

# CoNFies: Controllable Neural Face Avatars

Heng Yu<sup>1</sup>, Koichiro Niinuma<sup>2</sup>, László A. Jeni<sup>1</sup>

<sup>1</sup> Carnegie Mellon University, Pittsburgh, PA, USA

<sup>2</sup> Fujitsu Research of America, Pittsburgh, PA, USA

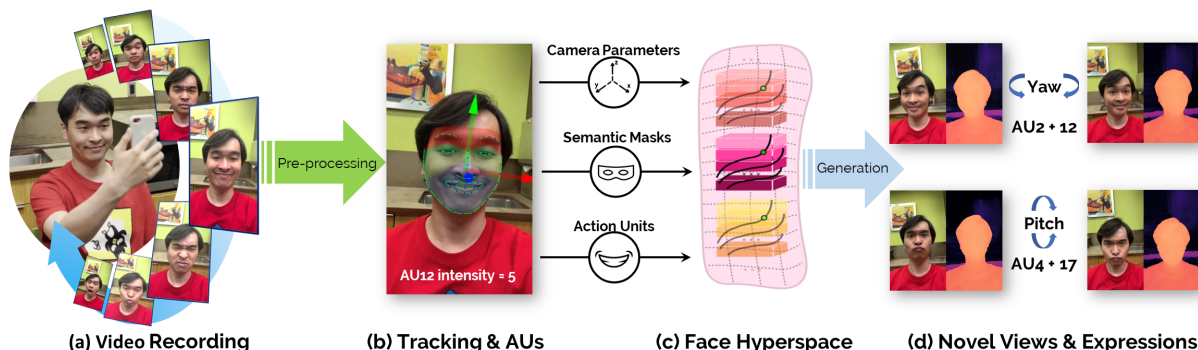


Fig. 1. 2D video of users recording themselves using a circular motion during a semi-structured facial expressions task (a) is processed by a person-independent face tracker that codes facial action units (AU) (b). The estimated camera parameters, semantic masks, action unit intensities, and the original 2D frames are used to build a disentangled face hyper-space (c). From this representation, novel views and unseen expressions can be generated along with their 3D depth (d).

## Abstract

Neural Radiance Fields (NeRF) are compelling techniques for modeling dynamic 3D scenes from 2D image collections. These volumetric representations would be well suited for synthesizing novel facial expressions but for two problems. First, deformable NeRFs are object agnostic and model holistic movement of the scene: they can replay how the motion changes over time, but they cannot alter it in an interpretable way. Second, controllable volumetric representations typically require either time-consuming manual annotations or 3D supervision to provide semantic meaning to the scene. We propose a controllable neural representation for face self-portraits (CoNFies), that solves both of these problems within a common framework, and it can rely on automated processing. We use automated facial action recognition (AFAR) to characterize facial expressions as a combination of action units (AU) and their intensities. AUs provide both the semantic locations and control labels for the system. CoNFies outperformed competing methods for novel view and expression synthesis in terms of visual and anatomic fidelity of expressions.

## I. INTRODUCTION

3D understanding of the world is crucial to the next round of technological innovations in creating digital representations of scenes, objects, and humans. However, the cost of getting 3D supervision for such systems is astronomically higher than those in 2D. Neural volumetric representations, such as Neural Radiance Fields (NeRF) [25], are compelling alternatives for building high fidelity representations from 2D image collections only.

Although, previous work in this direction has demonstrated promising results for modeling and synthesizing novel views of static scenes [44] [12] [9] [41], articulated

objects [38] [31] [13], have explored the use of coarse-grain controls over limited properties, such as color [15], material [46], and object editing [42], relatively neglected is the fine-level control of semantic scene attributes.

Our interest is in general social interactions, and thus we are interested in high-fidelity 3D modeling of facial appearance and dynamics. Previous neural representation based approaches for facial actions either required parametric models to encode facial expressions [13] [3] or were limited in the level of control over scene attributes [30].

In a recent work Kania et al. [17] proposed a controllable neural representation that can achieve simple manipulation, such as opening and closing the mouth and the eyes using a learned mapping between a segmentation mask and a control value that describes the state of that region. The method relies on temporally sparse and manual annotation of the regions of interest along with their control signals, which limits the application of the method.

We wish to make high fidelity 3D reconstruction and control of complex facial movements with a simplified camera setup and little or no manual annotation. A high-level summary of our method is shown on Fig. 1. First, we capture fine-scale transitions of facial movements during a semi-structured expression task. The recorded video is then processed with an automated facial action recognition system [11] [6] that provides anatomically correct action unit (AU) [33] intensities and facial landmarks. From these, semantic facial masks are generated automatically and frames are sub-sampled to build an AU-balanced set of training data. The selected 2D frames, semantic masks, AU intensities, and camera parameters then used to build a face hyper-space, that can be used to synthesize novel views and unseen expression combinations.

Our contributions are as follows:

- **Anatomically Correct Control.** Previous work was limited to holistic deformations or required manual annotation. We achieve an anatomically controllable neural avatar with no manual annotation. We show that this is achievable by using automated facial action coding that provides consistent facial key-points and semantic labels across subjects.
- **Multi-label Semantic Masks.** We achieve a completely disentangled representation in the feature space where the different semantic regions do not affect each other and each region has multiple semantic control variables. We demonstrate that this formulation correctly handles different action unit combinations and achieves better visual fidelity than previous methods.

## II. RELATED WORKS

Our work focuses on automatic control over avatar expressions and is closely related to several computer vision and graphics research domains such as neural rendering and avatar animation.

### A. Neural Rendering and Novel View Synthesis

Implicit neural representations represented by NeRF has become more and more popular recently. NeRF can synthesize high-quality rendering results from novel views and some following extensions further enhanced the algorithm in rendering quality improvement [44] [41], faster training and inference [14] [43] [26], generalization model [35] [27] and so on. NeRF and its variants achieve remarkable performance on static objects, and several of these variants extend it to dynamic scenes, which is the same scenario as ours. Park et al. [29] [30] introduce deformation fields along with a canonical NeRF to learning the movements. Some other works handle the dynamic scene through learning movement offset [32] [38] or scene flow [22] [40]. These methods achieve eye-catching results in dynamic scenes and can achieve the separation of static parts and dynamic parts to some extent through the learned deformation field/offset/flow. However, they are far from the fine-grained control over the dynamic scenes.

### B. Avatar Animation

Avatar animation is a well explored research area. Some works have attempted to manipulate or edit a face [10] [21] [36] while they are mainly image-based and fail to leverage 3d representation. Other works [20] [18] [37] utilize 3D Morphable Model (3DMM) as 3D face representation to achieve the head pose control and image or video reanimation. However, they have limited ability to synthesize novel views since they neglect scene geometry or appearance. Given that high-quality novel view synthesis and fine-grained avatar control is pretty challenging, some works take advantage of neural rendering to model a non-rigid 3d avatar. NerFACE [13] allows face expression and head pose control by modeling a 4D face avatar using neural radiance fields and a facial expression tracking algorithm while it assumes a

static background and fixed camera. Some other works either require professional equipment and training dataset [24] or impose parametric face models such as 3DMM [2] [3], which limit their application scenarios. HyperNeRF [30] introduces hyper space that can better fit dynamic face avatars and also control facial expressions to some extent through hyper space. However, it is far from fine-grained control and cannot achieve per-attribute control. CoNeRF [17] can achieve per-attribute control by imposing an attribute value and mask based on HyperNeRF while it can only achieve simple control over each attribute and different attributes may affect each other. It also requires manual labeling which is labor-intensive. In contrast, we propose an automatic system that enables fine-grained comprehensive control over a face avatar and novel-view synthesis simultaneously without any manual labels. Recently, Cao et al. [8] proposed an approach creating volumetric avatars using only a short phone capture. Though their approach can generate a high-fidelity avatar, a large-scale high resolution multi-view dataset is required to pre-train their model. Unlike their approach, our method requires only a single input video.

## III. METHOD

Our system consists of three parts: (i) data and annotation processing, (ii) network training, and (iii) avatar control. We will describe each component in detail in the following parts.

### A. Data and annotation processing

The data collection process of our system can be done using just a smart phone with a slo-mo video function. After collecting slo-mo video of changing expressions with moving cameras, we apply OpenFace [6] on every frame to detect the facial landmarks [1] [5] and facial action units (AUs) [4]. It is worth noting that other facial landmark/AU detection methods can also work and may have better results, but this is not the focus of our work. OpenFace can output 68 2D landmark locations and 17 AU intensities from 0 to 5 (as shown in Fig. 2). To mitigate the noise impact of AUs detection between adjacent frames, we apply Savitzky–Golay filter [34] to smooth the AUs and then sample the frames to reduce the computational load of the whole system. We found uniform sampling can lead to an extreme imbalance in AU intensity distribution (Fig. 3(a)) since there exist many neutral frames in the dataset. To alleviate this problem, we propose a balanced sampling strategy. We define each AU value and AU intensity pair as a AU-intensity block, which consists of frames with a corresponding AU value and

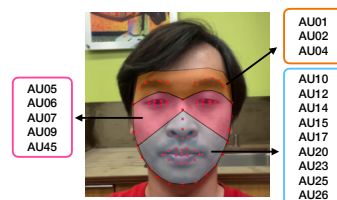


Fig. 2. Facial landmarks, AUs and mask annotations

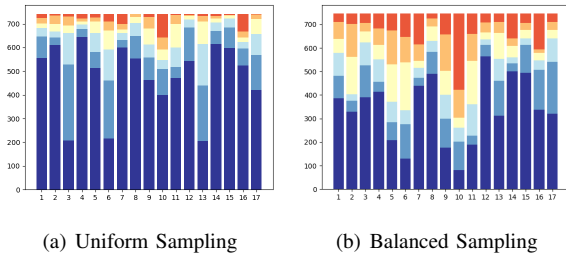


Fig. 3. AU intensity distribution using different sampling strategies

intensity. We ascendingly sort all the AU intensity blocks (total block number equals to intensity number times AU number) according to frame number in each block first. Then we uniformly sample frames from each block in order. Before sampling each block, we remove the frames that are already sampled. Using this strategy, we can get a more balanced sampling result shown in Fig. 3(b).

After obtaining the 2D facial landmarks and AUs, we generate attribute masks and controllable AUs values as annotations. We define three regions and assign each action unit (AU) to its corresponding mask as shown in Fig. 2. We calculate the middle points of eyebrow and eye key-points as the boundary between the first and second region. We also extend some distance along the eyebrow direction as the boundary of the first region to make sure the whole eyebrows are included in the region. The boundary of the third region consists of landmarks (#3 – #13, #28) detected by OpenFace. We also normalize each AU according to:

$$AU' = \min\left(\frac{AU - AU_{min}}{\alpha AU_{max} - AU_{min}} \times 2 - 1, 1\right) \quad (1)$$

where  $AU \in [0, 5]$ ,  $AU' \in [-1, 1]$ ,  $AU_{min}$  and  $AU_{max}$  are minimum and maximum and values for each AU among all frames, respectively, and  $\alpha$  is the factor that adjusts the maximum of AUs and we set  $\alpha$  as 0.8 for all the experiments.

### B. Network architecture

In this section, we briefly introduce NeRF [25], HyperNeRF [30] and CoNeRF [17] for completeness and then describe our approach in detail.

1) *Neural Radiance Field (NeRF)*. NeRF uses a fully-connected neural network to learn the implicit 3D scene volumetric representations through a partial set of 2D images and can generate novel views. The NeRF network takes a sample 3d position  $\mathbf{x} = (x, y, z)$  and a 2d view direction  $\mathbf{d} = (\theta, \phi)$  as input and outputs the emitted color  $\mathbf{c}$  and volume density  $\sigma$  at position  $\mathbf{x}$  with view direction  $\mathbf{d}$ . Then one can accumulate the densities and colors into image pixels  $\mathbf{C}$  in RGB using classical volume rendering techniques [16] as follows:

$$C(\mathbf{r}) = \int_{t_n}^{t_f} T(t) \sigma(\mathbf{r}(t)) \mathbf{c}(\mathbf{r}(t), \mathbf{d}) dt \quad (2)$$

where  $T(t) = \exp(-\int_{t_n}^t \sigma(\mathbf{r}(s)) ds)$ ,  $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$  is the camera ray with near bound  $t_n$  and far bound  $t_f$ .  $C(\mathbf{r})$  is the expected image pixel color of the ray  $\mathbf{r}(t)$ .

2) *HyperNeRF and CoNeRF*. Given that original NeRF can only model static scenes, HyperNeRF, which extend Nerfies [29], is proposed to model dynamic objects, especially face avatars, by introducing canonical hyper-space. The input of HyperNeRF includes sample point  $\mathbf{x}$  and view direction  $\mathbf{d}$ , which is similar to the template NeRF, and also a latent deformation code  $\omega_i$  and a latent appearance code  $\psi_i$ . The sample point  $\mathbf{x}$  concatenated with the image’s latent deformation code  $\omega_i$  are taken as input to the spatial deformation field  $T$  and the ambient slicing surface  $H$  as follows, which yields a warped coordinate  $\mathbf{x}'$  and a coordinate in ambient space  $\mathbf{w}$ , respectively.

$$\mathbf{x}' = T(\mathbf{x}, \omega_i); \quad \mathbf{w} = H(\mathbf{x}, \omega_i) \quad (3)$$

The density  $\sigma$  and color  $\mathbf{c}$  can be then obtained by taking  $\mathbf{x}'$ ,  $\mathbf{w}$  along with direction  $\mathbf{d}$  and latent appearance code  $\psi_i$  as input into template NeRF  $F$ :

$$(\sigma, \mathbf{c}) = F(\mathbf{x}', \mathbf{w}, \mathbf{d}, \psi_i) \quad (4)$$

HyperNeRF can model time-varying shapes even with topological changes and can render different expressions of face avatar by using setting specific ambient coordinates in hyper-space. However, it is far from fine-grained control and fails to achieve per attribute control. Inspired by HyperNeRF, CoNeRF introduces regressors  $A$  and  $M$  to regress the attribute  $\alpha$  and the corresponding mask  $\mathbf{m}$ . The attribute  $\alpha$  is generated from latent deformation code  $\omega_i$  and then is taken as input into ambient slicing surface  $H$  to generate a coordinate in ambient space  $\mathbf{w}$ :

$$\alpha = A(\omega_i); \quad \mathbf{w} = H(\mathbf{x}, \alpha) \quad (5)$$

The corresponding mask map  $\mathbf{m}$  is generated using the warped coordinate  $\mathbf{x}'$  and the ambient space  $\mathbf{w}$  and then is used to mask out  $\mathbf{w}$ :

$$\mathbf{m} = M(\mathbf{x}', \mathbf{w}); \quad \mathbf{w}' = \mathbf{w} \odot \mathbf{m} \quad (6)$$

The following density and color generation is the same as HyperNeRF. CoNeRF maps attribute  $\alpha$  to its corresponding area through mask  $\mathbf{m}$  so as to achieve the control over the corresponding area when rendering from novel views by assigning attribute  $\alpha$  with different values. However, it requires manual labeling of attribute  $\alpha$  and the corresponding mask  $\mathbf{m}$  which is labor-intensive and can only perform simple control over each area since only one attribute is related to each mask. CoNeRF can also lead to movement in another area when controlling one area since its network architecture does not achieve complete decoupling between attributes.

3) *CoNFies*. We propose CoNFies based on CoNeRF that can achieve more complex and accurate control over attributes. The network architecture of CoNFies is shown as Fig. 4. We learn  $K$  attributes  $\alpha_{1...K}$  from latent deformation code  $\omega_i$  using attribute network  $A$ . Different from CoNeRF, we use  $\tanh$  as activation function in the last layer of  $A$  to learn attributes whose range is  $[-1, 1]$ . We adopt the hyper-space  $W$  as proposed in HyperNeRF while our hyper-space

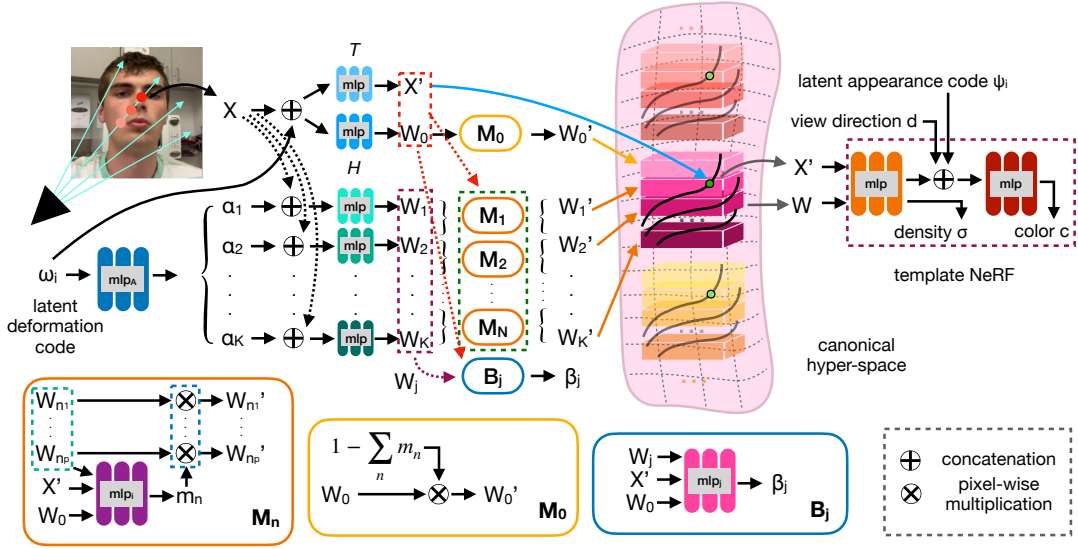


Fig. 4. CoNFies architecture.  $\alpha_i$  represents AU attribute learned from the latent code.  $\beta_i$  is uncertainty and  $m_i$  is mask.

$W$  consists of  $K+1$  ( $K$  attributes and 1 remaining part) components, which are attribute-specific. First,  $K$  original hyper-space components are generated using  $X$  and corresponding attribute  $\alpha$  as following:

$$\mathbf{w}_i = H_i(\mathbf{x}, \alpha_i); \quad i = 1 \dots K \quad (7)$$

We also generate one hyper-space  $\mathbf{w}_0$  for remaining avatar part and generate a warped coordinate  $\mathbf{x}'$  through deformation field  $T$ :

$$\mathbf{w}_0 = H_0(\mathbf{x}, \omega_i); \quad \mathbf{x}' = T(\mathbf{x}, \omega_i) \quad (8)$$

After obtaining the original hyper-space  $\mathbf{w}_{0 \dots K}$ , we then learn  $N$  masks according to a pre-defined correspondence between attributes and masks (many-to-one):

$$\mathbf{m}_n = M_n(\mathbf{x}', \mathbf{w}_0, \mathbf{w}_{n_1} \dots \mathbf{w}_{n_p}); \quad n = 1 \dots N \quad (9)$$

where  $\mathbf{w}_{n_1} \dots \mathbf{w}_{n_p}$  are the hyper-space generated from corresponding attributes  $\alpha_{n_1} \dots \alpha_{n_p}$  that are related to  $\mathbf{m}_n$ . The mask  $\mathbf{m}_0$  for hyper-space  $\mathbf{w}_0$  is  $1 - \sum_n \mathbf{m}_n$ . Final hyper-space  $\mathbf{w}'_{0 \dots K}$  are obtained by masking the original hyper-space using corresponding masks:

$$\mathbf{w}'_i = \mathbf{w}_i \odot \mathbf{m}_j; \quad i = 1 \dots K; \quad j = 1 \dots N \quad (10)$$

where  $\mathbf{m}_j$  is the mask which  $\alpha_i$  is related to. The whole hyper-space  $W$  is obtained by concatenating  $\mathbf{w}'_{0 \dots K}$  and the final density  $\sigma$  and color  $\mathbf{c}$  are obtained using (4), which is the same as template NeRF. We also render the mask field into image space using an analogous volume rendering process:

$$M(\mathbf{r}|\theta, \beta_c) = \int_{t_n}^{t_f} T(t) \sigma(\mathbf{r}(t)) \mathbf{m}(\mathbf{r}(t), d) dt \quad (11)$$

It is also worth noting that we learn an uncertainty  $\beta_i$  for each attribute  $\alpha_i$  as follows to help reduce the potential noises after AUs filtering during training.

$$\beta_i = B_i(\mathbf{x}', \mathbf{w}_0, \mathbf{w}_i); \quad i = 1 \dots K \quad (12)$$

4) *Training Losses.* Given a training set collection of  $C$  images, the losses of our method consist of two parts, reconstruction losses  $L_{rec}$  and control losses  $L_{ctrl}$ , which are similar to [17]:

$$\arg \min_{\theta, \{\mu_c\}} L_{rec}(\theta, \{\mu_c\}) + L_{ctrl}(\theta, \{\mu_c\}) \quad (13)$$

where  $\theta$  is network parameters and  $\mu_c$  represents the latent code (deformation/appearance) of image  $c$ . Reconstruction losses  $L_{rec}$  have two parts ( $L_{recon}$  and  $L_{reg}$ ). One is the primary reconstruction loss, which aims to reconstruct input observations  $\{C_c\}$  as follow (gt = ground truth):

$$L_{recon} = \sum_{\mathbf{r} \in R} \|C(\mathbf{r}|\theta, \beta_c) - C^{gt}(\mathbf{r})\|_2^2 \quad (14)$$

The other one a Gaussian prior on the latent codes  $\{\mu_c\}$  as proposed in [28]:

$$L_{reg} = \sum_c \|\mu_c\|_2^2 \quad (15)$$

Control losses  $L_{ctrl}$  also have two parts: attribute mask loss  $L_{mask}$  and attribute value loss  $L_{attr}$ , as proposed in [17]. For attribute mask loss, we first project 3D volumetric neural mask field  $\mathbf{m}$  into 2D mask image using (11) and the attribute mask loss can be written as:

$$L_{mask} = \sum_{\mathbf{r}, a} \delta_{c,a} CE(M(\mathbf{r}|\theta, \beta_c), M_{c,a}^{gt}(\mathbf{r})) \quad (16)$$

where  $CE(\cdot, \cdot)$  represents cross entropy and  $M_{c,a}^{gt}(\mathbf{r})$  is  $a$ -th attribute in the  $c$ -th image.  $\delta_{c,a}$  denotes an indicator, where  $\delta_{c,a} = 1$  means attribute  $a$  for image  $c$ , which  $\mathbf{r}$  is belong to, is provided, otherwise  $\delta_{c,a} = 0$ . We also stop gradients in (16) w.r.t  $\sigma$  and employ focal loss [23] in place of the standard cross entropy loss as in [17]. For attribute value loss, we employ the AUs after filtering as ground-truth and  $\beta_{c,a}$  learned from (12) to further reduce noises in AUs:

$$L_{attr} = \sum_c \sum_a \delta_{c,a} \frac{|\alpha_{c,a} - \alpha_{c,a}^{gt}|^2}{2\beta_{c,a}^2} + \frac{(\log \beta_{c,a})^2}{2} \quad (17)$$

where larger  $\beta_{c,a}$  values attenuate the importance of learned  $\alpha_{c,a}$  and the second term precludes the trivial minimum at  $\beta_{c,a} = \infty$ . Hence the network can better learn to adjust  $\alpha_{c,a}$  and reduce the negative effect of noises in AUs. The final loss is:

$$L = L_{recon} + w_{reg}L_{reg} + w_{mask}L_{mask} + w_{attr}L_{attr} \quad (18)$$

where  $w_{reg}$ ,  $w_{mask}$  and  $w_{attr}$  are weighting coefficients.

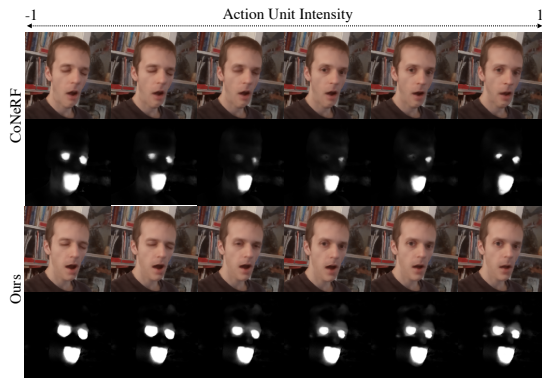
## IV. EXPERIMENTS

### A. Implementation details

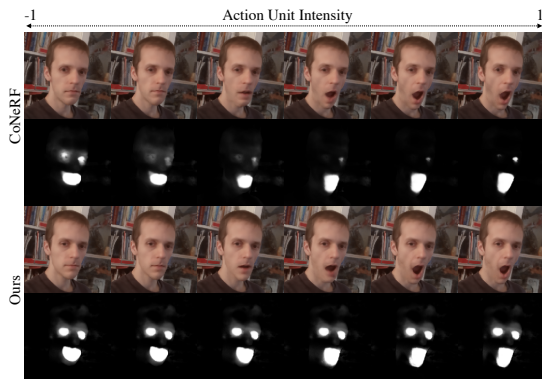
Our method is based on the JAX [7] implementation of CoNeRF [17]. Attribute network  $A$  has six layers, each of which is a 32 neuron multi-layer perceptron (MLP) and has a skip connection at the fifth layer following [29] [30] [17]. Deformation field  $T$  and ambient slicing surface  $H_i$  have the same architecture of those in [30] [17]. Mask network  $M_i$  and uncertainty network  $B_i$  have the same structure, which is a four-layer MLP with 128 neuron per layer and followed by an additional 64 neuron layer with a skip connection as in [17]. The template NeRF is the same as original NeRF [25] but with a different input dimension size. In our case, the number of attributes  $K$  is 17, which is the number of AUs and the number of mask  $N$  is 3 as shown in Fig. 2. We resize all the input images to  $480 \times 270$  and train our NeRF model for 250k iterations with 128 samples per ray and a batch size of 512 rays. We use Adam [19] with initial learning rate  $lr = 1e-4$  and set  $w_{reg} = 1e-4$ ,  $w_{mask} = 1e-2$  and  $w_{attr} = 0.1$  for all experiments. We introduce exponentially decaying on  $lr$  and  $w_{attr}$ , which decay to  $1e-5$  and 0, respectively. Exponentially decaying on  $w_{attr}$  can help further reduce the noises introduced by AU detection. We train our model on a NVIDIA A100 GPU with 80G memory and the whole process takes around 26 hours

### B. Dataset

We used both the real dataset provided in [17] and we collected our own video sequences using a smartphone. In our data collection each of the sequences was captured with an Apple iPhone 13 Pro using 120 fps slo-mo mode and is about 2 minutes long. In each video sequence, the person performs different facial expressions related to a single AU one by one first and then performs arbitrary facial expressions related to multiple AUs. Each video is extracted to frames with 120 fps and we perform OpenFace [6] and smoothing as mentioned above. Note that OpenFace can only provide 17 AU intensities and we use these in the following experiments. More AU intensities may be obtained from other methods, which is not the focus of our paper. We then undersample the sequences to give approximately 750 frames per capture. The frames along with the AUs and attribute masks generated automatically form the datasets we use in our experiments.



(a) eye control



(b) mouth control

Fig. 5. Controlling results of CoNFies and CoNeRF. Our CoNFies can perform a better control over one attribute without affecting other attributes.

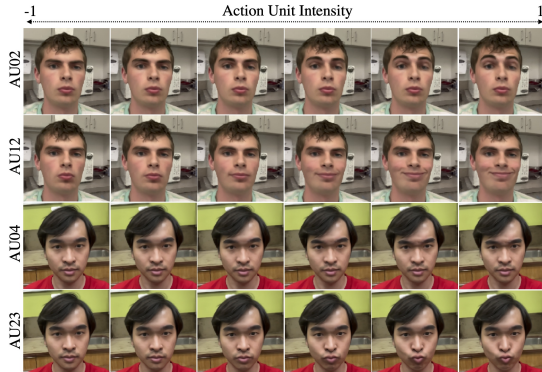


Fig. 6. Control using single AU. AU02 is outer brow raiser. AU04 is brow lowerer. AU12 is lip corner puller. AU23 is lip tightener.

### C. Decoupling Mask

We compare our method with CoNeRF using their dataset with manual attribute values and mask area labels to show the effectiveness of our decoupling mask structure. We control the eyes and mouth separately using attribute values and show the rendering image results along with the masks in Fig. 5. We can see from the mask results that when controlling one attribute, the masks of the other attribute are not affected in our method. But one attribute can affect the others in CoNeRF, which lead to unexpected movement and

TABLE I  
INTRACLASS CORRELATION (ICC) COMPARISON BETWEEN CoNeRF  
AND OUR METHOD.

AU	CoNeRF	Ours	AU	CoNeRF	Ours
01	0.52	0.86	14	0.17	0.77
02	0.54	0.73	15	-0.15	0.81
04	0.23	0.91	17	0.31	0.88
05	-0.11	0.40	20	0.03	0.63
06	0.00	0.00	23	0.49	0.82
07	-0.44	0.73	25	0.36	0.90
09	0.55	0.74	26	0.08	0.93
10	0.00	0.00	45	0.21	0.83
12	0.00	0.87	mean	0.16	0.69

artifacts in the rendering results.

#### D. Attribute Control

Our CoNFies can achieve attribute control using different AUs. We first show the controlling result using single AU (AU 02, 04, 12, 23) on two sequences as in Fig. 6. We also conducted quantitative evaluation to compare our method and CoNeRF.

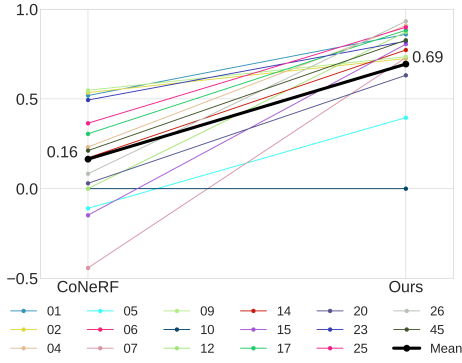


Fig. 7. Intra-class Correlation (ICC) comparison between CoNeRF and our method.

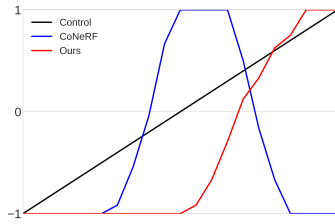


Fig. 8. Comparison of AU01 intensity transition between control values and synthesized images.

In this experiment, AU intensities were obtained from synthesis images generated by our method and CoNeRF, and compare them using Intra-class Correlation (ICC). OpenFace was used to obtain AU intensities from synthesis images. Only images with extreme AU values (near -1 or 1) were manually selected to train CoNeRF while our method automatically obtained images to train the model. Table I and Fig. 7 show that our method outperforms CoNeRF. Fig. 8 compares AU01 intensity transition between control values and synthesized images. The AU intensity transition of our

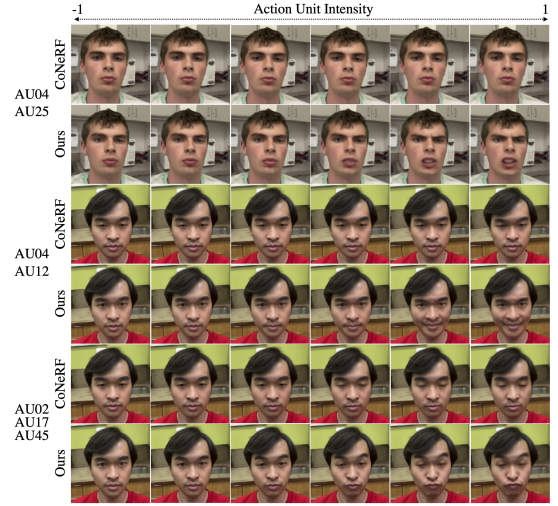


Fig. 9. Control using multiple AUs on different regions. AU02 is outer brow raiser. AU04 is brow lowerer. AU12 is lip corner puller. AU17 is chin raiser. AU25 is lips part. AU45 is blink.

method is much closer to the control than CoNeRF. This result also indicates that our method can handle subtle AU changes better than CoNeRF.



Fig. 10. Control using multiple AUs on the same region. AU12 is lip corner puller. AU25 is lips part.

Our method can perform complex control using multiple AUs simultaneously as shown in Fig. 9 and Fig. 10. From Fig. 9, we can see that our method can perform combined AUs control over different regions, e.g. eyebrow and mouth.

We show that our method can even perform more complicated control over the same region. As shown in Fig. 10, we can control the smile while keeping mouth open or control mouth open while keeping the smile. However, CoNeRF cannot render good results under such complex scenarios. More visualization results can be found in the supplementary material.

#### E. Novel View Synthesis

As a NeRF-based method, we can synthesize novel views with single or multiple AUs control. We show the novel view synthesis results along with corresponding masks in Fig. 11, where we keep the facial expression constant. To further show the power of our method, we show the rendering result

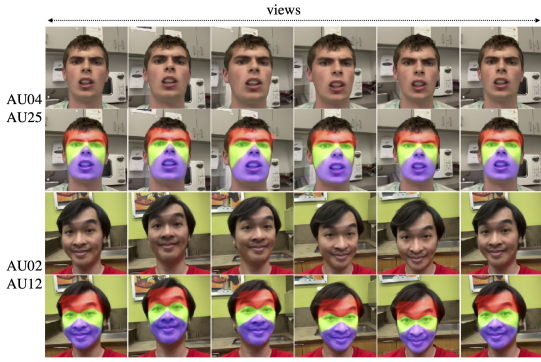


Fig. 11. Novel view synthesis under fixed AU setting. AU02 is outer brow raiser. AU04 is brow lowerer. AU12 is lip corner puller. AU25 is lips part.

under different views and control the AU values (AU02 and AU12) simultaneously in Fig. 12.

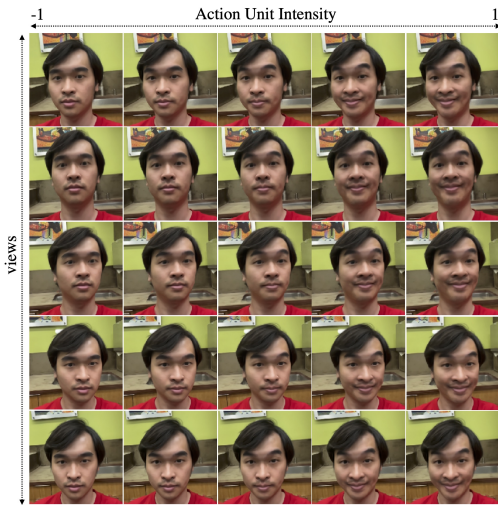


Fig. 12. Rendering results under different AU (AU02 and AU12) intensities and views

We also evaluate the rendering quality of our method on a frame interpolation task as proposed in [17] using the dataset it released. We interpolate every other frame and do not perform any attribute control. We use Peak Signal-to-Noise Ratio (PSNR), Multi-scale Structural Similarity (MS-SSIM) [39] and Learned Perceptual Image Patch Similarity (LPIPS) [45] to quantitatively evaluate our method compared with NeRF [25], NeRF + Latent, Nerfies [29], HyperNeRF [30], CoNeRF-*M* and CoNeRF [17], which is consistent with [17]. The result is shown in Table II, where we can see that our method can achieve comparable performance in the rendering quality from a novel view.

### F. Facial Expression Transfer

In addition to adjusting AU values manually to control the avatar’s expression, it is possible to copy them from another person’s face. In this case, first we detect the AU intensities from the source person’s face and use them to re-synthesize the same expression on the avatar. We show the rendering results using another face sequence to control the trained sequence in Fig. 13.

TABLE II  
QUANTITATIVE RESULTS

Method	PSNR $\uparrow$	MS-SSIM $\uparrow$	LPIPS $\downarrow$
NeRF	28.795	0.951	0.210
NeRF + Latent	32.653	0.981	0.182
NeRFies	32.274	0.981	0.180
HyperNeRF	32.520	0.981	0.169
CoNeRF- <i>M</i>	32.061	0.979	0.167
CoNeRF	32.342	0.981	0.168
Ours	32.356	0.982	0.166



Fig. 13. Facial expression transfer using reference sequence.

## V. CONCLUSIONS

We have proposed an automated approach for controllable neural face avatars. Once a 2D video is recorded using slo-mo mode, our network is automatically trained by utilizing AU intensities and facial landmarks. We have introduced the decoupling mask structure so that the different semantic regions do not affect each other and each region have multiple control variables. We have shown that our approach outperforms CoNeRF in terms of fine-grained AU control. Moreover, our approach has a capability to control multiple AUs and novel views simultaneously.

## VI. ACKNOWLEDGMENTS

This research was supported by Fujitsu. We thank Joel Julin from University of Pittsburgh for helping with data collection and comments that greatly improved the manuscript. We thank Xuxin Cheng from Carnegie Mellon University who provided the data acquisition equipment. We would also like to show our gratitude to Nian-Hsuan Tsai from Carnegie Mellon University who helped build the OpenFace system.

## REFERENCES

- [1] Z. Amir, B. Tadas, and M. Louis-Philippe. Convolutional experts constrained local model for facial landmark detection. *Proceedings of the IEEE CVPRW*, pages 2051–2059, 2017.
- [2] S. Athar, Z. Shu, and D. Samaras. Flame-in-nerf: Neural control of radiance fields for free view face animation. *arXiv preprint arXiv:2108.04913*, 2021.
- [3] S. Athar, Z. Xu, K. Sunkavalli, E. Shechtman, and Z. Shu. Rignerf: Fully controllable neural 3d portraits. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20364–20373, 2022.
- [4] T. Baltrušaitis, M. Mahmoud, and P. Robinson. Cross-dataset learning and person-specific normalisation for automatic action unit detection. In *2015 11th IEEE International Conference and Workshops on Automatic Face and Gesture Recognition (FG)*, volume 6, pages 1–6. IEEE, 2015.

- [5] T. Baltrusaitis, P. Robinson, and L.-P. Morency. Constrained local neural fields for robust facial landmark detection in the wild. In *Proceedings of the IEEE international conference on computer vision workshops*, pages 354–361, 2013.
- [6] T. Baltrusaitis, A. Zadeh, Y. C. Lim, and L.-P. Morency. Openface 2.0: Facial behavior analysis toolkit. In *2018 13th IEEE international conference on automatic face & gesture recognition (FG 2018)*, pages 59–66. IEEE, 2018.
- [7] J. Bradbury, R. Frostig, P. Hawkins, M. J. Johnson, C. Leary, D. Maclaurin, G. Necula, A. Paszke, J. VanderPlas, S. Wanderman-Milne, and Q. Zhang. JAX: composable transformations of Python+NumPy programs, 2018.
- [8] C. Cao, T. Simon, J. K. Kim, G. Schwartz, M. Zollhoefer, S. Saito, S. Lombardi, S.-E. Wei, D. Belko, S.-I. Yu, Y. Sheikh, and J. Saragih. Authentic volumetric avatars from a phone scan. *ACM Transactions on Graphics (TOG)*, 2022.
- [9] A. Chen, Z. Xu, A. Geiger, J. Yu, and H. Su. Tensorf: Tensorial radiance fields. *arXiv preprint arXiv:2203.09517*, 2022.
- [10] Y. Deng, J. Yang, D. Chen, F. Wen, and X. 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.
- [11] I. O. Ertugrul, L. A. Jeni, W. Ding, and J. F. Cohn. Afar: A deep learning based tool for automated facial affect recognition. In *2019 14th IEEE International Conference on Automatic Face & Gesture Recognition (FG 2019)*, pages 1–1. IEEE, 2019.
- [12] S. Fridovich-Keil, A. Yu, M. Tancik, Q. Chen, B. Recht, and A. Kanazawa. Plenoxels: Radiance fields without neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5501–5510, 2022.
- [13] G. Gafni, J. Thies, M. Zollhofer, and M. Nießner. Dynamic neural radiance fields for monocular 4d facial avatar reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8649–8658, 2021.
- [14] S. J. Garbin, M. Kowalski, M. Johnson, J. Shotton, and J. Valentin. Fastnerf: High-fidelity neural rendering at 200fps. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14346–14355, 2021.
- [15] W. Jang and L. Agapito. Codenerf: Disentangled neural radiance fields for object categories. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12949–12958, 2021.
- [16] J. T. Kajiya and B. P. Von Herzen. Ray tracing volume densities. *ACM SIGGRAPH computer graphics*, 18(3):165–174, 1984.
- [17] K. Kania, K. M. Yi, M. Kowalski, T. Trzcinski, and A. Tagliasacchi. Conerf: Controllable neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18623–18632, 2022.
- [18] H. Kim, P. Garrido, A. Tewari, W. Xu, J. Thies, M. Niessner, P. Pérez, C. Richardt, M. Zollhöfer, and C. Theobalt. Deep video portraits. *ACM Transactions on Graphics (TOG)*, 37(4):1–14, 2018.
- [19] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [20] M. R. Koujan, M. C. Doukas, A. Roussos, and S. Zafeiriou. Head2head: Video-based neural head synthesis. In *2020 15th IEEE International Conference on Automatic Face and Gesture Recognition (FG 2020)*, pages 16–23. IEEE, 2020.
- [21] M. Kowalski, S. J. Garbin, V. Estellers, T. Baltrušaitis, M. Johnson, and J. Shotton. Config: Controllable neural face image generation. In *European Conference on Computer Vision*, pages 299–315. Springer, 2020.
- [22] Z. Li, S. Niklaus, N. Snavely, and O. Wang. Neural scene flow fields for space-time view synthesis of dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6498–6508, 2021.
- [23] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 2980–2988, 2017.
- [24] S. Ma, T. Simon, J. Saragih, D. Wang, Y. Li, F. De La Torre, and Y. Sheikh. Pixel codec avatars. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 64–73, 2021.
- [25] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *European conference on computer vision*, pages 405–421. Springer, 2020.
- [26] T. Müller, A. Evans, C. Schied, and A. Keller. Instant neural graphics primitives with a multiresolution hash encoding. *arXiv preprint arXiv:2201.05989*, 2022.
- [27] M. Niemeyer and A. Geiger. Giraffe: Representing scenes as compositional generative neural feature fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11453–11464, 2021.
- [28] J. J. Park, P. Florence, J. Straub, R. Newcombe, and S. Lovegrove. Deepsdf: Learning continuous signed distance functions for shape representation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 165–174, 2019.
- [29] K. Park, U. Sinha, J. T. Barron, S. Bouaziz, D. B. Goldman, S. M. Seitz, and R. Martin-Brualla. Nerfies: Deformable neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5865–5874, 2021.
- [30] K. Park, U. Sinha, P. Hedman, J. T. Barron, S. Bouaziz, D. B. Goldman, R. Martin-Brualla, and S. M. Seitz. Hypernerf: A higher-dimensional representation for topologically varying neural radiance fields. *arXiv preprint arXiv:2106.13228*, 2021.
- [31] A. Pumarola, E. Corona, G. Pons-Moll, and F. Moreno-Noguer. D-NeRF: Neural radiance fields for dynamic scenes. <https://arxiv.org/abs/2011.13961>, 2020.
- [32] A. Pumarola, E. Corona, G. Pons-Moll, and F. Moreno-Noguer. D-nerf: Neural radiance fields for dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10318–10327, 2021.
- [33] E. L. Rosenberg and P. Ekman. *What the face reveals: Basic and applied studies of spontaneous expression using the Facial Action Coding System (FACS)*. Oxford University Press, 2020.
- [34] A. Savitzky and M. J. Golay. Smoothing and differentiation of data by simplified least squares procedures. *Analytical chemistry*, 36(8):1627–1639, 1964.
- [35] K. Schwarz, Y. Liao, M. Niemeyer, and A. Geiger. Graf: Generative radiance fields for 3d-aware image synthesis. *Advances in Neural Information Processing Systems*, 33:20154–20166, 2020.
- [36] A. Tewari, M. Elgharib, F. Bernard, H.-P. Seidel, P. Pérez, M. Zollhöfer, and C. Theobalt. Pie: Portrait image embedding for semantic control. *ACM Transactions on Graphics (TOG)*, 39(6):1–14, 2020.
- [37] A. Tewari, M. Elgharib, G. Bharaj, F. Bernard, H.-P. Seidel, P. Pérez, M. Zollhofer, and C. Theobalt. Stylerig: Rigging stylegan for 3d control over portrait images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6142–6151, 2020.
- [38] E. Tretschk, A. Tewari, V. Golyanik, M. Zollhöfer, C. Lassner, and C. Theobalt. Non-rigid neural radiance fields: Reconstruction and novel view synthesis of a dynamic scene from monocular video. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12959–12970, 2021.
- [39] Z. Wang, E. P. Simoncelli, and A. C. Bovik. Multiscale structural similarity for image quality assessment. In *The Thirty-Seventh Asilomar Conference on Signals, Systems & Computers, 2003*, volume 2, pages 1398–1402. Ieee, 2003.
- [40] W. Xian, J.-B. Huang, J. Kopf, and C. Kim. Space-time neural irradiance fields for free-viewpoint video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9421–9431, 2021.
- [41] Q. Xu, Z. Xu, J. Phillip, S. Bi, Z. Shu, K. Sunkavalli, and U. Neumann. Point-nerf: Point-based neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5438–5448, 2022.
- [42] B. Yang, Y. Zhang, Y. Xu, Y. Li, H. Zhou, H. Bao, G. Zhang, and Z. Cui. Learning object-compositional neural radiance field for editable scene rendering. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 13779–13788, 2021.
- [43] A. Yu, R. Li, M. Tancik, H. Li, R. Ng, and A. Kanazawa. Plenotrees for real-time rendering of neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5752–5761, 2021.
- [44] K. Zhang, G. Riegler, N. Snavely, and V. Koltun. Nerf++: Analyzing and improving neural radiance fields. *arXiv preprint arXiv:2010.07492*, 2020.
- [45] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. 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.
- [46] X. Zhang, P. P. Srinivasan, B. Deng, P. Debevec, W. T. Freeman, and J. T. Barron. Nerfactor: Neural factorization of shape and reflectance under an unknown illumination. *ACM Transactions on Graphics (TOG)*, 40(6):1–18, 2021.