# MetaHead: An Engine to Create Realistic Digital Head

Dingyun Zhang[1]    Chenglai Zhong[1]    Yudong Guo[2]    Yang Hong[1]    Juyong Zhang[1]

[1]University of Science and Technology of China    [2]Image Derivative Inc

## Abstract

*Collecting and labeling training data is one important step for learning-based methods because the process is time-consuming and biased. For face analysis tasks, although some generative models [25, 10, 45, 6, 5, 11, 20] can be used to generate face data, they can only achieve a subset of generation diversity, reconstruction accuracy, 3D consistency, high-fidelity visual quality, and easy editability. One recent related work is the graphics-based generative method [48], but it can only render low realism head with high computation cost. In this paper, we propose MetaHead, a unified and full-featured controllable digital head engine, which consists of a controllable head radiance field(MetaHead-F) to super-realistically generate or reconstruct view-consistent 3D controllable digital heads and a generic top-down image generation framework LabelHead to generate digital heads consistent with the given customizable feature labels. Experiments validate that our controllable digital head engine achieves the state-of-the-art generation visual quality and reconstruction accuracy. Moreover, the generated labeled data can assist real training data and significantly surpass the labeled data generated by graphics-based methods in terms of training effect. The project page is available at: https://ustc3dv.github.io/MetaHead/.*

## 1. Introduction

Generation and reconstruction of digital human are in increasing demand in virtual reality [42], game and movie character production, and metaverse(shared virtual 3D world) [33], as digital human with controllable geometry and appearance could bring realistic visual experience. On the other hand, semantic labels of heads such as 2D/3D landmarks [55], eye gaze angle [51] and hair color are valuable and necessary precondition for many facial analysis tasks [53, 7, 46, 47], and have various applications like face alignment [23] and registration, facial expression analysis [34], video conferencing [8], human computer interface [16]. Large well-labeled training data [17] is the key to the success of learning-based models. The more precise the labels, the better the learning-based models will perform. Labeled instances however are often difficult, expensive, or time-consuming to obtain, as they require the efforts of experienced human annotators. In addition, most of the data collected by the existing head label database is under medium pose and illumination, and there are very little challenging big pose and dark-illumination data. Therefore, learning-based models still have a lot of room for improvement in challenging scenarios and generalization.

A digital head engine is an all-in-one model that can reconstruct and control digital heads, and generate digital heads consistent with the customizable head feature labels. So far, little work has been done on the unified engine [48], despite the growing promise and demand for it. Graphics-based methods [48] can generate heads consistent with specified labels. However, there exists obvious gap between the synthesized texture distribution and real texture distribution, which results that the final synthesized heads look like cartoon images.

Recently, learning-based models get more and more attentions. Among them, heads generated by early 2D-based GAN models [25, 10, 45] lack view consistency. The 3D generative adversarial structure models based on neural radiance field (NeRF) representation [11, 5] can improve the view consistency of head geometry and appearance. However, since the latent space of GAN is hard to semantically manipulate with disentanglement, they do not allow flexible head control of the output images. The model introduced with 3DMM prior [20] can reconstruct the controllable heads(it is not a generative model), but its reconstruction quality is poor on photorealism and clarity. When changing the viewing angle, the heads obtained by the above models will have hair and teeth flickering phenomenon, which destroys the visual effect. In addition, their generation or control of heads only covers medium poses and simple expressions (smile and mouth opening) due to their limited training data and attributes controlling disentanglement.

We propose MetaHead, a super-realistic controllable head engine, which realizes the reconstruction, control, and generation of heads consistent with the given labels. A naive baseline solution is to represent the 3D head scene based on NeRF, and use the decoupled 3DMM coefficients

as shape and appearance prior conditions to input NeRF, and utilize real 2D images as supervision, similar to [20]. However, this end-to-end solution is actually the decoder structure of the Auto-encoder, and the generated figures are very blurry [30]. At the same time, the 3DMM coefficients are based on the Base Face Model(BFM) [37], and the identity bases are sparse(built from very limited 3D scans). The implicit semantics of identity and expression coefficients are also not good at direct 3D control over face geometry during volume learning procedure.

In MetaHead, we propose a controllable head radiance field(MetaHead-F) by designing a strategy to combine the decoder with the pre-trained GAN generator. This end-to-end structure ensures that the disorganized latent space can be decoupled via learning. GAN does not use point-wise loss during pre-training but distribution-matched GAN loss, so the convergence point is close to the data manifold surface [3] and thus MetaHead-F can generate visually clear human heads. The generator unit inspired by Style-GAN3 [24] can suppress aliasing information, in which we designed a hierarchical structure attention module to solve the chronic problem of hair and teeth flickering when the viewing angle changes. Furthermore, we add a parallel semantic network to implicitly learn hair and mouth contours. We use face recognition features and the 3D landmarks in datum space as the prior condition signal of the head geometry to directly input into MetaHead-F, so that MetaHead-F can accurately and effectively reconstruct, generate and control the head identity and expression, covering challenging poses and expressions. 3DMM texture coefficients and spherical harmonic illumination [13] coefficients are input to MetaHead-F as appearance condition signals.

MetaHead-F is the main body of the MetaHead engine model. We also propose LabelHead, a top-down paradigm in which MetaHead could generate heads consistent with the customizable head labels. The specified head features would be embedded in the latent space of MetaHead-F, leading it to a huge design space. We can assign label values to each feature to generate head images with various feature labels. We took landmark feature(see Sec. 4.5 and Sec. 4.6) and eye gaze angle feature(see Appendix A.4) as examples to verify that the labeled data generated by MetaHead can remarkably assist real data and significantly surpass the labeled data synthesized by graphics-based methods.

It is noteworthy that when reconstructing the heads, MetaHead-F would perform backward fitting on labels of head features. Therefore, MetaHead can bidirectionally inference labels of features such as landmark coordinates, eye gaze angle, hair color and so on that are difficult to annotate before. In summary, our main contributions include:

- We propose MetaHead, the unprecedented learning-based super-realistic controllable head engine, which combines and realizes the head generation, reconstruc-

tion, 3D control, and generating heads consistent with the given head labels. Furthermore, it can also bottom-up estimate the labels of head features bidirectionally.

- We propose a controllable head radiance field (MetaHead-F) to generate or reconstruct view-consistent 3D controllable digital heads, which unifies the advantages of end-to-end and GAN structures, and achieves state-of-the-art in head reconstruction accuracy, control accuracy and generation visual quality.

- We propose a generic top-down image generation framework LabelHead to generate heads with the customizable feature labels. It enables synthesizing large amount of labeled head images with various shapes and appearances. The developed framework Label-Head can be applied to any shape-appearance related fields other than human head images.

## 2. Related Work

**Learning-based Head Reconstruction, Generation and Control.** The generation visual quality of Generative Adversarial Net(GAN) [12] is impressive. For 2D-based GANs, there are two main lines to control the head attributes. Since the well-designed StyleGAN [25] contains a semantically rich latent space, a commonly-used approach [43, 15, 44] is to explore and find "walking" directions that control a specific attribute of interest. Other works [10, 45] propose to use contrastive learning [21], collecting attribute-decoupled image pairs(different in only one attribute) to allow the model to learn the control of interested attributes. However, these methods rely on 2D-based generative models and thus the generated heads lack 3D view consistency.

NeRF [35]-based generative models represent a 3D scene as a radiance field parameterized by a Multi-Layer Perceptron (MLP). pi-GAN [6] has adopting a SIREN-based NeRF representation as the generator. Because of the use of discriminant loss and the pure NeRF structure, pi-GAN cannot be trained at high resolution. GRAM [11] propose to regulate point sampling and radiance field learning on 2D manifolds to enhance the image quality. EG3D [5] introduce a novel NeRF representation tri-plane and a dual-discriminator to improve the generation quality and rendering efficiency. These models improve the 3D consistency, but only randomly output heads and do not support head controllability.

3D Morphable Model(3DMM) [4] is a parametric face model which represent the face as a linear combination of a set of principle components derived from 3D scans via Principle Components Analysis (PCA). HeadNeRF [20] is a NeRF-based reconstruction model that introduces 3DMM
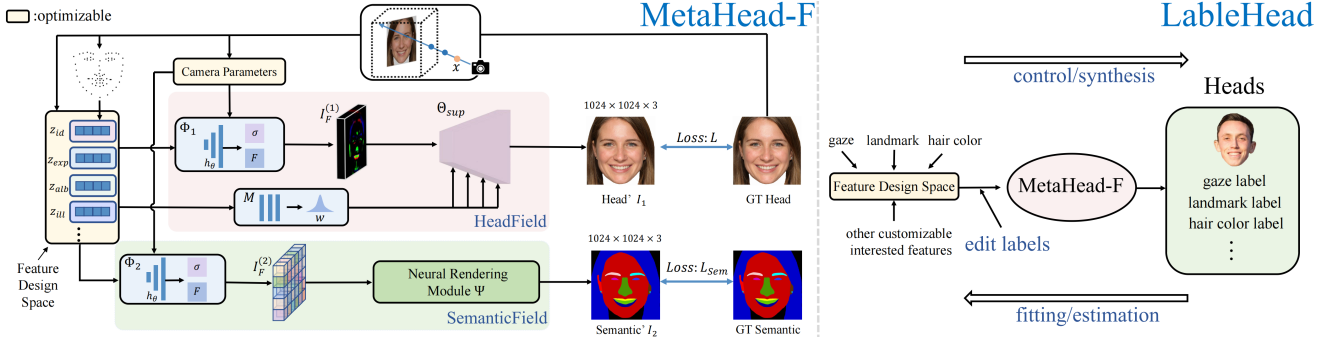
Figure 1. Overview of digital head engine MetaHead. It consists of a controllable head radiance field (MetaHead-F) to super-realistically reconstruct or generate view-consistent 3D controllable digital heads and a generic top-down image generation framework LabelHead to generate heads consistent with the given customizable feature labels.
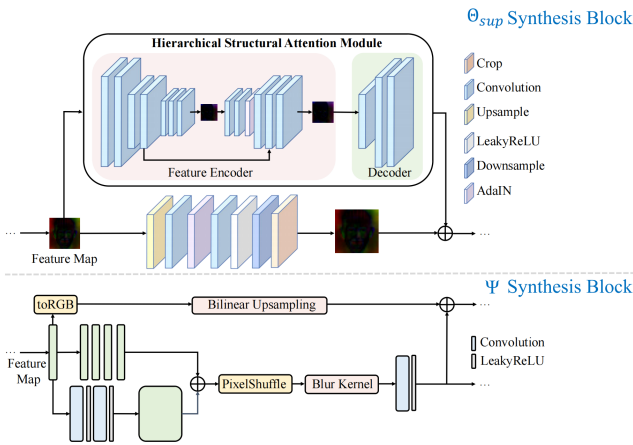


Figure 2. Synthesis block of super-resolution module $\Theta_{sup}$(Top) and neural rendering module $\Psi$(Bottom).

conditional priors. Its visual quality lacks realism and clarity, probably because it is essentially a decoder. At the same time, there are obvious discrepancy in expression, identity, hair contour and so on between the reconstructed head and the original head.

**Labeled Head Synthesis and Digital Head Engine.** There are many human head features such as landmark coordinates, eye gaze angle, and hair color that are difficult to precisely label. Existing labeled database such as [31, 40, 55, 22, 41, 51, 28] rely on experienced manual annotators, but due to the difficulty of collection, the data lacks geometry and appearance variation. This constraints the effectiveness of estimation models for head features. An increasing research [48, 49, 2] seeks to replace or supplement real data with synthetic heads in some head-related tasks such as eye gaze estimation and face recognition, depending on graphics-based methods. The synthetic texture distribution however suffers a gap with real texture distribution and thus the method is far from satisfactory. Existing learning-based head models has a potential to address this issue, but they are not yet capable to generate a head consistent with the

given label of the interested feature.

To address the above issues, and go one step further, we hope to propose a unified framework to fully realize the hyper-realistic head reconstruction, control, generation of heads consistent with the given label of the interested feature, and inversely estimate the head label. To the best of our knowledge, there is no similar work to our digital head engine MetaHead yet.

## 3. Method

In this section, we introduce how we design the hyper-realistic digital human engine MetaHead, including its reconstruction and synthesis framework controllable head radiance field(MetaHead-F), the loss terms, strategy of integrating the end-to-end decoder structure with pre-trained generator, and image generation framework LabelHead that generates heads consistent with the given customizable feature labels. The feature design space of MetaHead-F has many customizable options in addition to identity, illumination and texture. Since the expression prior is a necessary prerequisite in head decoupling control, here we take landmark feature as an example. Experiments illustrate that these four priors have enabled the model to achieve the state-of-the-art reconstruction and generation results. Examples of adding other features to the feature design space are given in Sec. 4.3 and Fig. 5 (d), (f)(please see the project page for video demos and more examples), which also verify that LabelHead allows MetaHead-F to precisely control more attributes such as gaze and hair color in addition to the common attributes.

### 3.1. Controllable Head Radiance Field

As shown in the left part of Fig. 1, MetaHead-F consists of HeadField and SemanticField. We set the conditional prior signals inputted to MetaHead-F to be optimizable, so that they could act as a bridge between HeadField and SemanticField, and the loss of SemanticField could exert semantic correction on HeadField through the optimization in

latent space.

3DMM is just a PCA face model regressed on limited 3D scans, and thus the widely used 3DMM identity and expression priors do not capture the fine-grained shape or geometric details of real faces. For a given 2D head image, we pretrained an encoder to map the extracted feature with the face recognition model AdaFace [29] to a lower dimension 128, which acts as the 3D identity prior signal $\mathbf{z}_{id}$.

For expression prior signal, we first randomly initialize the 3DMM coefficients and generate the corresponding face mesh. We specify the vertices on mesh corresponding to 68-point landmarks. Then we project them onto the 2D head image, and optimize the 3DMM coefficients and the corresponding mesh by reducing the distance between the projection points and the ground truth landmark. After that, we project the optimized vertices to the 3D datum space by performing the inverse camera transformation, and denote the result vertices as 3D-HeadPoints, so as to guarantee that they does not contain camera pose information and can be decoupled from the latter. Furthermore, mesh (controlled by implicit 3DMM expression coefficient) cannot close eyes, but 3D-HeadPoints could have precise control on challenging expressions during training since it could be explicitly optimization. Thus, we specify 3D-HeadPoints as expression prior signal $\mathbf{z}_{exp}$. In fact, we can incorporate more points as 3D-HeadPoints. Experiments illustrate that 68 points could already guarantee precise 3D reconstruction and control of expressions(see Tab. 1 and Tab. 2).

Besides, we input the 3DMM texture coefficient from the above optimization as the head texture condition signal $\mathbf{z}_{alb}$. Similarly, we regress to get 27-dimension spherical harmonic illumination coefficient as the head illumination condition signal $\mathbf{z}_{ill}$.

Given a 2D head image, first we randomly sample the points along the casted camera rays, denoted as $\mathbf{x} \in \mathbb{R}^3$, and perform positional encoding to it. The result $\gamma(\mathbf{x})$ is input to the MLP-based implicit neural function $h_\theta$:

$$h_\theta : (\gamma(\mathbf{x}), \mathbf{z}_{id}, \mathbf{z}_{exp}, \mathbf{z}_{alb}, \mathbf{z}_{ill}) \mapsto (\sigma, F), \quad (1)$$

where $\theta$ represents the network parameters, $\sigma$ is the density value at $\mathbf{x}$, and $F$ is an intermediate feature related to the radiance color at $\mathbf{x}$. After that, the final pixel color of feature map $I_F \in \mathbb{R}^{1024 \times 36 \times 36}$ is given by volume rendering:

$$I_F(r) = \int_0^\infty w(t) \cdot F(r(t)) dt$$
$$\text{where} \quad w(t) = exp(-\int_0^t \sigma(r(s)) ds) \cdot \sigma(r(t)). \quad (2)$$

$t$ defines a sample point within near and far bounds and $r(t)$ represents a ray emitted from the camera center. The above structure are the volume rendering modules in HeadField

and SemanticField, denoting them as $\Phi_i$ for $i \in \{1, 2\}$ respectively, and the corresponding feature maps are $I_F^{(i)}$ for $i \in \{1, 2\}$.

We designed a super-resolution module $\Theta_{sup}$ in Head-Field. We then impose the conditional supervision signal into mapping network $\mathbf{M}$ and perform the result $\mathbf{w}$ to strengthen the 3D head prior signal in $\Theta_{sup}$. Through the combination of end-to-end decoder structure and GAN generator(pre-trained $\Theta_{sup}$), we greatly improved the visual quality of the output heads.

Existing state-of-the-art face reconstruction and generation models suffer the dynamic-scene problem that the output head texture such as hair, eyebrow and teeth would visibly flickering during changing camera poses. To address this artifact, we design a hierarchical structural attention module(Fig. 2) customized for 3D head vision tasks in each synthesis block of $\Theta_{sup}$. With the above design, we efficiently eliminate the structure distortion and 3D texture flickering(please watch the video demo in project page). For SemanticField, the architecture of synthesis block of neural rendering upsampler $\Psi$ is shown in Fig. 2. $\Theta_{sup}$ and $\Psi$ ensure that the resolution of the output heads $I_1$ and semantic mask $I_2$(a by-product) increases gradually, and the outputs could possess multi-view consistency.

### 3.2. Field Learning & Decoder-GAN Combination

The training loss functions include the $L_1$ loss $\mathcal{L}_{\text{photo}}$ and the perceptual LPIPS [50] loss $\mathcal{L}_{\text{perc}}$ between the reconstructed head $I_1$ and the real head $I_{\text{GT}}^1$, and the $L_1$ semantic loss $\mathcal{L}_{\text{sem}}$ between the generated mask $I_2$ and the ground truth mask $I_{\text{GT}}^2$. We also impose an identity loss:

$$\mathcal{L}_{\text{ID}} = 1 - \langle R(I_1), R(I_{\text{GT}}^1) \rangle, \quad (3)$$

where $R$ is an InsightFace [14] face recognition network and $\langle \cdot, \cdot \rangle$ computes the cosine similarity between the arguments. Our proposed conditional supervision signals are inherently semantically decoupled and MetaHead-F would gradually learn the disentanglement control during optimization in latent space, which is guided by $\mathcal{L}_{\text{photo}}$ and $\mathcal{L}_{\text{perc}}$, thus requiring no additional decoupling loss or contrastive learning. To preserve the visual attributes of the input heads, we minimize the $L_2$ norm of the optimized step in the feature space:

$$\mathcal{L}_{\text{reg}} = \sum_* \|\mathbf{z}_* - \mathbf{z}_*^0\|_2, \quad (4)$$

where $*$ stand for the four head prior conditions, and $\mathbf{z}_*^0$ denotes the original value.

Previous work are usually a single GAN or decoder structure, because organizing the pre-trained generator after another network and training together would only lead to disordered signal transmission, thus generating mosaic-like images full of colored blocks. Our pre-trained gen-
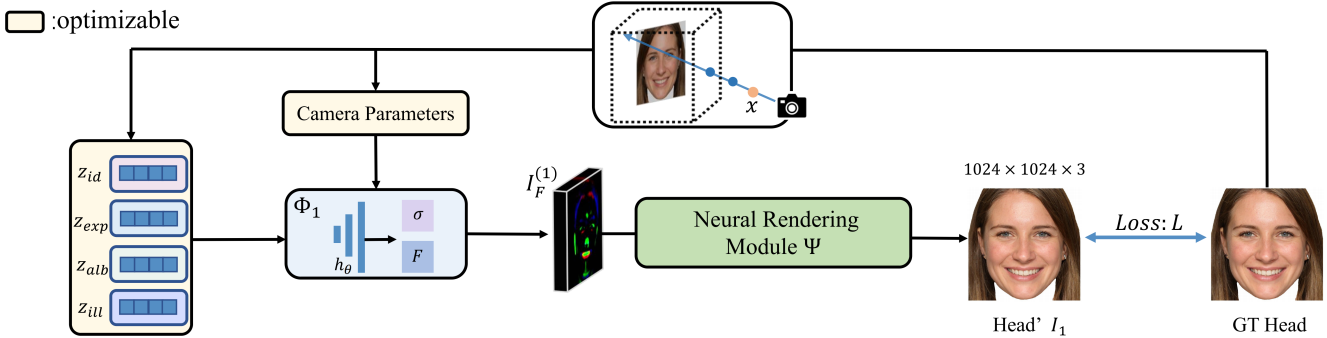
Figure 3. The model architecture of the baseline in ablation studies on MetaHead-F model designs. Our baseline removes the SemanticField and the super-resolution module $\Theta_{sup}$. In order to ensure that the baseline could output the same resolution as the HeadField, we replace $\Theta_{sup}$ with $\Psi$, which exactly raises the performance of the baseline. Meanwhile, we replace the identity and expression conditional supervisory signals $\mathbf{z}_{id}$ and $\mathbf{z}_{exp}$ with the traditional 3DMM fitting coefficients in baseline.

erator $\Theta_{sup}$ requires an input latent feature $\mathbf{w}$ and a low-resolution Fourier feature. The interface of this Fourier feature makes it possible to combine the end-to-end structure and the generator. We analyzed the reason for the signal disorder, which is in that the distribution of input feature map $I_F^{(1)}$ is far from the pre-trained distribution of the generator. Based on this assumption, we design a multi-stage training strategy. In the first stage, we use $L_2$ loss denoted as $\mathcal{L}_{dist}$ to reduce the distribution deviation between the output distribution of $\Phi_1$ and the input distribution of pre-trained $\Theta_{sup}$. In the second stage, we reduce the weight of $\mathcal{L}_{dist}$ and the following formula is used to mildly replace the input feature map $M^{in}$:

$$M^{in} = (1 - \lambda) \cdot \varphi + \lambda \cdot I_F^{(1)}, \tag{5}$$

where $\varphi$ is a Fourier-distributed random feature and $\lambda$ is a weight parameter.

This simple yet effective method greatly improves the output visual quality of HeadField, and enables it to use real images as supervision, successfully combining the advantages of end-to-end training methods(supervised learning) and generative adversarial networks(unsupervised learning). Benefit from this, we bring forward an unparalleled GAN generation module HeadField with decoupled attributes control.

Therefore, the total learning objective is:

$$\mathcal{L}_{total} = \lambda_{photo}\mathcal{L}_{photo} + \lambda_{perc}\mathcal{L}_{perc} + \lambda_{sem}\mathcal{L}_{sem} + \\ \lambda_{ID}\mathcal{L}_{ID} + \lambda_{reg}\mathcal{L}_{reg} + \lambda_{dist}\mathcal{L}_{dist}, \tag{6}$$

where $\lambda_*$ are the loss weights.

### 3.3. LabelHead: Top-down Head Image Synthesis

Thanks to our controllable end-to-end head model MetaHead-F, and its huge feature design space, we can add any customizable head-related attributes to the feature design space, such as landmark(arbitrary number of points),

| Model | CelebAMask-HQ [32] | | | | FFHQ [25] | | | |
|---|---|---|---|---|---|---|---|---|
| | AED↓ [39] | ID↑ [5] | PSNR↑ | SSIM↑ | AED↓ | ID↑ | PSNR↑ | SSIM↑ |
| baseline | 0.1103 | 0.5210 | 19.6 | 0.702 | 0.1009 | 0.5439 | 21.6 | 0.755 |
| w/ SemanticField | 0.0939 | 0.5788 | 22.85 | 0.8236 | 0.0826 | 0.6051 | 23.01 | 0.8607 |
| w/ our identity code | 0.1012 | 0.8044 | 21.45 | 0.8433 | 0.0908 | 0.8667 | 22.08 | 0.8878 |
| w/ our expression code | 0.0355 | 0.5410 | 21.80 | 0.8218 | 0.0222 | 0.5643 | 22.48 | 0.8751 |
| w/ $\Theta_{sup}$ | 0.0904 | 0.5443 | 28.43 | 0.8429 | 0.0857 | 0.5632 | 29.70 | 0.9032 |
| MetaHead-F | 0.0289 | 0.9058 | 29.98 | 0.8632 | 0.0150 | 0.9145 | 30.20 | 0.9247 |

Table 1. Quantitative ablation study on model designs of MetaHead-F.

eye gaze angle, pixel-wise semantic category [32] and hair color. We train on image-label pairs. Then we freely edit the label value of interested attributes and MetaHead-F could generate a 3D controllable head which is consistent with the given label. Except for this, we can further fix one feature and manipulate other feature values, such as dimming the illumination and enlarging the pose, and thus to generate image-label paired training data under challenging scenes. LabelHead enables us to precisely 3D control more customizable attributes in addition to the common ones. Besides, as shown in the right part of Fig. 1, we can impose the photometric loss($\mathcal{L}_{photo}$ and $\mathcal{L}_{perc}$) to bottom-up fit the various attributes of interest of the input head, and estimate the label bidirectionally(see Sec. 5). This is unreachable for [48] under graphics-based forward rendering framework.

## 4. Experiments

### 4.1. Implementation Details

We train MetaHead-F on FFHQ dataset [25] and remain 6000 for reconstruction test. We also processed monocular videos of 500 people, and collected other videos containing exaggerated expressions and extreme poses with professional equipment as training data. For super-resolution module $\Theta_{sup}$, we pretrained it with the discriminator of StyleGAN2 [26] on FFHQ training data. The final model is trained for 120 hours with four 3090 GPUs.
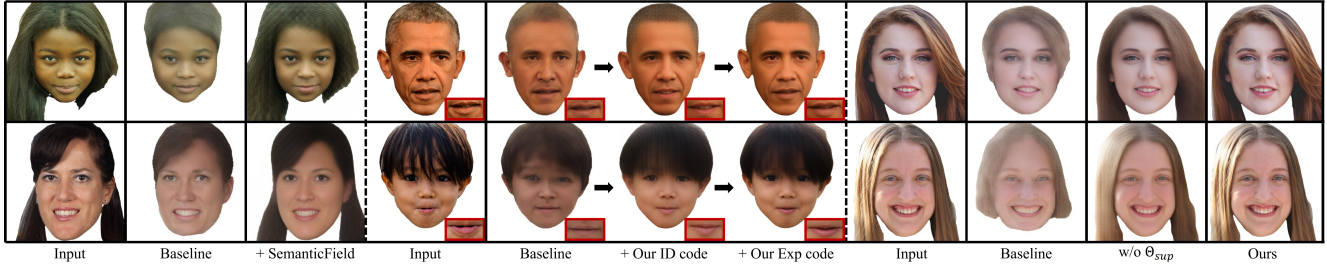
Figure 4. Qualitative ablation study on MetaHead-F model designs. Expressions are highlighted in red. Id: identity, Exp: expression. Refer to Sec. 4.2 for details.
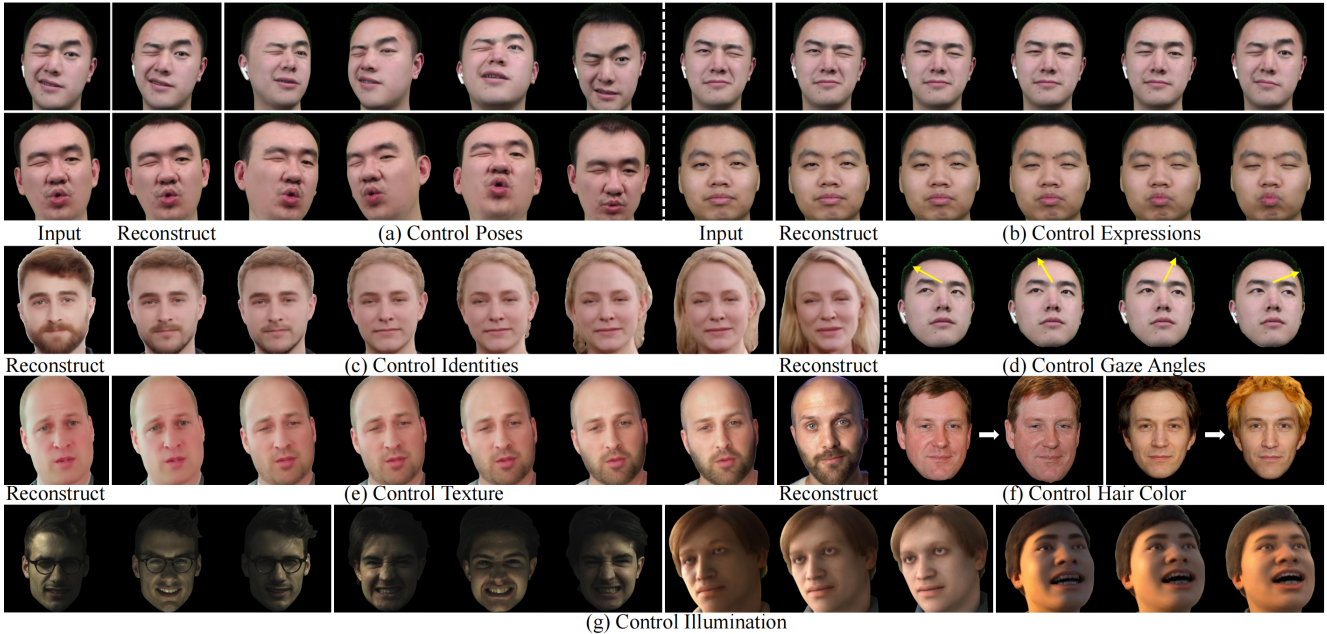


Figure 5. Attributes controlling results of MetaHead-F. (c) and (e) demonstrate progressive interpolation(from left to right) over the identity and texture of the reconstructions. MetaHead-F can achieve highly disentangled and precise 3D control over the identity, expression, texture, illumination, gaze angle, hair color, and camera pose of heads, and can also generate heads in a 3D-aware stylization manner.

## 4.2. Ablation Study on MetaHead-F Model Designs

We conduct ablation studies on FFHQ [25] testing data to validate the effectiveness of each component proposed in our head model MetaHead-F. Fig. 3 illustrates the model architecture of the baseline, which shares the same training data and number of training epochs as MetaHead-F. Our baseline removes the SemanticField and the super-resolution module $\Theta_{sup}$. In order to ensure that the baseline could output the same resolution as the HeadField, we replace $\Theta_{sup}$ with $\Psi$, which exactly raises the performance of the baseline. Meanwhile, we replace the identity and expression conditional supervisory signals $\mathbf{z}_{id}$(identity code) and $\mathbf{z}_{exp}$(expression code) with the traditional 3DMM fitting coefficients in baseline. The reconstruction results are shown in Fig. 4 and Tab. 1. We add SemanticField to the baseline in the left part of Fig. 4 and in row 4 of Tab. 1. We replace 3DMM coefficients with our proposed condi-

tional supervisory signals $\mathbf{z}_{id}$ and $\mathbf{z}_{exp}$ to the baseline in rows 5 and 6 of Tab. 1, respectively. Besides, we replace 3DMM coefficients with our proposed signals one by one in the middle part of Fig. 4. Moreover, we replace $\Psi$ with $\Theta_{sup}$ in baseline in row 7 of Tab. 1, and we replace $\Theta_{sup}$ in HeadField with $\Psi$ in the full model MetaHead-F in the right part of Fig. 4. ID [5] and AED [39] are common metrics to measure the identity and expression 3D reconstruction and control accuracy, respectively, calculated between the test and reconstructed heads. PSNR and SSIM are both visual quality metrics.

The studies show that SemanticField precisely control the hair and mouth shape, our proposed identity and expression conditional prior signal enhance the control over fine-grained geometric details(identity and expression), and super-resolution module $\Theta_{sup}$ significantly improves the output visual quality. Ablation studies in dynamic scenes are given in the video demo(please see the project page),
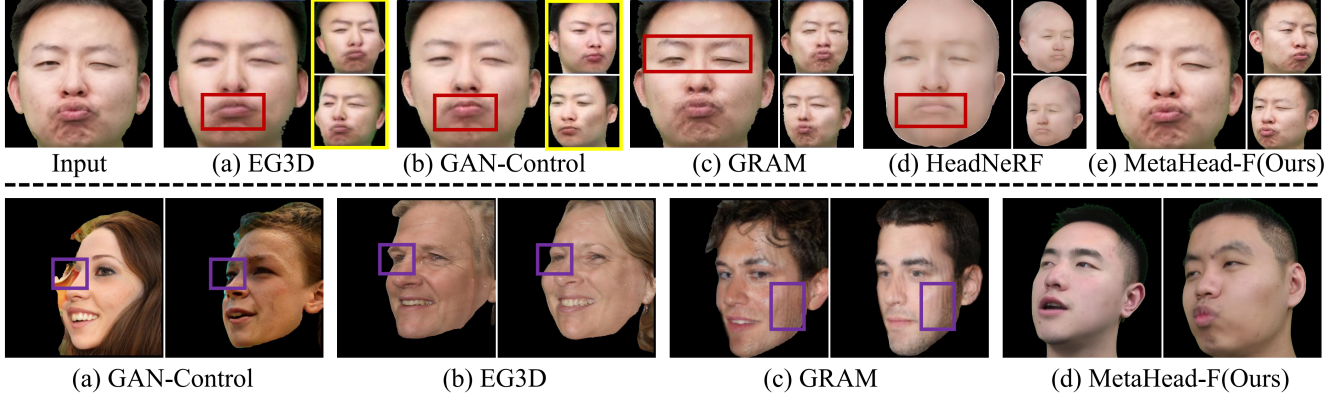
Figure 6. (Top)Reconstruction and attributes control accuracy comparison with existing methods. Red: expression inaccuracy. Yellow: identity and expression inconsistency. (Bottom)Heads generated by MetaHead-F in extreme poses, compared with existing methods. Previous methods generate heads with unreasonable artifacts such as blurry or missing eyes, and layered skin, while MetaHead-F produces robust results.

| Methods | Venue | FID↓ [19] | $DS_\alpha \uparrow$ [10] | $DS_\beta \uparrow$ | $DS_\gamma \uparrow$ | $DS_\theta \uparrow$ |
|---|---|---|---|---|---|---|
| GIRAFFE [36] | CVPR'21 | 32.6 | - | - | - | 36.2 |
| pi-GAN [6] | CVPR'21 | 55.2 | - | - | - | 34.6 |
| GRAM [11] | CVPR'22 | 17.9 | - | - | - | 37.6 |
| EG3D [5] | CVPR'22 | **4.7** | - | - | - | 39.2 |
| DiscoFaceGAN [10] | CVPR'20 | 56.6 | 7.85 | 80.4 | 489 | 36.7 |
| PIRenderer [39] | ICCV'21 | 73.4 | 8.16 | 64.3 | - | 30.2 |
| GAN-Control [45] | ICCV'21 | 14.6 | 9.26 | 69.2 | 496 | 39.8 |
| HeadNeRF[20] | CVPR'22 | 160.5 | 7.91 | 52.1 | 471 | 41.5 |
| MetaHead-F(Ours) | - | <u>6.7</u> | **15.52** | **93.3** | **511** | **47.2** |

Table 2. Quantitative comparison of visual quality, and disentanglement and accuracy of attributes control with existing methods.

which shows that $\Theta_{sup}$ and our proposed hierarchical structure attention module nicely solve the chronic problem of hair and teeth flickering when the viewing angle changes.

## 4.3. Qualitative Evaluation on MetaHead-F

**Attributes Controlling Results** Fig. 5 shows the reconstruction results of MetaHead-F and its separate control over identity, expression, texture, illumination, and camera poses. The last two showcases in (g) illustrate that MetaHead-F can stylizedly generate and control digital heads. As shown in Fig. 5, MetaHead-F provides precise and view-consistent 3D control over the shape and appearance of heads in a highly decoupled manner, including challenging expressions and illumination. In comparison, previous models [5, 11, 10, 45] could only output or control moderate expressions. Furthermore, if we add the gaze angle, hair color and semantic label into the feature design space of MetaHead-F, MetaHead-F could further precisely control these three head features(see (d), (f) in Fig. 5 and Appendix A.5), which cannot be achieved by any existing models.

**Comparison on Reconstruction and Disentanglement** We present results on reconstruction accuracy and disentanglement of attributes control in row 1 of Fig. 6, com-

paring our method against four state-of-the-art head synthesis( [5, 45, 11]) or reconstruction-only( [20]) methods. More examples are displayed in Fig. 12 of the supplementary material. For each sub-figure, the left column corresponds to the reconstruction result, and the right column corresponds to the results of controlling camera poses. In Fig. 6 row 1 (a)-(d), these four methods fail to reconstruct the right expression, and a few examples are highlighted in red. EG3D and GAN-Control fail to preserve the identity and expression consistency when only changing camera poses, which are highlighted in yellow. Besides, HeadNeRF missing the fine-grained appearance details such as chin wrinkle. In contrast, Fig. 6 row 1 (e) shows that MetaHead-F precisely reconstructs all factors(covering shape and appearance) of the heads with highly consistent and precise pose control. Only poses change during the pose editing while other head properties remain unchanged, which shows a good and robust control disentanglement.

**Comparison on Large Pose Generation** We present large pose generation results of GAN-Control [45], EG3D [5], GRAM [11] and ours in row 2 of Fig. 6. Previous methods fail to generate robust results, and suffer the unreasonable or blurry artifacts, which are highlighted in purple. A possible reason is that the data distribution of extreme poses in their training data is far from enough. Besides, GRAM produces layered artifacts in large pose, which is perhaps caused by the surface manifold learning. In comparison, MetaHead-F could generate natural novel-view heads in the extreme camera pose.

## 4.4. Quantitative Evaluation on MetaHead-F

We compare our method to several state-of-the-art generative or reconstruction-only methods. GIRAFFE [36], pi-GAN [6], GRAM [11] and EG3D [5] are only able to
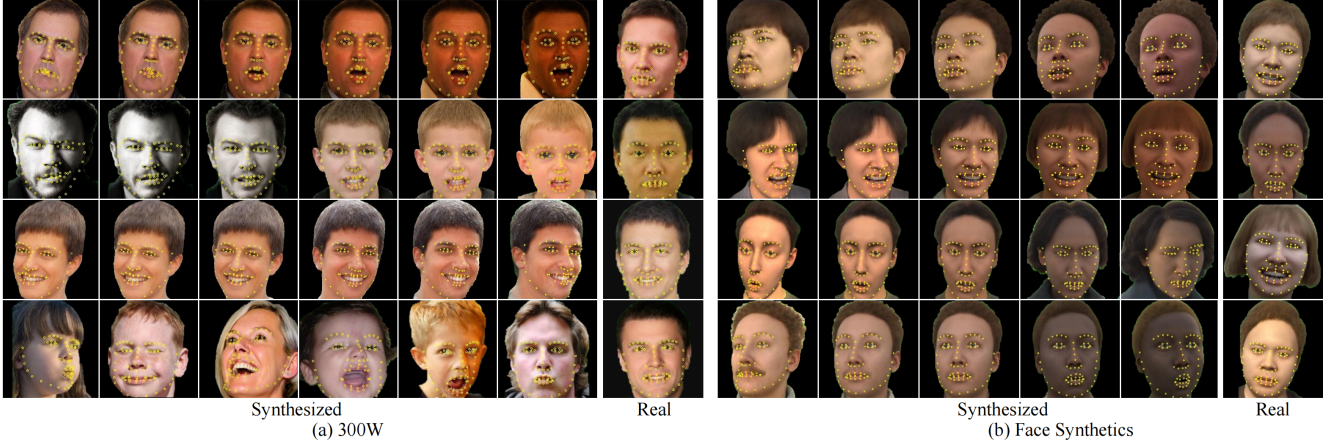
Figure 7. We fine-tuned MetaHead-F on **300W** [40] and Microsoft Face Synthetics [48] respectively. Right side of sub-figures show real images. MetaHead-F synthesizes heads(left side of sub-figures) indistinguishable from real images with accurate feature(landmark) labels(marked in yellow) and much larger shape and appearance variation.
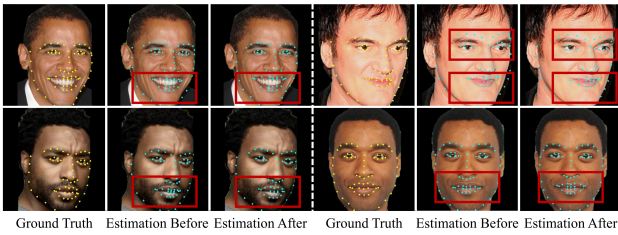


Figure 8. Landmark estimations on **300W** [40] testing data before and after adding the synthesized heads.
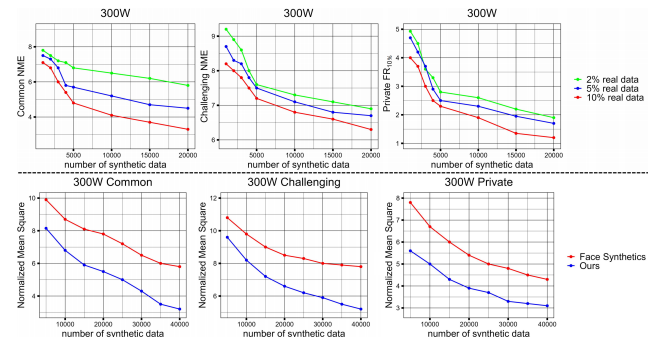


Figure 9. (Top) Our synthetic data significantly reduce the landmark estimation error. (Bottom) Landmark estimation error comparison between using state-of-the-art method(Face Synthetics Dataset [48]) and our synthetic data.

control pose while DiscoFaceGAN [10], PIRenderer [39], GAN-Control [45] and HeadNeRF [20] are controllable methods.

We compare the visual quality of head generation using the Frèchet Inception Distance(FID) scores [19](see column 3 of Tab. 2). We use FFHQ [25] and CelebAMask-HQ [32] as the real data, and measure FID scores between them and 50k randomly generated images. The results show that our result is only second to EG3D. It is worth mentioning that since the output resolution of MetaHead-F is 1024, to be fair to the above methods, we sacrificed to downsample the output to resolution of 512 with bilinear interpolation, which would result in the degradation of visual quality, and thus causing a higher FID score. Besides, EG3D couldn't control any attributes except for pose. We compare the disentanglement and accuracy of attributes control on FFHQ [25] testing data using the Disentanglement Score(DS)(see columns 4-7 of Tab. 2), which was proposed in [10]. $DS_\alpha$, $DS_\beta$, $DS_\gamma$ and $DS_\theta$ stand for DS score of identity, expression, illumination and pose, respectively. The results show that MetaHead-F achieves significantly better disentanglement and more precise control compared to existing methods.

## 4.5. Qualitative Evaluation on LabelHead

To equip the proposed MetaHead-F with the capability to generate heads consistent with the customizable head labels, we offer a generic top-down image generation framework LabelHead, which is able to be applied on any shape-appearance related head features. Due to limited space, we only discuss landmark feature in detail here, and experiment paradigm for other features are in the same way. We also discuss eye gaze angle and semantic label(category) feature quantitatively and qualitatively in Appendix A.4 and Appendix A.5, including synthesizing heads with consistent feature labels and the direct control of eye appearance and head shape.

**Shape and Appearance Variation** We embed the landmark feature 3D-HeadPoints into the feature design space and fine-tune the pre-trained MetaHead-F on **300W** [40] training images for 2D face alignment task and Face Syn-

thetics [48] for 3D face alignment task, respectively, which have ground truth landmark labels. By providing a sequence of head attributes such as identity, illumination, texture, pose and 3D-HeadPoints as conditional signals, we can easily generate heads with diverse shape and appearance variation as shown in Fig. 7, while the latter is very important in label-annotated head generation. We then project 3D-HeadPoints onto the corresponding 2D images to obtain the landmark labels. We show 2D landmark labels in Fig. 7 (a) rows 1-3, 3D landmark labels in Fig. 7 (a) row 4 and Fig. 7 (b). Fig. 7 (a) shows that LabelHead could generate super-realistic heads under challenging scenes with accurate labels. Fig. 7 (b) shows that MetaHead-F could synthesize stylized heads with equally accuracy. Besides, Fig. 7 also indicates that LabelHead can generate multi-task(2D and 3D) landmark labels with good generalization.

**Qualitative Evaluation on Label-Estimation**  We first qualitatively evaluate whether synthesized heads help label estimation with small amount of real data. We train a landmark estimator using widely known backbone ResNet34 [18] on image-landmark pairs of **300W** [40] training data and test on **300W** [40] testing data. As shown in Fig. 8, the model is under-fitted with purely real data which cause poor performance. However, after we add our generated heads which are 10 times larger than the amount of real data, the performance remarkably improves. Some accuracy comparison examples are highlighted in red in Fig. 8. This demonstrates that our synthesized heads indeed capture the correlation between landmark, head shape and head appearance, and can be great useful for applications with small amount of real data.

## 4.6. Quantitative Evaluation on LabelHead

**Quantitative Evaluation on Label-Estimation**  We also quantitatively evaluate whether our generated heads help label estimation with small amount of real data. The experiments are conducted on **300W** [40] following the common protocol in [54], where we perform testing on three parts: the common, challenging and private subsets. The alignment accuracy is evaluated by the Normalized Mean Error(NME) and Failure Rate below a $10\%$ error threshold($FR_{10\%}$), while the lower means the better. We train a landmark estimator using ResNet34 [18] on image-landmark pairs of training data. Training data consists of k% of real images and n generated heads. Fig. 9 row 1 shows that as we continuously add the synthetic heads, the performance keeps improving until saturation.

Next we compare our generated data with the state-of-the-art landmark-labeled image generation method [48] on head synthesis. Our synthesized heads and Microsoft Face Synthetics generated by [48] are used to train the landmark estimator ResNet34 respectively and the result models are

| | Hue(0-360°) | Saturation(0-1) | Value(0-1) |
|---|---|---|---|
| Error(MSE [1]) | 0.113 | 0.049 | 0.011 |

Table 3. Bottom-up estimation error of the hair HSV color values using LabelHead.

tested on **300W** common, challenging and private subsets. Since the landmark labels are 3D in Microsoft Face Synthetics, we utilize InsightFace [14] to perform a translation on jawline from 3D to 2D for fairness as guided in [48]. As shown in Fig. 9 row 2, our training data achieve superior results on each subset. The key reason is that MetaHead-F could generate super-realistic heads which is approximate to the distribution of real data and further cover the diverse challenging scenes.

## 5. Application

Due to the limited space, we demonstrate applications including one-shot facial retargeting, text-to-head generation and text-based 3D head manipulation in Appendix B in the supplementary material.

**Attributes Fitting Using LabelHead**  We take the fitting/estimation of the hair color feature as an example of the bottom-up label estimation here. We first embed the HSV(Hue, Saturation and Value) color values of hair as a feature into the feature design space of MetaHead-F and train MetaHead-F on FFHQ [25]. We then test on CelebAMask-HQ [32]. We randomly initialize the label values of all features including hair color, then using the photometric loss between the test image and the generated head of MetaHead-F to optimize the label value, thus getting the HSV estimation. Tab. 3 shows that the bottom-up fitting results are of high accuracy, which indicates that hair color feature captures the latent semantic position of hair and demonstrates the effectiveness of the bidirectional label estimation using LabelHead.

## 6. Conclusion

We have presented a digital head engine Meta-Head, which consists of a controllable head radiance field(MetaHead-F) to super-realistically generate or reconstruct view-consistent controllable digital heads, enabling much more precise and decoupled 3D controllability over 3D identity, expression, texture, illumination and pose of the generated heads than existing state-of-the-art methods, and a generic top-down image generation framework LabelHead to generate heads consistent with the given customizable feature labels, which also enables MetaHead-F to control any shape-appearance related head features and bidirectionally estimate the labels of head features.

# References

[1] David M Allen. Mean square error of prediction as a criterion for selecting variables. *Technometrics*, 13(3):469–475, 1971. 9

[2] Gwangbin Bae, Martin de La Gorce, Tadas Baltrušaitis, Charlie Hewitt, Dong Chen, Julien Valentin, Roberto Cipolla, and Jingjing Shen. Digiface-1m: 1 million digital face images for face recognition. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pages 3526–3535, January 2023. 3

[3] Urs Bergmann, Nikolay Jetchev, and Roland Vollgraf. Learning texture manifolds with the periodic spatial GAN. *CoRR*, abs/1705.06566, 2017. 2

[4] Volker Blanz and Thomas Vetter. A morphable model for the synthesis of 3d faces. In *Proceedings of the 26th annual conference on Computer graphics and interactive techniques*, pages 187–194, 1999. 2

[5] Eric R Chan, Connor Z Lin, Matthew A Chan, Koki Nagano, Boxiao Pan, Shalini De Mello, Orazio Gallo, Leonidas J Guibas, Jonathan Tremblay, Sameh Khamis, et al. Efficient geometry-aware 3d generative adversarial networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16123–16133, 2022. 1, 2, 5, 6, 7, 13

[6] 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 *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5799–5809, 2021. 1, 2, 7, 13

[7] Roddy Cowie, Ellen Douglas-Cowie, Nicolas Tsapatsoulis, George Votsis, Stefanos Kollias, Winfried Fellenz, and John G Taylor. Emotion recognition in human-computer interaction. *IEEE Signal processing magazine*, 18(1):32–80, 2001. 1

[8] Antonio Criminisi, Jamie Shotton, Andrew Blake, and Philip HS Torr. Gaze manipulation for one-to-one teleconferencing. In *Computer Vision, IEEE International Conference on*, volume 2, pages 191–191. IEEE Computer Society, 2003. 1

[9] Jiankang Deng, Jia Guo, Xue Niannan, and Stefanos Zafeiriou. Arcface: Additive angular margin loss for deep face recognition. In *CVPR*, 2019. 13

[10] Yu Deng, Jiaolong Yang, Dong Chen, Fang Wen, and Xin Tong. Disentangled and controllable face image generation via 3d imitative-contrastive learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5154–5163, 2020. 1, 2, 7, 8, 13

[11] Yu Deng, Jiaolong Yang, Jianfeng Xiang, and Xin Tong. Gram: Generative radiance manifolds for 3d-aware image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10673–10683, 2022. 1, 2, 7, 13

[12] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications of the ACM*, 63(11):139–144, 2020. 2

[13] Robin Green. Spherical harmonic lighting: The gritty details. In *Archives of the game developers conference*, volume 56, page 4, 2003. 2

[14] Jia Guo, Jiankang Deng, Alexandros Lattas, and Stefanos Zafeiriou. Sample and computation redistribution for efficient face detection. *arXiv preprint arXiv:2105.04714*, 2021. 4, 9, 20

[15] Erik Härkönen, Aaron Hertzmann, Jaakko Lehtinen, and Sylvain Paris. Ganspace: Discovering interpretable gan controls. *Advances in Neural Information Processing Systems*, 33:9841–9850, 2020. 2

[16] H Rex Hartson and Deborah Hix. Human-computer interface development: concepts and systems for its management. *ACM Computing Surveys (CSUR)*, 21(1):5–92, 1989. 1

[17] Trevor Hastie, Robert Tibshirani, Jerome Friedman, Trevor Hastie, Robert Tibshirani, and Jerome Friedman. Overview of supervised learning. *The elements of statistical learning: Data mining, inference, and prediction*, pages 9–41, 2009. 1

[18] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 9

[19] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017. 7, 8

[20] Yang Hong, Bo Peng, Haiyao Xiao, Ligang Liu, and Juyong Zhang. Headnerf: A real-time nerf-based parametric head model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20374–20384, 2022. 1, 2, 7, 8, 13

[21] Ashish Jaiswal, Ashwin Ramesh Babu, Mohammad Zaki Zadeh, Debapriya Banerjee, and Fillia Makedon. A survey on contrastive self-supervised learning. *Technologies*, 9(1):2, 2020. 2

[22] Shashank Jaiswal, Timur Almaev, and Michel Valstar. Guided unsupervised learning of mode specific models for facial point detection in the wild. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pages 370–377, 2013. 3

[23] Xin Jin and Xiaoyang Tan. Face alignment in-the-wild: A survey. *Computer Vision and Image Understanding*, 162:1–22, 2017. 1

[24] Tero Karras, Miika Aittala, Samuli Laine, Erik Härkönen, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Alias-free generative adversarial networks. *Advances in Neural Information Processing Systems*, 34:852–863, 2021. 2

[25] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4401–4410, 2019. 1, 2, 5, 6, 8, 9, 13

[26] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8110–8119, 2020. 5

[27] Petr Kellnhofer, Adria Recasens, Simon Stent, Wojciech Matusik, , and Antonio Torralba. Gaze360: Physically unconstrained gaze estimation in the wild. In *IEEE International Conference on Computer Vision (ICCV)*, October 2019. 15, 16, 18

[28] Petr Kellnhofer, Adria Recasens, Simon Stent, Wojciech Matusik, and Antonio Torralba. Gaze360: Physically unconstrained gaze estimation in the wild. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6912–6921, 2019. 3

[29] Minchul Kim, Anil K Jain, and Xiaoming Liu. Adaface: Quality adaptive margin for face recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18750–18759, 2022. 4

[30] Anders Boesen Lindbo Larsen, Søren Kaae Sønderby, Hugo Larochelle, and Ole Winther. Autoencoding beyond pixels using a learned similarity metric. In *International conference on machine learning*, pages 1558–1566. PMLR, 2016. 2

[31] Vuong Le, Jonathan Brandt, Zhe Lin, Lubomir Bourdev, and Thomas S Huang. Interactive facial feature localization. In *European conference on computer vision*, pages 679–692. Springer, 2012. 3

[32] Cheng-Han Lee, Ziwei Liu, Lingyun Wu, and Ping Luo. Maskgan: Towards diverse and interactive facial image manipulation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5549–5558, 2020. 5, 8, 9, 18, 19

[33] Lik-Hang Lee, Tristan Braud, Pengyuan Zhou, Lin Wang, Dianlei Xu, Zijun Lin, Abhishek Kumar, Carlos Bermejo, and Pan Hui. All one needs to know about metaverse: A complete survey on technological singularity, virtual ecosystem, and research agenda. *arXiv preprint arXiv:2110.05352*, 2021. 1

[34] Jakob Lovén, David A Orlando, Alla A Sigova, Charles Y Lin, Peter B Rahl, Christopher B Burge, David L Levens, Tong Ihn Lee, and Richard A Young. Revisiting global gene expression analysis. *Cell*, 151(3):476–482, 2012. 1

[35] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 65(1):99–106, 2021. 2

[36] Michael Niemeyer and Andreas Geiger. Giraffe: Representing scenes as compositional generative neural feature fields. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2021. 7, 13

[37] Pascal Paysan, Reinhard Knothe, Brian Amberg, Sami Romdhani, and Thomas Vetter. A 3d face model for pose and illumination invariant face recognition. In *2009 sixth IEEE international conference on advanced video and signal based surveillance*, pages 296–301. Ieee, 2009. 2

[38] 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 *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 20

[39] Yurui Ren, Ge Li, Yuanqi Chen, Thomas H Li, and Shan Liu. Pirenderer: Controllable portrait image generation via semantic neural rendering. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 13759–13768, 2021. 5, 6, 7, 8, 13

[40] Christos Sagonas, Epameinondas Antonakos, Georgios Tzimiropoulos, Stefanos Zafeiriou, and Maja Pantic. 300 faces in-the-wild challenge: Database and results. *Image and vision computing*, 47:3–18, 2016. 3, 8, 9

[41] Christos Sagonas, Georgios Tzimiropoulos, Stefanos Zafeiriou, and Maja Pantic. A semi-automatic methodology for facial landmark annotation. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pages 896–903, 2013. 3

[42] Martijn J Schuemie, Peter Van Der Straaten, Merel Krijn, and Charles APG Van Der Mast. Research on presence in virtual reality: A survey. *Cyberpsychology & behavior*, 4(2):183–201, 2001. 1

[43] Yujun Shen, Jinjin Gu, Xiaoou Tang, and Bolei Zhou. Interpreting the latent space of gans for semantic face editing. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9243–9252, 2020. 2

[44] Yujun Shen, Ceyuan Yang, Xiaoou Tang, and Bolei Zhou. Interfacegan: Interpreting the disentangled face representation learned by gans. *IEEE transactions on pattern analysis and machine intelligence*, 44(4):2004–2018, 2020. 2

[45] Alon Shoshan, Nadav Bhonker, Igor Kviatkovsky, and Gerard Medioni. Gan-control: Explicitly controllable gans. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14083–14093, 2021. 1, 2, 7, 8, 13, 17

[46] Justus Thies, Michael Zollhofer, Marc Stamminger, Christian Theobalt, and Matthias Nießner. Face2face: Real-time face capture and reenactment of rgb videos. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2387–2395, 2016. 1

[47] Jian-Gang Wang and Eric Sung. Study on eye gaze estimation. *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)*, 32(3):332–350, 2002. 1

[48] Erroll Wood, Tadas Baltrušaitis, Charlie Hewitt, Sebastian Dziadzio, Thomas J Cashman, and Jamie Shotton. Fake it till you make it: face analysis in the wild using synthetic data alone. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 3681–3691, 2021. 1, 3, 5, 8, 9

[49] Erroll Wood, Tadas Baltrušaitis, Louis-Philippe Morency, Peter Robinson, and Andreas Bulling. Learning an appearance-based gaze estimator from one million synthesised images. In *Proceedings of the Ninth Biennial ACM Symposium on Eye Tracking Research & Applications*, pages 131–138, 2016. 3

[50] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 586–595, 2018. 4

[51] Xucong Zhang, Yusuke Sugano, Mario Fritz, and Andreas Bulling. Appearance-based gaze estimation in the wild. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4511–4520, 2015. 1, 3

[52] Xucong Zhang, Yusuke Sugano, Mario Fritz, and Andreas Bulling. It's written all over your face: Full-face appearance-based gaze estimation. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pages 51–60, 2017. 15, 16, 18

[53] Wenyi Zhao, Rama Chellappa, P Jonathon Phillips, and Azriel Rosenfeld. Face recognition: A literature survey. *ACM computing surveys (CSUR)*, 35(4):399–458, 2003. 1

[54] Shizhan Zhu, Cheng Li, Chen Change Loy, and Xiaoou Tang. Face alignment by coarse-to-fine shape searching. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4998–5006, 2015. 9

[55] Xiangxin Zhu and Deva Ramanan. Face detection, pose estimation, and landmark localization in the wild. In *2012 IEEE conference on computer vision and pattern recognition*, pages 2879–2886. IEEE, 2012. 1, 3

# Supplementary Material

In this supplementary material, we provide both qualitative and quantitative results that were not included in our main manuscript. In order to demonstrate our dynamic results, we also provide the video demos, which can be found in the project page: https://ustc3dv.github.io/MetaHead/.

## A. Additional Experiments

### A.1. Ablation Study on MetaHead-F model designs

**Additional Qualitative Ablation Studies**  We show more ablation experiments on SemanticField, super-resolution module $\Theta_{sup}$ and our proposed conditional supervision head signals in Fig. 10 and Fig. 11. The experimental settings are the same as in Sec. 4.2 of the main manuscript. Studies show that SemanticField precisely control the hair and mouth shape, our proposed identity and expression conditional prior signal enhance the control over fine-grained geometric details(identity and expression), and super-resolution module $\Theta_{sup}$ significantly improves the output visual quality.

**We strongly recommend to watch the video demos in the project page**, which contains ablation studies on the super-resolution module $\Theta_{sup}$ with our proposed hierarchical structural attention module to solve head texture flickering problem in dynamic scenes and the qualitative comparison with the existing state-of-the-art head models.

### A.2. Qualitative Evaluation on MetaHead-F

**Additional Comparison on Reconstruction and Control Accuracy**  We present more results on reconstruction and attributes control accuracy in Fig. 12, comparing our method against four state-of-the-art head synthesis( [5, 45, 11]) or reconstruction-only( [20]) methods. For each sub-figure, column 3 corresponds to the reconstruction results, while columns 2 and 4 correspond to the results obtained by controlling the camera poses of the reconstructions. In Fig. 12, these four methods fail to reconstruct the right expression, and a few examples are highlighted in red. EG3D [5], GRAM [11] and GAN-Control [45] fail to preserve the identity and expression consistency when only changing camera poses, which are highlighted in yellow. Besides, HeadNeRF [20] missing the fine-grained appearance details such as acne and moles. The above four models fail to render the details of head accessories such as headphones, which are highlighted in blue.

In contrast, Fig. 12 shows that MetaHead-F precisely reconstructs all factors(covering shape and appearance) of the heads with highly view-consistent and precise pose control. During pose editing, only the poses change while the other

| Methods | Venue | ID↑ [5] |
|---|---|---|
| GIRAFFE [36] | CVPR'21 | 0.64 |
| pi-GAN [6] | CVPR'21 | 0.67 |
| GRAM [11] | CVPR'22 | 0.74 |
| EG3D [5] | CVPR'22 | 0.77 |
| DiscoFaceGAN [10] | CVPR'20 | 0.62 |
| PIRenderer [39] | ICCV'21 | 0.64 |
| GAN-Control [45] | ICCV'21 | 0.66 |
| HeadNeRF[20] | CVPR'22 | 0.71 |
| Ours | - | **0.79** |

Table 4. Quantitative comparison between MetaHead-F and existing methods in terms of 3D view consistency.

head properties remain unchanged. This demonstrates effective and robust control disentanglement.

**Additional Comparison on Expression Control Disentanglement**  We present more results on the expression control, comparing our method against the state-of-the-art controllable head synthesis method GAN-Control [45]. Fig. 13 shows the generated heads using GAN-Control and our method MetaHead-F respectively. For each sub-figure, column 1 shows a reference image, columns 2 to 7 show images generated with random expression control. Each row corresponds to the same person.

In Fig. 13 (a), we can see that GAN-Control fails to preserve the identity when changing only the facial expression code. A few examples are highlighted in red. Moreover, we observe that with the original range of the expression parameters, the model results in only a small variation of expressions.

Fig. 13 (b) shows that MetaHead-F generates compelling photo-realistic heads with highly consistent, precise expression control. When controlling expressions for the heads in each row, the other head attributes, such as identity and texture, remain unchanged. Besides, MetaHead-F could enable diverse expression variation including frowning, pouting, curling lips, etc.

### A.3. Additional Quantitative Evaluation on MetaHead-F

We quantitatively evaluate the 3D view consistency assessed by multi-view facial identity consistency(ID), following the method proposed in  [5]. We calculate the mean Arcface [9] cosine similarity score between pairs of views of the same synthesized head rendered from two random camera poses. Tab. 4 reports the comparison results on the FFHQ [25] testing data. Our method achieves the state-of-the-art view consistency.
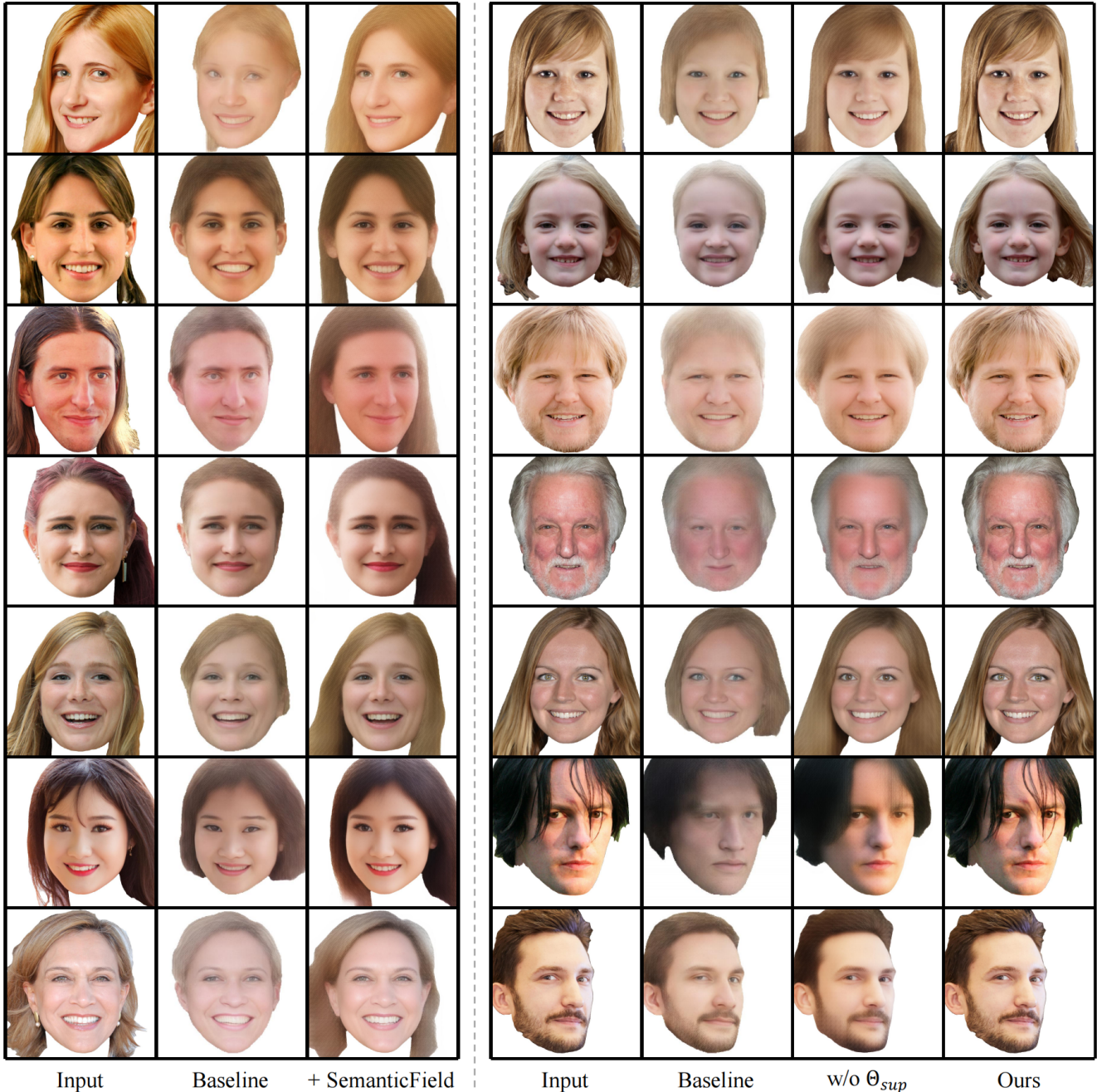
Figure 10. Qualitative ablation study on MetaHead-F model designs for SemanticField and the super-resolution module $\Theta_{sup}$. Please refer to Sec. 4.2 of the main manuscript for details.

| Input | Baseline | + SemanticField | | Input | Baseline | w/o $\Theta_{sup}$ | Ours |

## A.4. Additional Quantitative Evaluation on Label-Head

Due to limited space, we only discuss landmark feature in detail in the main manuscript. We discuss synthesizing heads with consistent eye gaze labels here and exhibit the direct control of eye appearance using gaze label in the video demo(please see the project page).

**Quantitative Evaluation on Label-Estimation**   In order to better understand humans – their desires, intents and states of mind – one need to be able to observe and perceive certain behavioral cues. Eye gaze direction is one such cue: it is a strong form of non-verbal communication, signalling engagement, interest and attention during social interactions. In recent years, methods for gaze estimation have not yet reached a satisfactory level of performance.

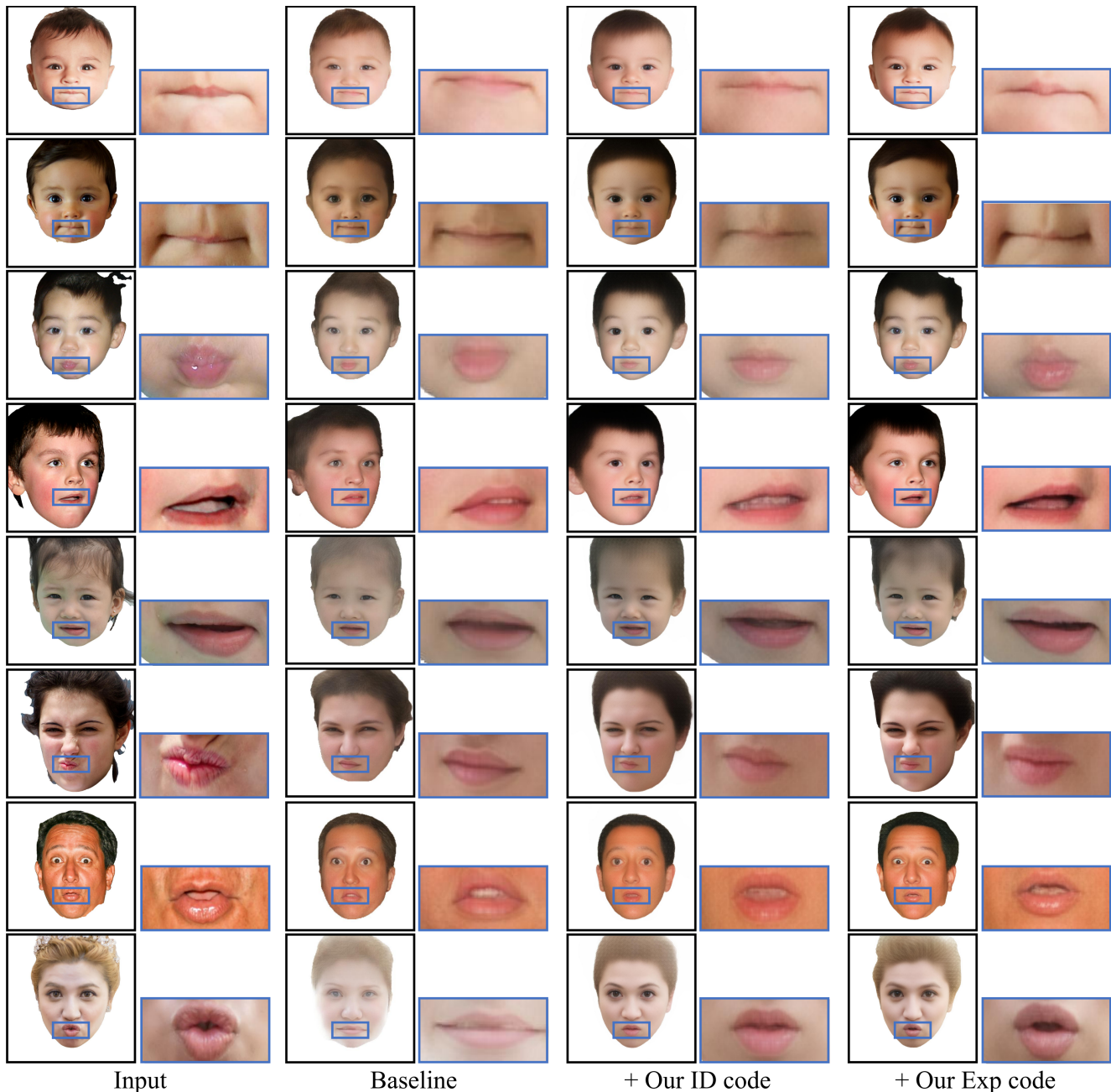| Input | Baseline | + Our ID code | + Our Exp code |

Figure 11. Qualitative ablation study on MetaHead-F model designs for our proposed conditional supervision head signals. Please refer to Sec. 4.2 of the main manuscript for details. Expressions are highlighted in blue. Id: identity, Exp: expression.

This is primarily due to the lack of sufficiently large and diverse labeled training data for the task. Collecting precise and highly varied gaze data with ground truth, particularly outside of the lab, is a challenging task.

We embed the eye gaze angle as the gaze feature into the feature design space of MetaHead-F, and fine-tune the pretrained MetaHead-F on the training images of two widely used gaze datasets **Gaze360** [27] and **MPIIFaceGaze** [52] respectively, which include ground truth gaze labels.

By providing a sequence of head attributes such as identity, illumination, texture, pose and eye gaze angle as conditional signals, we can easily generate heads, which are consistent with the gaze angle label, with diverse shape and appearance variation, covering the diverse challenging scenes.

We quantitatively evaluate whether our generated heads help label estimation with small amount of real data. First, the experiment is conducted on **Gaze360** [27]. We train an appearance-based gaze estimator Pinball LSTM [27] on
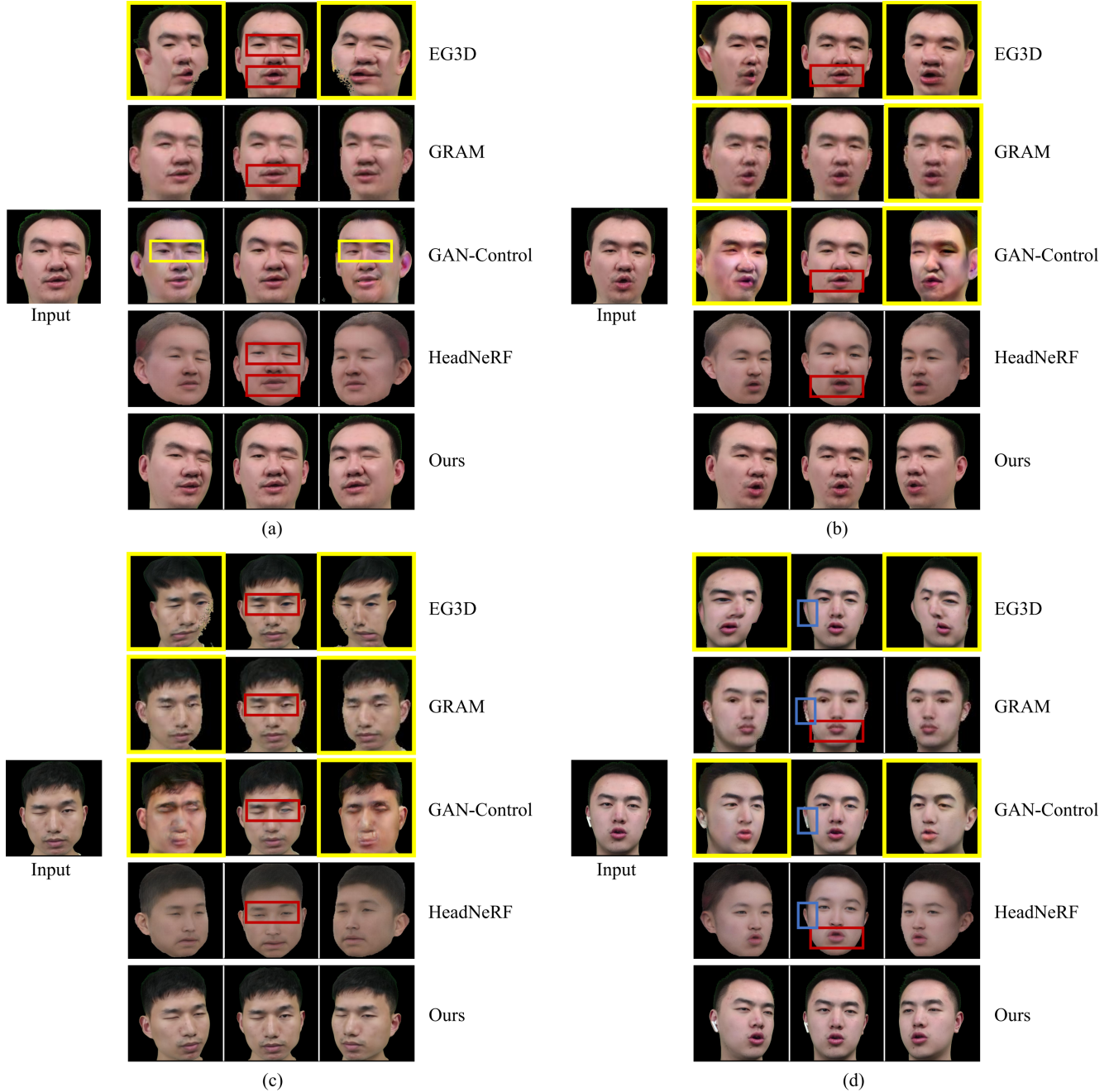
Figure 12. Reconstruction and attributes control accuracy comparison with existing methods. For each sub-figure, column 3 corresponds to the reconstruction results, while columns 2 and 4 correspond to the results obtained by controlling the poses of the reconstructions. Red: expression inaccuracy. Yellow: identity and expression inconsistency. Blue: head accessories inaccuracy. Please refer to Appendix A.2 for details.

image-gaze pairs of training data and test on **Gaze360** [27] testing data. Training data consists of k% of real images and n generated heads.

The estimation accuracy is evaluated by the mean angular error which are provided separately for samples where the subject is looking within 90°(Front 180°) and 20°(Front facing) of the camera direction. Fig. 14 shows that as we

continuously add the synthetic heads, the performance of the estimator keeps improving until saturation. Experiments on **MPIIFaceGaze** [52] also provide similar results(see Fig. 14). This suggests that our synthetic heads indeed capture the correlation between gaze, eye shape, and eye appearance, and are helpful for applications with little real data or less well annotated data.
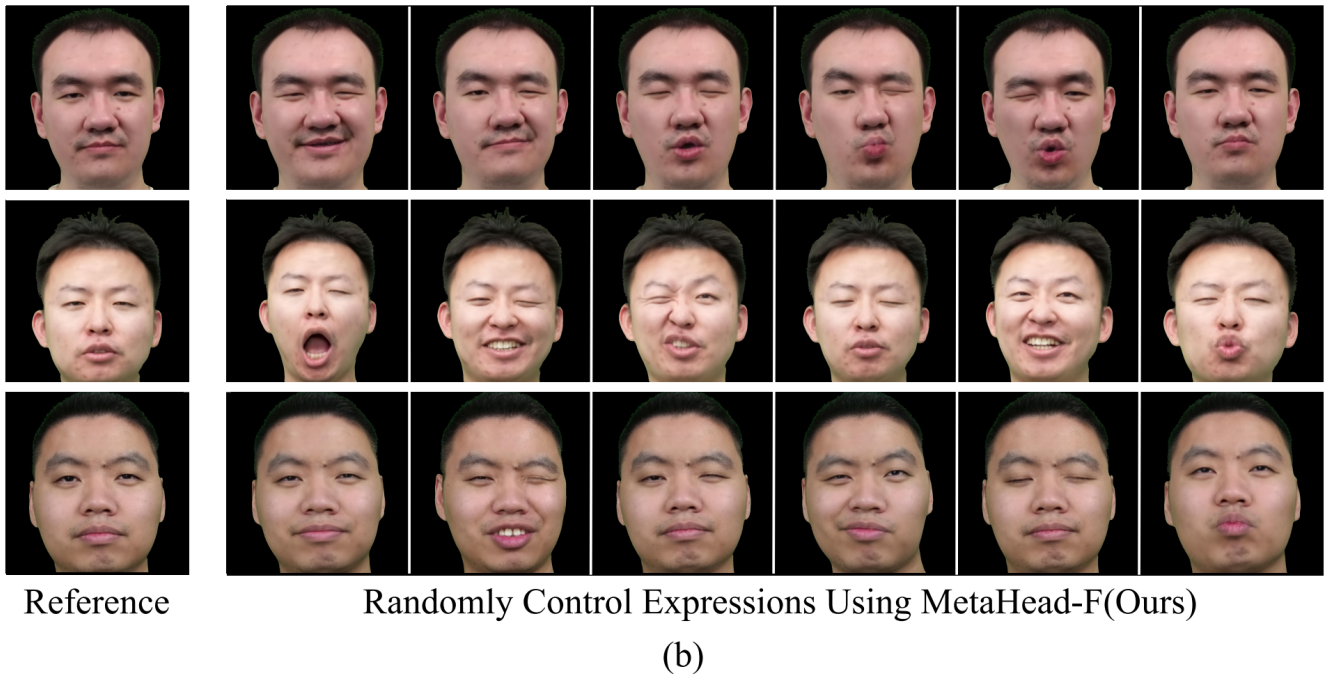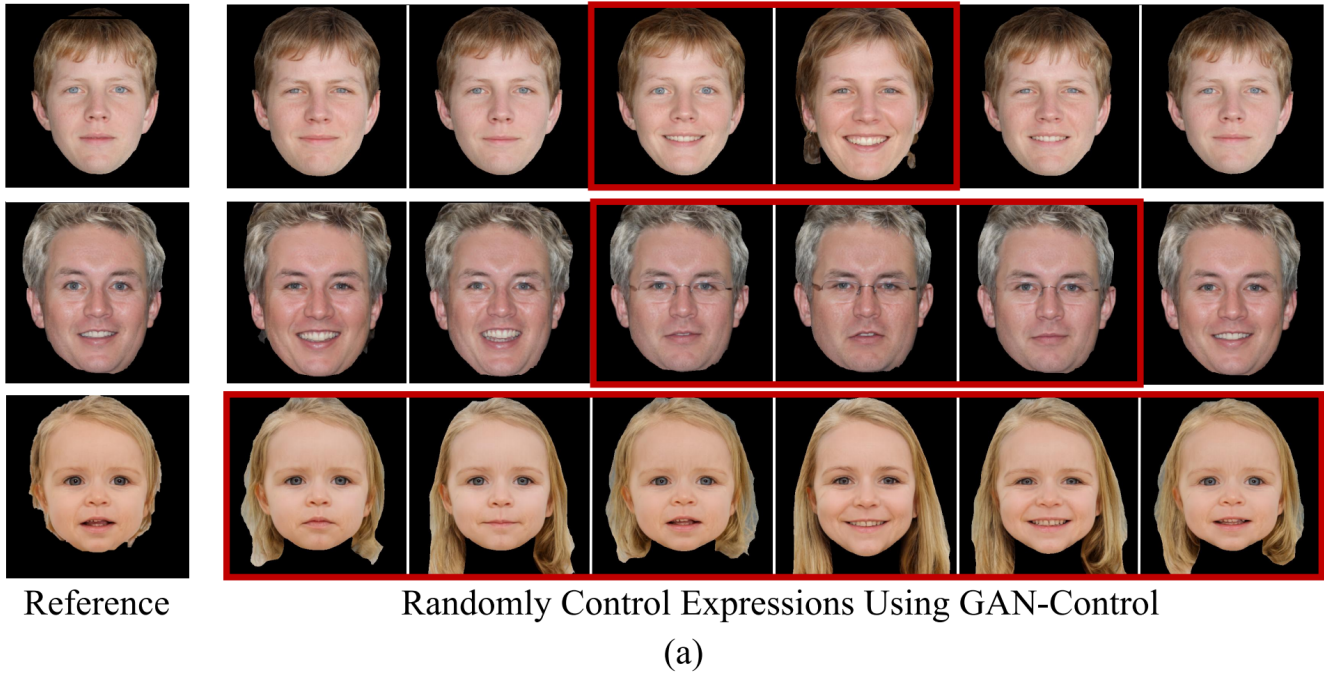
Reference            Randomly Control Expressions Using GAN-Control

(a)



Reference            Randomly Control Expressions Using MetaHead-F(Ours)

(b)

Figure 13. (a)Heads generated by GAN-Control [45] with different expressions. Each row is generated with the same parameters except the expression code. Some examples where GAN-Control produces identity inconsistencies when only the expression is supposed to change are highlighted in red. (b)For a reference head, MetaHead-F randomly control its expression in a highly consistent and photo-realistic manner.

Besides, since we embed the gaze angle feature into the feature design space, MetaHead-F could precisely control eye gaze angle(direction). We demonstrate dynamic control results on eye gaze in the video demo(please see the project page).

### A.5. Additional Qualitative Evaluation on Label-Head

Apart from the landmark and gaze features, we also discuss synthesizing heads with consistent semantic labels and the control of head shape using semantic labels here.
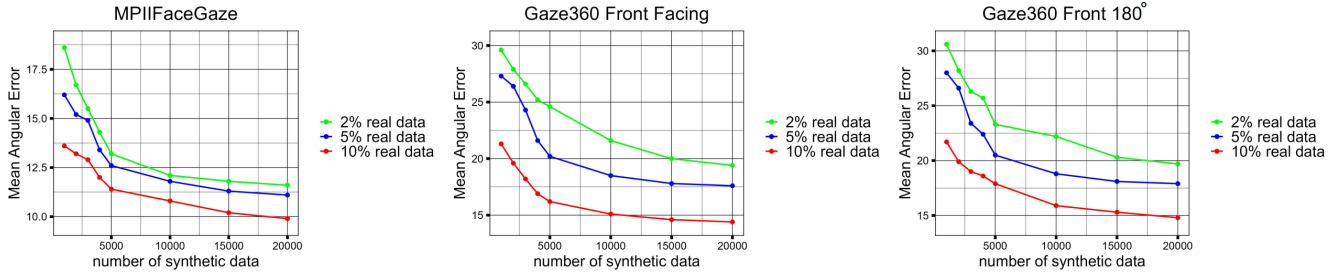
Figure 14. Performance comparison on **Gaze360** [27] and **MPIIFaceGaze** [52] using real and synthetic data. Our synthetic data significantly reduce the gaze estimation error.



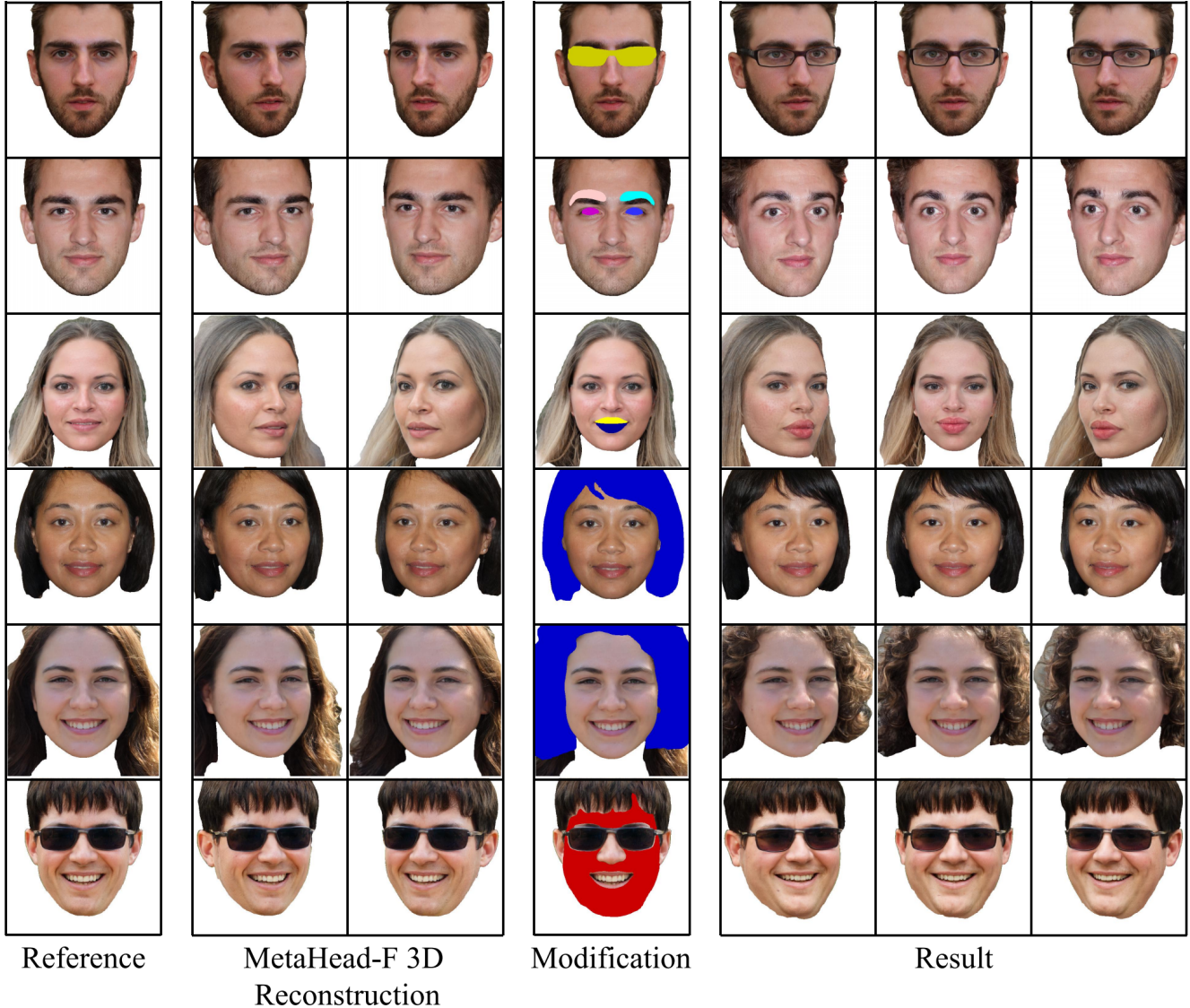| Reference | MetaHead-F 3D Reconstruction | Modification | Result |
|---|---|---|---|

Figure 15. Results of local head editing using semantic labels. LabelHead allows users to perform locally 3D fine-grained head shape control and editing in a disentangled and view-consistent manner. Please refer to Appendix A.5 for details.

**Head Shape Control Using Semantic Labels**    As is classified in CelebAMask-HQ [32], each pixel in the head semantic mask falls into 19 distinct categories including skin, eyebrows, ears, mouth, lip, etc. Given a head image, we flatten the semantic category corresponding to each pixel into a latent representation as the semantic label. We then embed
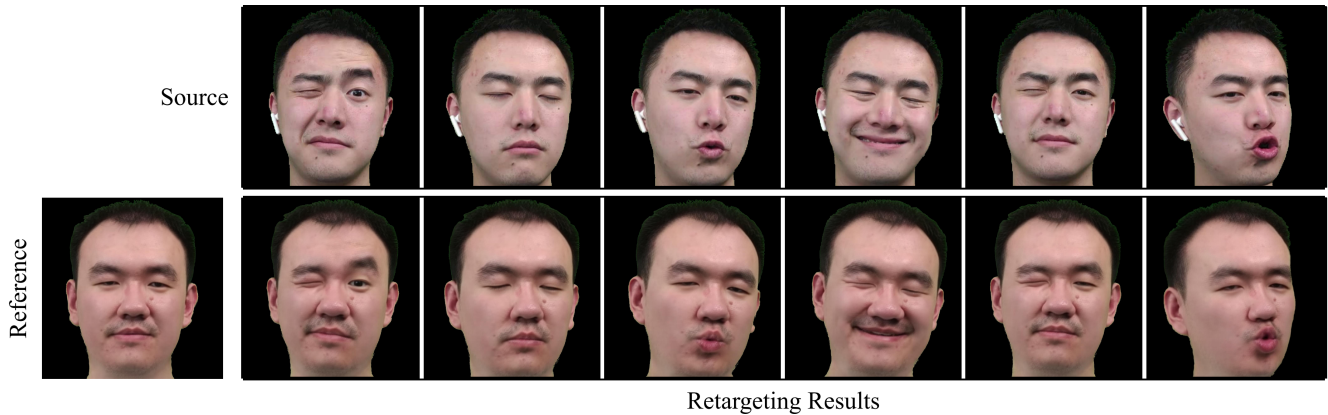
Source

Reference

Retargeting Results

Figure 16. Application of head model MetaHead-F: one-shot facial retargeting.



(a) "a woman with lipstick"   (b) "a man with blonde and curly hair"   (c) "an old man with beard"   (d) "a child with orange hair"

Figure 17. Application of head model MetaHead-F: text-to-head generation.



Reference   MetaHead-F 3D Reconstruction   "Thick Beard"       Reference   MetaHead-F 3D Reconstruction   "Dark Skin"

Reference   MetaHead-F 3D Reconstruction   "Lipstick"       Reference   MetaHead-F 3D Reconstruction   "Curly Hair"

Reference   MetaHead-F 3D Reconstruction   "Child"       Reference   MetaHead-F 3D Reconstruction   "Middle Aged"

Reference   MetaHead-F 3D Reconstruction   "Sunglasses"       Reference   MetaHead-F 3D Reconstruction   "Bigger Eyes and Beard"

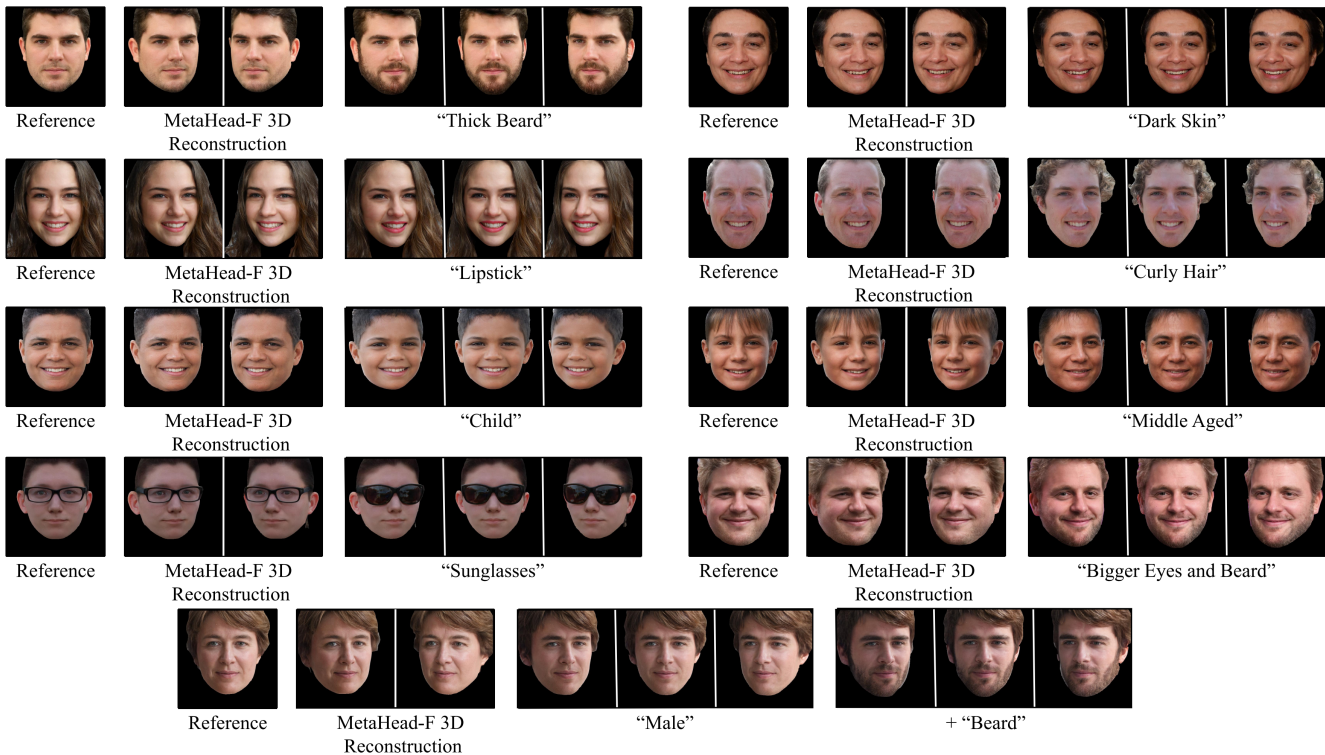Reference   MetaHead-F 3D Reconstruction   "Male"   + "Beard"

Figure 18. Application of head model MetaHead-F: text-based 3D head manipulation.

the semantic label feature into the feature design space of MetaHead-F.

We fine-tune the pre-trained head model MetaHead-F on CelebAMask-HQ [32], in which each image has a seg-

mentation mask of 19-class facial attributes. After training, given the reference head image, we randomly initialize the label values of all features including semantic label, then using the photometric loss between the test image and the generated head of MetaHead-F to optimize the label value, thus getting the semantic mask estimation.

We conduct interactive head shape manipulation by locally editing(drawing) on the obtained semantic mask and leverage MetaHead-F to generate the corresponding modified free-view 3D heads. As is shown in Fig. 15, the head shape control using semantic labels is disentangled and view-consistent.

## B. Additional Applications

**One-shot Facial Retargeting**   Fig. 16 shows the application of MetaHead-F for one-shot facial retargeting. Our expression and identity condition head signals are decoupled and can precisely control the head. Therefore, for a source subject's video or monocular images, we extract and replace the identity signals of the source subject with that of the test reference head. Then, the pre-trained MetaHead-F could perform high-fidelity single-view 3D reconstruction and facial retargeting(reenactment).

**Text-to-head Generation**   A natural way to customize 3D heads is to use language guidance. However, discovering semantically meaningful latent manipulations usually requires painstaking human examination of the many degrees of freedom. We leverage the power of Contrastive Language-Image Pre-training (CLIP) [38] models in order to develop a text-based interface for MetaHead-F head generation and manipulation that does not require such manual effort. Our simple yet efficient approach for leveraging CLIP to guide image generation is to perform the direct latent code optimization with the CLIP loss:

$$\mathcal{L}_{\text{CLIP}}(\mathbf{z}) = D_{\text{CLIP}}(M(\mathbf{z}), t), \tag{7}$$

where $M$ is the pre-trained MetaHead-F, $t$ is the text prompt, and $D_{\text{CLIP}}$ is the cosine distance between its image and text CLIP embedding arguments. As shown in Fig. 17, one can finely customize the super-realistic 3D heads using text descriptions.

**Text-based 3D Head Manipulation**   We can further semantically 3D edit generated heads using text prompts through solving the optimization problem of minimize the cosine distance between the CLIP embeddings of its text and image inputs. Similarity to the input head is enforced by the identity loss:

$$\mathcal{L}_{\text{ID}}(\mathbf{z}) = 1 - \langle R(M(\mathbf{z})), R(M(\mathbf{z}_s)) \rangle, \tag{8}$$

where $\mathbf{z}_s$ is the source latent code, $R$ is an InsightFace [14] face recognition network and $\langle \cdot, \cdot \rangle$ computes the cosine similarity between the arguments. As shown in Fig. 18, we can achieve a wide variety of disentangled, meaningful and view-consistent 3D head control faithful to the text prompt.

## C. Discussion

**Ethics Statement**   Our digital head engine MetaHead focuses on technical development. Our approach can be used to super-realistically generate or reconstruct view-consistent 3D controllable digital heads, generate digital heads consistent with the given customizable feature labels and bidirectionally estimate the labels of features such as landmark coordinates, eye gaze angle, hair color and so on that are difficult to annotate before. However, since our digital head engine could generate heads at a quality that some might find difficult to differentiate from real human heads, we believe that it is essential to develop safeguarding measures to mitigate the potential for misuse.