of 42% and 40% on FID with only an additional 10K training images. When the training size reaches 40K for both resolutions, the FID curves start to

²<https://github.com/NVlabs/stylegan2-ada-pytorch>

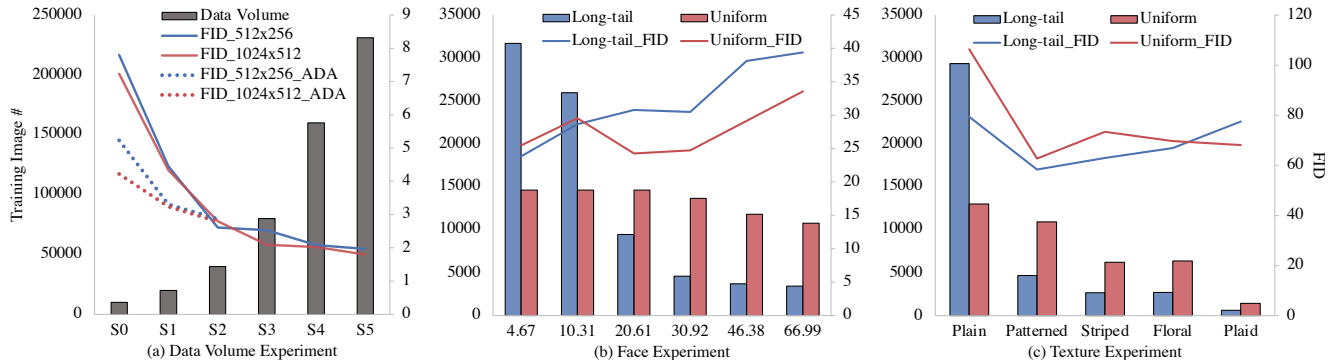


Figure 4. **Experiment results.** (a) FID scores for experiments $S0 - S5$ in 1024×512 and 512×256 resolutions. Dotted lines shows the FID scores of the models trained with ADA. (b) Bin-wise FIDs of long-tailed and uniform distribution in terms of facial yaw angle along with the number of training images (c) Bin-wise texture FIDs of long-tailed and uniform distribution along with the number of training images.

converge to a certain extent.

The dotted lines indicate the results of ADA experiments with subsets $S0 - S2$. The employed data augmentation strategy helps to reduce FID when training data is less than $40K$. Table 2 and Figure 12 in Appendix B show the detailed quantitative results of FID and IS scores, from where the IS increases with the size of training data and slows down after the amount of data reaches $40K$.

Discussion. The experiments confirm that ADA can improve the generation quality for datasets smaller than $40K$ images, in terms of FID and IS. However, ADA still cannot fully compensate for the impact of insufficient data. Besides, when the amount of data is less than $40K$, the relationship between image quality and data size is close to linear. As the amount of data increases to $40K$ and more, the improvement in the quality of the resulting images slows down and is less significant.

4.2. Data Distribution

Motivation. The nature of GAN makes the model inherits the distribution of the training dataset and introduces generation bias due to dataset imbalance [49]. This bias severely affects the performance of GAN models. To address this issue, studies for unfairness mitigation [15, 25, 55, 75, 96] have attracted substantial research interest. In this work, we explore the question of data distribution in human generation and conduct experiments to verify whether a uniform data distribution can improve the performance of a human generation model.

Experimental Settings. This study decomposes the distribution of the human body into Face Orientation and Clothing Texture, since face fidelity has a significant impact on visual perception and clothing occupies a large portion of the full-body image. Figure 5 depicts the distribution of face orientation angle in SHHQ. The general features of human faces are relatively symmetrical; thus, we fold yaw distribution vertically along 0° and get the long-tailed distribution. For the face and clothing experiments, we collect

an equal number of long-tailed and uniformly distributed datasets from SHHQ for face rotation angle and upper-body clothing texture, respectively.

Results. To evaluate the image quality in terms of different distributions, the cropped faces and clothing regions are used to calculate FID, and FID is calculated separately for each bin. Result can be found in Figure 4 (b) and (c).

1) Face Orientation: As for the long-tailed experiment (blue curve in Figure 4 (b)), the FID progressively grows as the face yaw angle increases and remains high when the facial rotation angle is too large. By contrast, the upward trend for the face FID in the uniform experiment (red) is more gradual. In addition, the amount of the training data of the first two bins in the uniform set is greatly reduced compared to the long-tail experiment, but the damage to FID is slight.

Figure 5 presents the random samples belonging to different bins in face experiments. It can be visually observed that the right-most samples in the uniform experiment have better image quality. More cropped faces in different bins from both experiments can be found in Appendix B.

2) Clothing Texture: From Figure 4 (c), except for the first bin (“plain” pattern), the FID curve climbs steadily as the amount of training data for the long-tailed experiment decreases, and the FID curve for the uniform experiment also shows a near-uniform pattern. In particular, FID of the last bin for the uniform experiment is lower than that in the long-tailed setting. We infer that the training samples for “plaid” clothing texture in the long-tailed experiment are too few to be learned by the model.

As for the “plain” bin results, the long-tailed distribution has a lower FID score in this bin. The reason may lie in that the number of plain textures in the long-tailed distribution is considerably higher than that in the uniform distribution. Also, it can be observed that the training patches in this bin are mainly textureless color blocks (see Appendix B), where such patterns may be easier to capture by models.

Discussion. Based on the above analysis, we conclude that

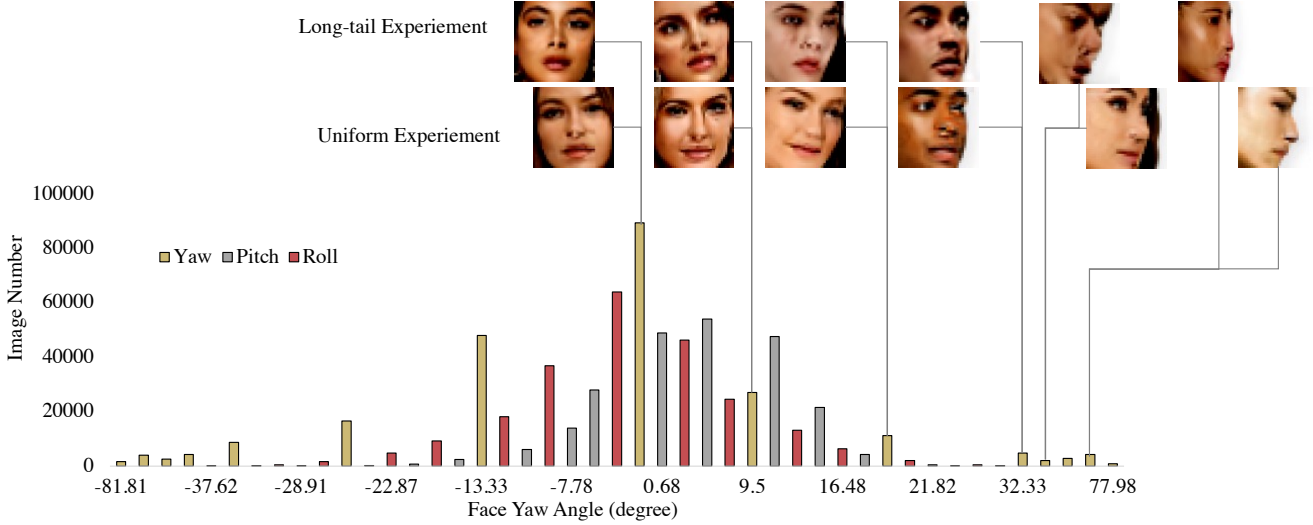


Figure 5. **Distribution of face orientation with face visualization.** The distribution of head rotation angles along the yaw/pitch/roll dimensions in the proposed SHHQ dataset, where yaw represents head rotation about the vertical axis. 0° in yaw represents the face facing straight ahead. The face patches are cropped from the *generated* images at 1024×512 resolution (see Section 4.2). Face patches in the first row are randomly sampled from the long-tailed experiment, and the second row is for the uniform experiment. Face yaw angle increases from left to right.

the uniform distribution of face rotation angles can effectively reduce the FID of rare training faces while maintaining acceptable image quality for the dominant faces. However, simply balancing the distribution of texture patterns does not always reduce the corresponding FID effectively. This phenomenon raises an interesting question that can be further explored: is the relation between image quality and data distribution also entangled with other factors, *e.g.*, image pattern and data size? Additionally, due to the nature of GAN-based structures, a GAN model memorizes the entire dataset, and usually, the discriminator tends to overfit those poorly sampled images at the tail of the distribution. Consequently, the long-tailed situation accumulated as “tail” images is barely generated. From this perspective, it also can be seen that the uniform distribution preserves the diversity of faces/textures and partially alleviates this problem.

4.3. Data Alignment

Motivation. Recently, researchers have drawn attention to spatial bias in generation tasks. Several works [39, 41] align face images with keypoints for face generation, and other studies propose different alignment schemes to preprocess non-rigid objects [6, 14, 32, 50, 76]. In this paper, we study the relationship between the spatial deviation of the entire human body and the generated image quality.

Experimental Settings. We randomly sample a set of $50K$ images from the SHHQ dataset and align every image separately using three different alignment strategies: aligning the image based on the face center, pelvis, and the midpoint of the whole body, as shown in Figure 6.

Following are the reasons for selecting these three posi-

tions as alignment centers. 1) For the face center, we hypothesize that faces contain rich semantic information that is valuable for learning and may account for a heavy proportion in human generation. 2) For the pelvis, studies related to human pose estimation [56, 60, 65, 88] conventionally predict the body joint coordinates relative to the pelvis. Thus we employ the pelvis as the alignment anchor. 3) For the body’s midpoint, the leg-to-body ratio (the proportion of upper and lower body length) may vary among different people; therefore, we try to find the mean coordinates of the full body with the help of the segmentation mask.

Results. Human images are complex and easily affected by various extrinsic factors such as body poses and camera viewpoints. The FID scores for the face-aligned, pelvis-aligned, and mid-body-aligned experiments are 3.5, 2.8, and 2.4, respectively. Figure 6 further interprets this perspective as the human bodies in (b) and (c) are tilted, and the overall image quality is degraded. The example shown in Figure 6 (c) also presents the inconsistent human positions caused by different leg-to-body ratios.

Discussion. Both FID scores and visualizations suggest that the human generative models gain more stable spatial semantic information through the mid-body alignment method than face- and pelvis-centered methods. We believe this observation could benefit later studies on human generation.

4.4. Experimental Insights

Now the questions can be answered based on the above investigations:

For **Question-1** (Data Size): A large dataset with more than

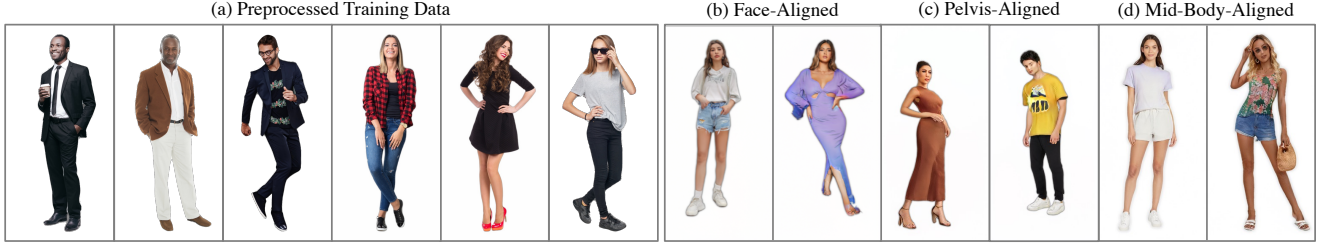


Figure 6. **Example of Preprocessed Data and Different Alignment Schemes.** Part (a) illustrates processed training data with the consideration of resolution, body position, body-part occlusion, human pose, multi-person and background. (b) - (d) display random sample results from baseline models with three different alignment strategies. (b) Face-aligned: aligned with averaged coordinate from a face bounding box. (c) Pelvis-aligned: the keypoint [7] of the pelvis is used as the alignment center. (d) Mid-body-aligned: body center is defined by averaging coordinates of extreme points (topmost, bottommost, leftmost, and rightmost points) on the boundaries of the segmentation mask.

40K images helps to train a high-fidelity unconditional human generation model, for both 512×256 and 1024×512 resolution.

For **Question-2** (Data Distribution): The uniform distribution of face rotation angles helps reduce the FID of rare faces while maintaining a reasonable quality of dominant faces. But simply balancing the clothing texture distribution does not effectively improve the generation quality.

For **Question-3** (Data Alignment): Aligning the human by the center of the full body presents a quality improvement over aligning the human by face or pelvis centers.

5. Model Zoo and Editing Benchmark

5.1. Model Zoo

In the field of face generation, a pre-trained StyleGAN [39] model has shown remarkable potential and success in various downstream tasks, including editing [2, 94], neural rendering [26], and super-resolution [12, 57], which have spawned a series of compelling research. Nevertheless, a publicly available pre-trained model is still lacking for the human generation task. To fill this gap, we train our baseline model on the collected 230K images (SHHQ) using the StyleGAN2 [40] framework. The training takes about 5 days on 8 Tesla V100 GPUs, and the model provides the best FID of 1.57. As seen in Figures 1 and 13, our model has the ability to generate full-body images with diverse poses and clothing textures under satisfactory image quality. The other two models are fully trained with StyleGAN [39] and StyleGAN3 [38], both in a resolution of 1024×512 . To adapt various application scenarios, we train models with different StyleGAN architectures [38–40] in a lower image resolution (512×256) as well. In total, all the 6 models will be released for future research. We believe they can contribute to the exploration of various tasks related to human generation and continuously benefit the community.

Furthermore, the style mixing results of the baseline model show the interpretability of the corresponding latent space, which can be seen in Figure 7. As seen in the figure,

source and reference images are sampled from the baseline model, and the rest images are the style-mixing results. We see that copying low layers from reference images to source images brings changes in geometry features (pose) from reference to the source. In contrast, other features such as skin colors, garment colors, and personal identities in source images are preserved. When copying middle styles, clothing type and identical appearance are copied from reference to the source. Finally, we observe that fine styles from high-resolution layers control the clothing color. More examples are displayed in Appendix E.1. These style mixing results suggest that the provided model’s geometry and appearance information are well disentangled.

5.2. Editing Benchmark

StyleGAN has presented remarkable editing capabilities over faces. In this section, we extend it to the full-scale human by using off-the-shelf inversion and editing methods, in which we validate the potential of our proposed model zoo. In addition, we re-implement the concurrent human generation method, InsetGAN [18], to further demonstrate another practical usage with the provided model zoo.

First, we leverage several SOTA StyleGAN-based facial editing techniques, such as InterFaceGAN [79], StyleSpace [94], and SeFa [80], with multiple editing directions: garment length for tops and bottoms, and global pose orientation. To examine the ability of editing real images with the provided model, PTI [71] is adopted to invert images before editing.

As illustrated in Figure 8, PTI presents the ability to invert real full-body human images. For attributes manipulation, StyleSpace [94] expresses better disentanglement compared to InterFaceGAN [79] and SeFa [80], as only the attribute-targeted region has been changed. However, as for the regions to be edited, the results of InterFaceGAN [79] are more natural and photo-realistic. It turns out that the latent space of the human body is more complicated than other domains such as faces, objects, and scenes, and more attention should be paid to disentangle human attributes. More editing results are shown in Appendix E.2.



Figure 7. **Style-mixing Results.** The sets of images are randomly sampled from the provided baseline model, with latent codes recorded. Those images are treated as reference and source images. The rest of the images are generated by style-mixing: borrowing low/mid/high layers in the corresponding reference images’ latent codes and combining them with the rest layers of latent code in source images. Low/mid/high corresponds to coarse/middle/fine spatial resolution. It can be seen that low layers in the latent code control coarse features such as poses, middle layers are related to clothing type, identical appearance, and higher layers convey fine-grained features, for example, clothing color.

Moreover, InsetGAN [18] proposes a multi-GAN optimization method to fuse face and body generated from separate GAN models. We re-implement this process by iteratively optimizing the latent codes for random faces and bodies generated by the FFHQ [40] and our baseline model, respectively. In Figure 9, we show the fused full-body images of six human postures with different male and female faces. The optimization procedure blends diverse faces and bodies in a graceful manner. Our baseline model can be adapted to different faces and generate more complicated and diverse full-body images.

In this section, we adopt a representative facial editing method on humans generated by our pre-trained model and obtain impressive results. We also show that images generated from the released human model can be further locally

optimized using existing GAN models. All these works demonstrate the effectiveness and convenience of our provided model zoo and verify its potential in human-centric tasks.

6. Future Work

In this study, we take a preliminary step towards the exploration of the human generation/editing problem. We believe many future works can be further explored based on the SHHQ dataset and the provided model zoo. In the following, we discuss three interesting directions, *i.e.*, Human Generation/Editing, Neural Rendering, and Multi-modal Generation.

Human Generation / Editing. Studies in unconditional



Figure 8. **Image editing results with different methods.** Image editing techniques are tested in both real images with PTI inversion (first row) and generated images from our baseline model (second row). The reference image is on the left, the second column shows the inverted result, and the remaining shows the results of modifying sleeve length and dress length. We shorten the length of upper and lower garments for the real image, and increase the upper/lower garment length for the generated image.

human generation [18, 100], human editing [3, 24, 72], virtual try-on [16, 44, 59, 89], and motion transfer [9, 47] heavily rely on large datasets to train or use existing pre-trained models as the first step of transfer learning. Furthermore, editing benchmarks show that disentangled editing of the human body remains challenging for existing methods [79, 94]. In this context, the released model zoo could expedite such research progress. Additionally, we further analyze failure cases generated by the provided model and discuss corresponding potential efforts that could be made to human generation tasks in Appendix D.

Neural Rendering. Another future research direction is to improve 3D consistency and mitigate artifacts in full-body human generation through neural rendering [10, 11, 26, 61, 62, 77]. Similar to work such as EG3D [10], StyleNeRF [26], and StyleSDF [62], we encourage researchers to use our human models to facilitate human generation with multi-view consistency.

Multi-modal Generation. Cross-modal representation is an emerging research trend, such as CLIP [68] and Imagebert [67]. Hundreds of studies are made on text-driven image generation and manipulation [33, 42, 63, 67, 70, 90, 95, 98], *e.g.*, DALLE [70] and AttnGAN [98]. In the meantime, several studies show interest in probing the transfer learning benefits of large-scale pre-trained models [21, 68, 93]. Most of these works focus on faces and objects, whereas research

fields related to full-scale humans could be explored more, for example, text-to-human generation and text-driven human attributes manipulation, with the help of the provided full-body human models.

7. Conclusion

This work mainly probes how to train unconditional human-based GAN models to generate photo-realistic images from a data-centric perspective. By leveraging the 230K SHHQ dataset, we analyze three fundamental yet critical issues that the community cares most about: data size, data distribution, and data alignment. While experimenting with StyleGAN and large-scale data, we obtain several empirical insights. Apart from these, we create a model zoo, consisting of six human-GAN models, and the effectiveness of the model zoo is demonstrated by employing several state-of-the-art face editing methods.

Acknowledgements. We thank Hao Zhu, Zhaoyang Liu and Zhuoqian Yang for their feedback and discussions. This study is partly supported by NTU NAP, MOE AcRF Tier 1 (2021-T1-001-088), and under the RIE2020 Industry Alignment Fund Industry Collaboration Projects (IAF-ICP) Funding Initiative, as well as cash and in-kind contribution from the industry partner(s).



Figure 9. **InsetGAN Results.** We show the combined results of six different human bodies generated from the given baseline model and six faces generated from the FFHQ [40] model. For every single face, shown on the left corner of each grid, we jointly optimize it with three different bodies.

References

- [1] Rameen Abdal, Yipeng Qin, and Peter Wonka. Image2StyleGAN++: How to edit the embedded images? In *CVPR*, 2020. 3
- [2] Rameen Abdal, Peihao Zhu, Niloy J Mitra, and Peter Wonka. StyleFlow: Attribute-conditioned exploration of StyleGAN-generated images using conditional continuous normalizing flows. *ACM TOG*, 2021. 2, 3, 8
- [3] Badour Albahar, Jingwan Lu, Jimei Yang, Zhixin Shu, Eli Shechtman, and Jia-Bin Huang. Pose with style: Detail-preserving pose-guided image synthesis with conditional StyleGAN. *ACM TOG*, 2021. 3, 10
- [4] Martin Arjovsky and Léon Bottou. Towards principled methods for training generative adversarial networks. *arXiv preprint*, arXiv:1701.04862, 2017. 3
- [5] Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein generative adversarial networks. In *ICML*, 2017. 2, 3
- [6] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale GAN training for high fidelity natural image synthesis. In *ICLR*, 2019. 2, 3, 5, 7
- [7] Zhe Cao, Tomas Simon, Shih-En Wei, and Yaser Sheikh. Realtime multi-person 2D pose estimation using part affinity fields. In *CVPR*, 2017. 4, 8
- [8] Lucy Chai, Jun-Yan Zhu, Eli Shechtman, Phillip Isola, and Richard Zhang. Ensembling with deep generative views. In *CVPR*, 2021. 3
- [9] Caroline Chan, Shiry Ginosar, Tinghui Zhou, and Alexei A Efros. Everybody dance now. In *ICCV*, 2019. 2, 10
- [10] Eric R Chan, Connor Z Lin, Matthew A Chan, Koki Nagano, Boxiao Pan, Shalini De Mello, Orazio Gallo, Leonidas Guibas, Jonathan Tremblay, Sameh Khamis, et al. Efficient geometry-aware 3D generative adversarial networks. *arXiv preprint*, arXiv:2112.07945, 2021. 10
- [11] Eric R Chan, Marco Monteiro, Petr Kellnhofer, Jiajun Wu, and Gordon Wetzstein. Pi-GAN: Periodic implicit generative adversarial networks for 3D-aware image synthesis. In *CVPR*, 2021. 10
- [12] Kelvin CK Chan, Xintao Wang, Xiangyu Xu, Jinwei Gu, and Chen Change Loy. GLEAN: Generative latent bank for large-factor image super-resolution. In *CVPR*, 2021. 8
- [13] MMSegmentation Contributors. MMSegmentation: OpenMMLab semantic segmentation toolbox and benchmark. <https://github.com/open-mmlab/mms Segmentation>, 2020. 4
- [14] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat GANs on image synthesis. *NeurIPS*, 2021. 2, 7
- [15] Nikolaos Dionelis, Mehrdad Yaghoobi, and Sotirios A Tsaftaris. Tail of distribution GAN (TailGAN): GenerativeAdversarial-network-based boundary formation. In *SSPD*, 2020. 2, 6
- [16] Haoye Dong, Xiaodan Liang, Xiaohui Shen, Bochao Wang, Hanjiang Lai, Jia Zhu, Zhiting Hu, and Jian Yin. Towards multi-pose guided virtual Try-On network. In *ICCV*, 2019. 2, 10

- [17] Patrick Esser, Ekaterina Sutter, and Björn Ommer. A variational U-Net for conditional appearance and shape generation. In *CVPR*, 2018. 3
- [18] Anna Frühstück, Krishna Kumar Singh, Eli Shechtman, Niloy J Mitra, Peter Wonka, and Jingwan Lu. Inset-GAN for full-body image generation. *arXiv preprint*, arXiv:2203.07293, 2022. 2, 8, 9, 10
- [19] Raghudeep Gadde, Qianli Feng, and Aleix M Martinez. Detail me more: Improving GAN’s photo-realism of complex scenes. In *ICCV*, 2021. 15
- [20] Andrew Gahan. *3ds Max Modeling for Games: Insider’s guide to game character, vehicle, and environment modeling*. 2012. 2
- [21] Deepti Ghadiyaram, Du Tran, and Dhruv Mahajan. Large-scale weakly-supervised pre-training for video action recognition. In *CVPR*, 2019. 10
- [22] Ke Gong, Xiaodan Liang, Dongyu Zhang, Xiaohui Shen, and Liang Lin. Look into person: Self-supervised structure-sensitive learning and a new benchmark for human parsing. In *CVPR*, 2017. 3, 4
- [23] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *NeurIPS*, 2014. 2
- [24] Artur Grigorev, Karim Isakov, Anastasia Ianina, Renat Bashirov, Ilya Zakharkin, Alexander Vakhitov, and Victor Lempitsky. StylePeople: A generative model of fullbody human avatars. In *CVPR*, 2021. 3, 10
- [25] Aditya Grover, Jiaming Song, Ashish Kapoor, Kenneth Tran, Alekh Agarwal, Eric J Horvitz, and Stefano Ermon. Bias correction of learned generative models using likelihood-free importance weighting. *NeurIPS*, 2019. 2, 6
- [26] Jiatao Gu, Lingjie Liu, Peng Wang, and Christian Theobalt. StyleNERF: A style-based 3D-aware generator for high-resolution image synthesis. *arXiv preprint*, arXiv:2110.08985, 2021. 8, 10
- [27] Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron C Courville. Improved training of wasserstein GANs. In *NeurIPS*, 2017. 2, 3
- [28] Xintong Han, Zuxuan Wu, Zhe Wu, Ruichi Yu, and Larry S Davis. VITON: An image-based virtual Try-On network. In *CVPR*, 2018. 3, 4
- [29] Erik Härkönen, Aaron Hertzmann, Jaakko Lehtinen, and Sylvain Paris. GanSpace: Discovering interpretable GAN controls. *NeurIPS*, 2020. 3
- [30] Rui Huang, Shu Zhang, Tianyu Li, and Ran He. Beyond face rotation: Global and local perception GAN for photo-realistic and identity preserving frontal view synthesis. In *ICCV*, 2017. 15
- [31] Liming Jiang, Bo Dai, Wayne Wu, and Chen Change Loy. Deceive D: Adaptive Pseudo Augmentation for GAN training with limited data. In *NeurIPS*, 2021. 2
- [32] Yuming Jiang, Kelvin CK Chan, Xintao Wang, Chen Change Loy, and Ziwei Liu. Robust reference-based super-resolution via C2-Matching. In *CVPR*, 2021. 7
- [33] Yuming Jiang, Ziqi Huang, Xingang Pan, Chen Change Loy, and Ziwei Liu. Talk-to-edit: Fine-grained facial editing via dialog. In *ICCV*, 2021. 3, 10
- [34] Nebojsa Jojic, Jin Gu, TS Shen, and Thomas S Huang. Computer modeling, analysis, and synthesis of dressed humans. *TCSVT*, 1999. 2
- [35] Wang-Cheng Kang, Chen Fang, Zhaowen Wang, and Julian McAuley. Visually-aware fashion recommendation and design with generative image models. In *ICDM*, 2017. 2
- [36] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of GANs for improved quality, stability, and variation. *ICLR*, 2017. 2, 3
- [37] Tero Karras, Miika Aittala, Janne Hellsten, Samuli Laine, Jaakko Lehtinen, and Timo Aila. Training generative adversarial networks with limited data. *NeurIPS*, 2020. 2, 5
- [38] Tero Karras, Miika Aittala, Samuli Laine, Erik Härkönen, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Alias-free generative adversarial networks. *NeurIPS*, 2021. 2, 3, 8
- [39] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *CVPR*, 2019. 2, 3, 4, 7, 8, 15
- [40] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of StyleGAN. In *CVPR*, 2020. 2, 3, 8, 9, 11
- [41] Vahid Kazemi and Josephine Sullivan. One millisecond face alignment with an ensemble of regression trees. In *CVPR*, 2014. 2, 7
- [42] Umut Kocasari, Alara Dirik, Mert Tiftikci, and Pinar Yarnardag. StyleMC: Multi-channel based fast text-guided image generation and manipulation. In *WACV*, 2022. 10
- [43] Chenyi Lei, Dong Liu, Weiping Li, Zheng-Jun Zha, and Houqiang Li. Comparative deep learning of hybrid representations for image recommendations. In *CVPR*, 2016. 2
- [44] Kathleen M Lewis, Srivatsan Varadharajan, and Ira Kemelmacher-Shlizerman. TryOnGAN: Body-aware Try-On via layered interpolation. *ACM TOG*, 2021. 3, 10
- [45] Zhong Li, Lele Chen, Celong Liu, Fuyao Zhang, Zekun Li, Yu Gao, Yuanzhou Ha, Chenliang Xu, Shuxue Quan, and Yi Xu. Animated 3D human avatars from a single image with GAN-based texture inference. *CNG*, 2021. 2
- [46] Xiaodan Liang, Si Liu, Xiaohui Shen, Jianchao Yang, Luoqi Liu, Jian Dong, Liang Lin, and Shuicheng Yan. Deep human parsing with active template regression. *PAMI*, 2015. 3, 4
- [47] Wen Liu, Zhixin Piao, Jie Min, Wenhan Luo, Lin Ma, and Shenghua Gao. Liquid warping GAN: A unified framework for human motion imitation, appearance transfer and novel view synthesis. In *ICCV*, 2019. 2, 10
- [48] Ziwei Liu, Ping Luo, Shi Qiu, Xiaogang Wang, and Xiaoou Tang. DeepFashion: Powering robust clothes recognition and retrieval with rich annotations. In *CVPR*, 2016. 3, 4, 5
- [49] Ziwei Liu, Zhongqi Miao, Xiaohang Zhan, Jiayun Wang, Boqing Gong, and Stella X Yu. Large-scale long-tailed recognition in an open world. In *CVPR*, 2019. 2, 6
- [50] Ziwei Liu, Sijie Yan, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Fashion landmark detection in the wild. In *ECCV*, 2016. 7
- [51] Liqian Ma, Xu Jia, Qianru Sun, Bernt Schiele, Tinne Tuytelaars, and Luc Van Gool. Pose guided person image generation. *NeurIPS*, 2017. 3
- [52] Liqian Ma, Qianru Sun, Stamatios Georgioulis, Luc Van Gool, Bernt Schiele, and Mario Fritz. Disentangled person image generation. In *CVPR*, 2018. 3

- [53] Qianli Ma, Jinlong Yang, Anurag Ranjan, Sergi Pujades, Gerard Pons-Moll, Siyu Tang, and Michael J Black. Learning to dress 3D people in generative clothing. In *CVPR*, 2020. 2
- [54] Xudong Mao, Qing Li, Haoran Xie, Raymond Y.K. Lau, Zhen Wang, and Stephen Paul Smolley. Least squares generative adversarial networks. In *ICCV*, 2017. 2
- [55] Giovanni Mariani, Florian Scheidegger, Roxana Istrate, Costas Bekas, and Cristiano Malossi. BaGAN: Data augmentation with balancing GAN. *arXiv preprint*, arXiv:1803.09655, 2018. 2, 5, 6
- [56] Julieta Martinez, Rayat Hossain, Javier Romero, and James J Little. A simple yet effective baseline for 3D human pose estimation. In *ICCV*, 2017. 7
- [57] Sachit Menon, Alexandru Damian, Shijia Hu, Nikhil Ravi, and Cynthia Rudin. PULSE: Self-supervised photo upsampling via latent space exploration of generative models. In *CVPR*, 2020. 8
- [58] Takeru Miyato, Toshiaki Kataoka, Masanori Koyama, and Yuichi Yoshida. Spectral normalization for generative adversarial networks. *ICLR*, 2018. 3
- [59] Assaf Neuberger, Eran Borenstein, Bar Hilleli, Eduard Oks, and Sharon Alpert. Image based virtual Try-On network from unpaired data. In *CVPR*, 2020. 2, 10
- [60] Xuecheng Nie, Jiashi Feng, Jianfeng Zhang, and Shuicheng Yan. Single-stage multi-person pose machines. In *ICCV*, 2019. 7
- [61] Michael Niemeyer and Andreas Geiger. GIRAFFE: Representing scenes as compositional generative neural feature fields. In *CVPR*, 2021. 10
- [62] Roy Or-El, Xuan Luo, Mengyi Shan, Eli Shechtman, Jeong Joon Park, and Ira Kemelmacher-Shlizerman. StyleSDF: High-resolution 3D-consistent image and geometry generation. *arXiv preprint*, arXiv:2112.11427, 2021. 10
- [63] Or Patashnik, Zongze Wu, Eli Shechtman, Daniel Cohen-Or, and Dani Lischinski. StyleCLIP: Text-driven manipulation of StyleGAN imagery. In *ICCV*, 2021. 3, 10
- [64] Chaitanya Patel, Zhouyingcheng Liao, and Gerard Pons-Moll. TailorNet: Predicting clothing in 3D as a function of human pose, shape and garment style. In *CVPR*, 2020. 2
- [65] Georgios Pavlakos, Xiaowei Zhou, Konstantinos G Derpanis, and Kostas Daniilidis. Coarse-to-fine volumetric prediction for single-image 3D human pose. In *CVPR*, 2017. 7
- [66] Albert Pumarola, Jordi Sanchez-Riera, Gary Choi, Alberto Sanfeliu, and Francesc Moreno-Noguer. 3DPeople: Modeling the geometry of dressed humans. In *CVPR*, 2019. 2
- [67] Di Qi, Lin Su, Jia Song, Edward Cui, Taroon Bharti, and Arun Sacheti. ImageBERT: Cross-modal pre-training with large-scale weak-supervised image-text data. *arXiv preprint*, arXiv:2001.07966, 2020. 10
- [68] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *ICML*, 2021. 10
- [69] Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint*, arXiv:1511.06434, 2015. 2
- [70] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *ICML*, 2021. 10
- [71] Daniel Roich, Ron Mokady, Amit H Bermano, and Daniel Cohen-Or. Pivotal tuning for latent-based editing of real images. *arXiv preprint*, arXiv:2106.05744, 2021. 2, 3, 8
- [72] Kripasindhu Sarkar, Vladislav Golyanik, Lingjie Liu, and Christian Theobalt. Style and pose control for image synthesis of humans from a single monocular view. *arXiv preprint*, arXiv:2102.11263, 2021. 3, 10
- [73] Kripasindhu Sarkar, Lingjie Liu, Vladislav Golyanik, and Christian Theobalt. HumanGAN: A generative model of human images. In *3DV*, 2021. 3
- [74] Kripasindhu Sarkar, Dushyant Mehta, Weipeng Xu, Vladislav Golyanik, and Christian Theobalt. Neural re-rendering of humans from a single image. In *ECCV*, 2020. 2
- [75] Prasanna Sattigeri, Samuel C Hoffman, Vijil Chenthamarakshan, and Kush R Varshney. Fairness GAN: Generating datasets with fairness properties using a generative adversarial network. *IBM JRD*, 2019. 2, 6
- [76] Axel Sauer, Katja Schwarz, and Andreas Geiger. StyleGAN-XL: Scaling StyleGAN to large diverse datasets. *arXiv preprint*, arXiv:2202.00273, 2022. 2, 7
- [77] Katja Schwarz, Yiyi Liao, Michael Niemeyer, and Andreas Geiger. GRAF: Generative radiance fields for 3D-aware image synthesis. *NeurIPS*, 2020. 10
- [78] Sefik Ilkin Serengil and Alper Ozpinar. HyperExtended LightFace: A facial attribute analysis framework. In *ICEET*, 2021. 5
- [79] Yujun Shen, Ceyuan Yang, Xiaoou Tang, and Bolei Zhou. InterFaceGAN: Interpreting the disentangled face representation learned by GANs. *PAMI*, 2020. 2, 3, 8, 10
- [80] Yujun Shen and Bolei Zhou. Closed-form factorization of latent semantics in GANs. In *CVPR*, 2021. 2, 3, 8
- [81] Aliaksandr Siarohin, Enver Sangineto, Stéphane Lathuilière, and Nicu Sebe. Deformable GANs for pose-based human image generation. In *CVPR*, 2018. 3
- [82] Dan Song, Ruofeng Tong, Jian Chang, Xiaosong Yang, Min Tang, and Jian Jun Zhang. 3D body shapes estimation from dressed-human silhouettes. In *CGF*, 2016. 2
- [83] Ayush Tewari, Mohamed Elgharib, Florian Bernard, Hans-Peter Seidel, Patrick Pérez, Michael Zollhöfer, and Christian Theobalt. PIE: Portrait image embedding for semantic control. *ACM TOG*, 2020. 3
- [84] Ayush Tewari, Mohamed Elgharib, Gaurav Bharaj, Florian Bernard, Hans-Peter Seidel, Patrick Perez, Michael Zollhofer, and Christian Theobalt. StyleRig: Rigging StyleGAN for 3D control over portrait images. In *CVPR*, 2020. 2
- [85] Jamal Toutouh, Erik Hemberg, and Una-May O'Reilly. Data dieting in GAN training. In *Deep Neural Evolution*. 2020. 2, 5
- [86] Omer Tov, Yuval Alaluf, Yotam Nitzan, Or Patashnik, and Daniel Cohen-Or. Designing an encoder for StyleGAN image manipulation. *ACM TOG*, 2021. 3
- [87] Christos Tzelepis, Georgios Tzimiropoulos, and Ioannis Patras. WarpedGANSpace: Finding non-linear RBF paths in GAN latent space. In *ICCV*, 2021. 3

- [88] Márton Végés and András Lőrincz. Absolute human pose estimation with depth prediction network. In *IJCNN*, 2019. 7
- [89] Bochao Wang, Huabin Zheng, Xiaodan Liang, Yimin Chen, Liang Lin, and Meng Yang. Toward characteristic-preserving image-based virtual Try-On network. In *ECCV*, 2018. 2, 10
- [90] Tianren Wang, Teng Zhang, and Brian Lovell. Faces a la carte: Text-to-face generation via attribute disentanglement. In *WACV*, 2021. 10
- [91] Tengfei Wang, Yong Zhang, Yanbo Fan, Jue Wang, and Qifeng Chen. High-fidelity GAN inversion for image attribute editing. *arXiv preprint*, arXiv:2109.06590, 2021. 3
- [92] Chunpeng Wu and Hai Li. Conditional transferring features: Scaling GANs to thousands of classes with 30% less high-quality data for training. In *IJCNN*, 2020. 5
- [93] Qingyang Wu, Lei Li, and Zhou Yu. TextGAIL: Generative adversarial imitation learning for text generation. *arXiv preprint*, arXiv:2004.13796, 2020. 10
- [94] Zongze Wu, Dani Lischinski, and Eli Shechtman. StyleSpace analysis: Disentangled controls for StyleGAN image generation. In *CVPR*, 2021. 2, 3, 8, 10, 15
- [95] Weihao Xia, Yujiu Yang, Jing-Hao Xue, and Baoyuan Wu. TediGAN: Text-guided diverse face image generation and manipulation. In *CVPR*, 2021. 10
- [96] Depeng Xu, Shuhan Yuan, Lu Zhang, and Xintao Wu. FairGAN: Fairness-aware generative adversarial networks. In *IEEE BigData*, 2018. 2, 6
- [97] Hongyi Xu, Eduard Gabriel Bazavan, Andrei Zanfir, William T Freeman, Rahul Sukthankar, and Cristian Sminchisescu. Ghum & ghuml: Generative 3D human shape and articulated pose models. In *CVPR*, 2020. 2
- [98] Tao Xu, Pengchuan Zhang, Qiuyuan Huang, Han Zhang, Zhe Gan, Xiaolei Huang, and Xiaodong He. AttnGAN: Fine-grained text to image generation with attentional generative adversarial networks. In *CVPR*, 2018. 10
- [99] Yanbo Xu, Yueqin Yin, Liming Jiang, Qianyi Wu, Chengyao Zheng, Chen Change Loy, Bo Dai, and Wayne Wu. TransEditor: Transformer-based dual-space GAN for highly controllable facial editing. In *CVPR*, 2022. 3
- [100] Gokhan Yildirim, Nikolay Jetchev, Roland Vollgraf, and Urs Bergmann. Generating high-resolution fashion model images wearing custom outfits. In *ICCV Workshops*, 2019. 3, 10
- [101] Dan Zhang and Anna Khoreva. Progressive augmentation of GANs. *NeurIPS*, 2019. 3
- [102] Shengyu Zhao, Zhijian Liu, Ji Lin, Jun-Yan Zhu, and Song Han. Differentiable augmentation for data-efficient GAN training. In *NeurIPS*, 2020. 2, 5
- [103] Liang Zheng, Liyue Shen, Lu Tian, Shengjin Wang, Jingdong Wang, and Qi Tian. Scalable person re-identification: A benchmark. In *ICCV*, 2015. 3, 4

Appendix

A. SHHQ: StyleGAN-Human Datasets

The dataset we collected consists of 230K high-quality images of humans that vary in clothing appearance, ethnicity, and pose. Several training samples are shown in Figure 10. Please note that all these images are unprocessed. As shown in Figure 11, we also conduct qualitative comparisons with other human datasets to demonstrate the superiority of our clean, high-quality data. Besides, we display more generated human images from the baseline model trained with our SHHQ in Figure 13.

B. Experiment Results

Table 2 and Figure 12 display the results of data size experiment elaborated in Section 4.1. The results align with our expectation that increment training data will improve IS scores and reduce FID scores. Figure 15 and 16 depict the comparison between cropped faces and textures generated by the long-tail and uniform experiments. Due to privacy concerns, cropped training faces are not shown.

Table 2. **FID and IS for experiments of data size.** Quantitative comparisons at resolutions of 512×256 and 1024×512 .

	Data Size	512×256		1024×512	
		FID	IS	FID	IS
S0	10K	7.80	3.87	7.23	3.93
S1	20K	4.46	4.40	4.33	4.56
S2	40K	2.61	4.81	2.80	4.92
S3	80K	2.53	4.90	2.09	5.01
S4	160K	2.09	4.92	2.02	5.04
S5	230K	1.97	5.04	1.57	5.02

C. Training Scheme

We adopt the official NVIDIA Pytorch version of StyleGAN2-ADA as our codebase, and use the architecture of StyleGAN2. Here are several settings we use to accommodate this human generation task: (a) The input human-image has a width-to-height ratio of 1 : 2, and the input resolution in the script is changed accordingly. (b) We adopt the same eight mapping layers as the original StyleGAN [39]. (c) There is no such a pretrained model for human images, so all the experiments are trained from scratch with the corresponding subset. (d) All other training hyperparameters adopt the default values.

D. Limitations

Compared to face generation, training an unconditional human GAN is an arduous task because the semantic features of the full-body are much more complicated than a

single face. Figure 14 shows some failure cases generated by the baseline model, suggesting several directions that can be strengthened in future human generation work. Artifacts caused by entangled features of faces/hands and clothing accessories are revealed in Figures 14 (a) - (c). Case (d) exhibits three hands on a person, which indicates that the global perception of the model needs to be improved [30]. We observe inferior hand quality in rare poses such as (e) and (f). To address this, the potential work could be augmenting training with such extreme poses, changing the data distribution, or implementing independent networks (i.e. fine-grained discriminators) to enhance local details [19]. The face and texture quality in cases (b) and (g) could be enhanced by local refinement as well.

E. Visualization of the Applications

E.1. Style-Mixing

Here we provide more examples of images generated by style-mixing on our baseline model. Figure 17, 18, and 19 represent the results of style-mixing on coarse, middle, and high resolution respectively. It shows that the latent at different scales control different high-level attributes of the clothed human, which is similar to face images.

E.2. Human Editing

Figure 20 displays the rotation of the human from the front view to the back view. The editing is done in W space. Figure 21 demonstrates the editing results in the length of sleeves and bottoms, based on StyleSpace [94].



Figure 10. Examples of raw data in the training dataset.



Figure 11. Samples from different dataset with diverse resolution.

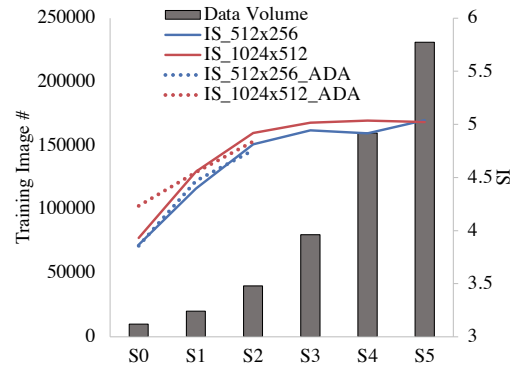


Figure 12. **IS scores.** IS scores for experiments $S_0 - S_5$ in 1024×512 and 512×256 resolutions. Dotted lines represents the IS scores with ADA strategies.



Figure 13. **Samples from our baseline model.** The model has shown the ability of generating random person with diverse clothing types, poses, genders, races, and hair types.



Figure 14. **Failure cases from baseline model.** (a) - (c): Features such as face, texture and accessories are entangled. (d) Three hands detected on a single person. (e) - (f): Inferior generated hand quality. (g): Face quality could be better.

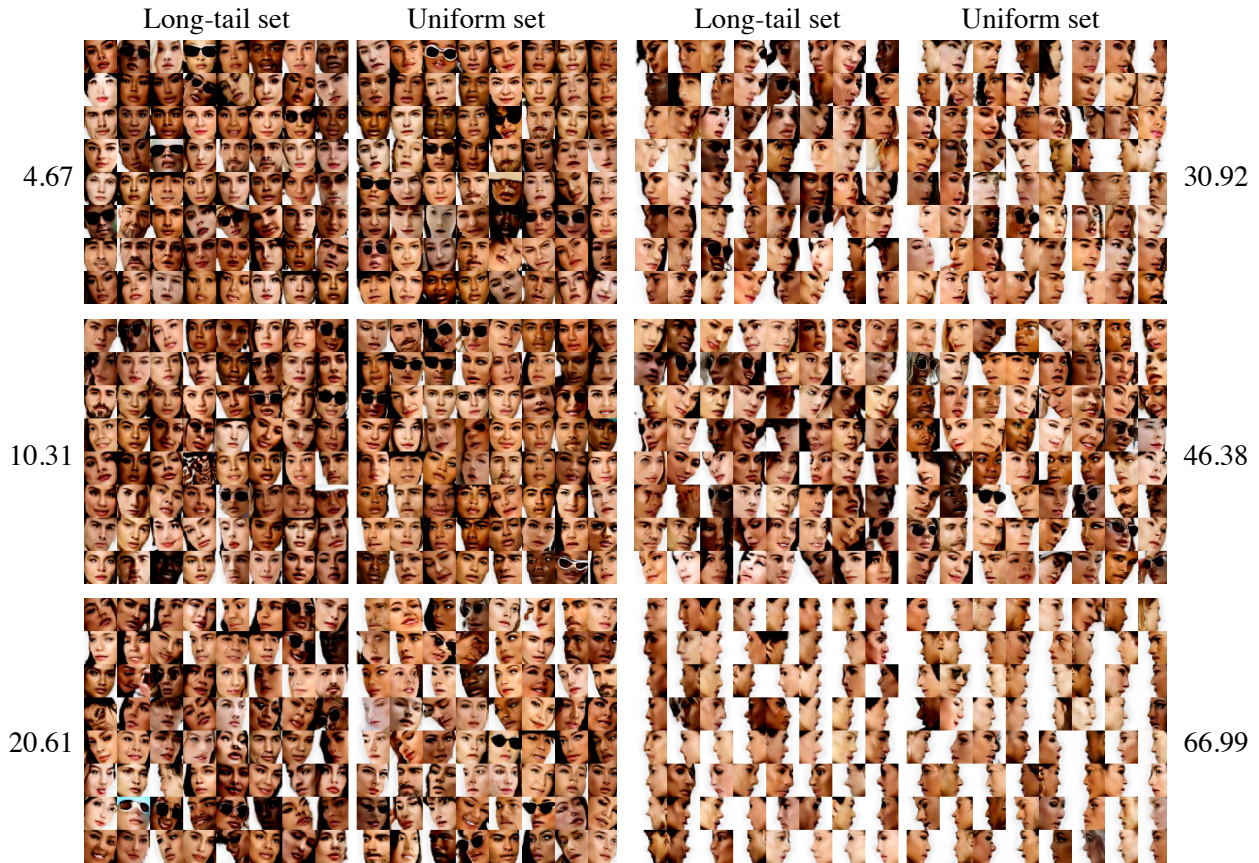


Figure 15. **Cropped faces with different face yaw angles from each bin.** All the images are generated from the long-tail and uniform experiments.



Figure 16. Random cropped texture patches from each bin for both long-tail and uniform experiments.



Figure 17. **Style-mixing** with copying styles of *coarse* resolutions from reference images (top row), and rest spatial information are used from source images (first column).



Figure 18. **Style-mixing** with copying styles of *middle* resolutions from reference images (top row), and rest spatial information are used from source images (first column).



Figure 19. **Style-mixing** with copying styles of *fine* resolutions from reference images (top row), and rest spatial information are used from source images (first column).



Figure 20. Human Editing on orientation.

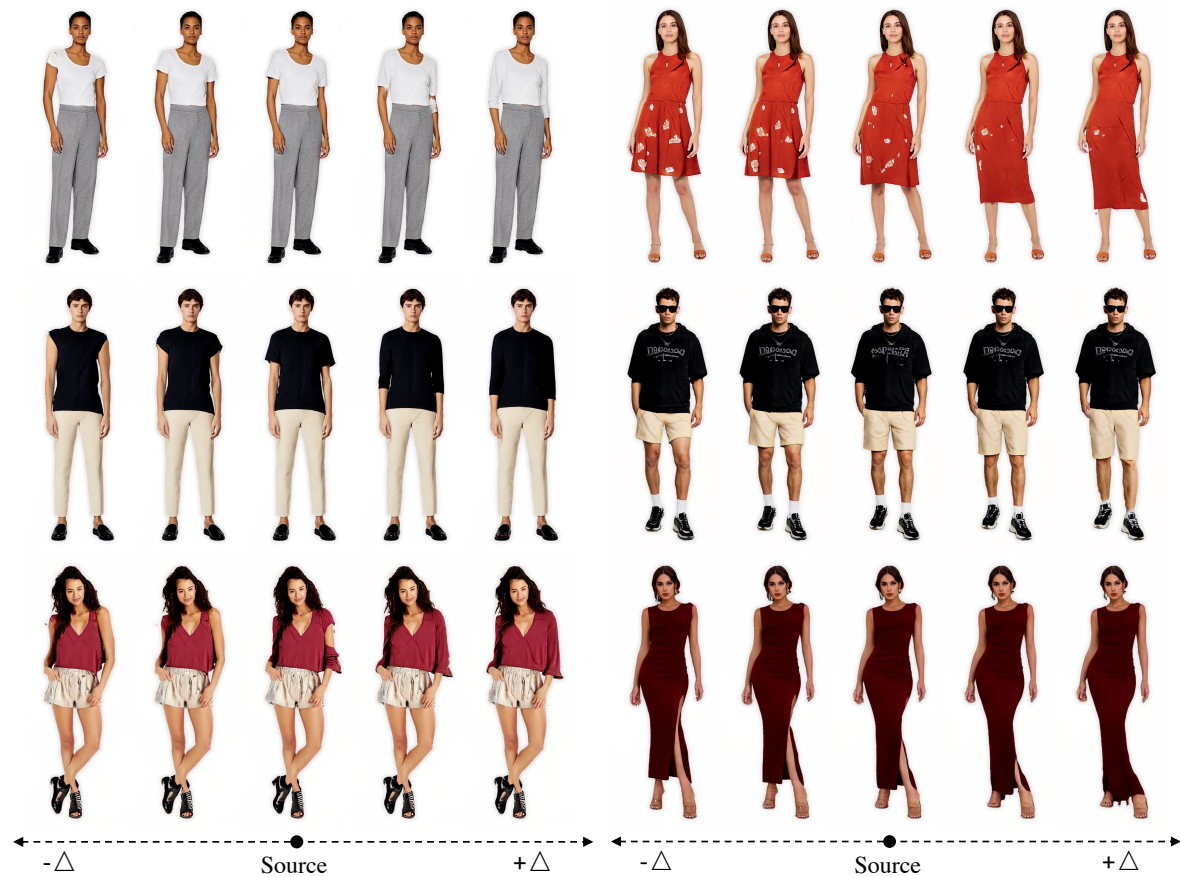


Figure 21. Human Editing on human sleeve length (left) and bottom length (right).