

# EMOPortraits: Emotion-enhanced Multimodal One-shot Head Avatars

Nikita Drobyshev

Antoni Bigata Casademunt  
Stavros PetridisKonstantinos Vougioukas  
Maja Pantic

Zoe Landgraf

{n.drobyshev, a.bigata-casademunt22, k.vougioukas,  
zoe.landgraf15, stavros.petridis04, m.pantic}@imperial.ac.uk

Imperial College London

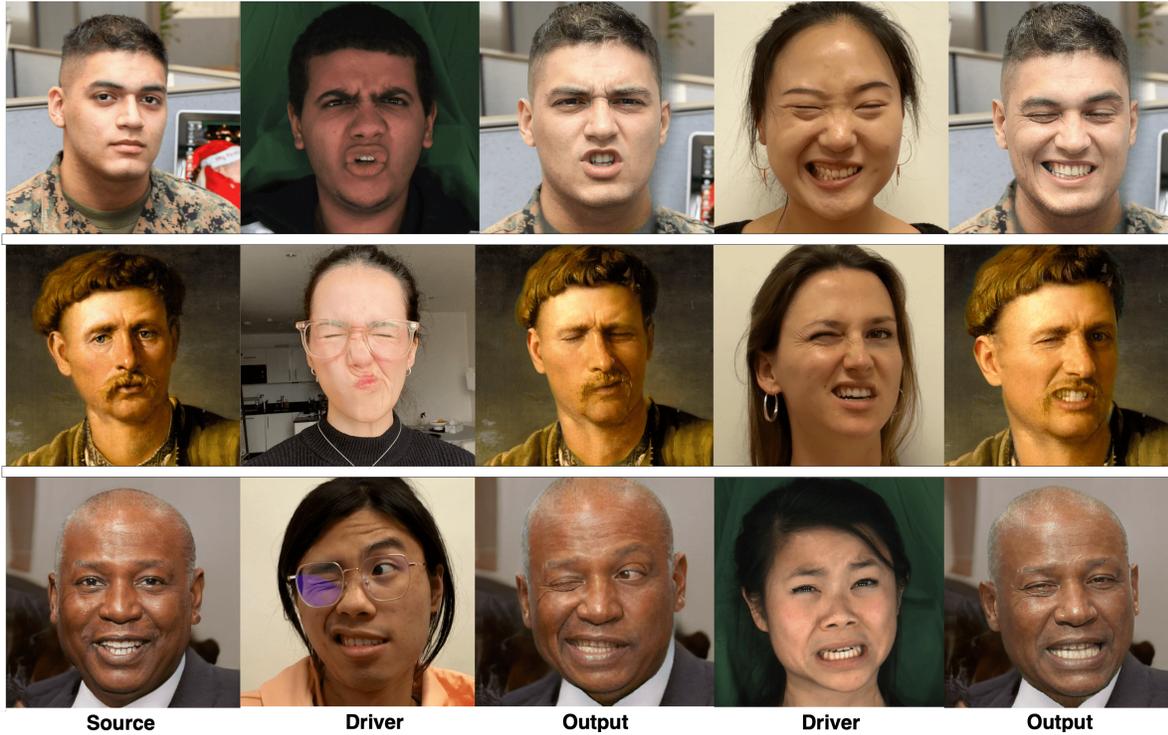


Figure 1. Selected animation results for image-driven mode of our method.

## Abstract

Head avatars animated by visual signals have gained popularity, particularly in cross-driving synthesis where the driver differs from the animated character, a challenging but highly practical approach. The recently presented *MegaPortraits* model has demonstrated state-of-the-art results in this domain. We conduct a deep examination and evaluation of this model, with a particular focus on its latent space for facial expression descriptors, and uncover several limitations with its ability to express intense face motions. To address these limitations, we propose substantial changes in both training pipeline and model architecture, to introduce our *EMOPortraits* model, where we:

*Enhance the model’s capability to faithfully support intense, asymmetric face expressions, setting a new state-of-the-art result in the emotion transfer task, surpassing previous methods in both metrics and quality.*

*Incorporate speech-driven mode to our model, achieving top-tier performance in audio-driven facial animation, making it possible to drive source identity through diverse modalities, including visual signal, audio, or a blend of both.*

*We propose a novel multi-view video dataset featuring a wide range of intense and asymmetric facial expressions, filling the gap with absence of such data in existing datasets. For dataset and video examples please visit [project page](#).*

# 1. Introduction

## 1.1. Representation of facial expression

Advancements in neural head technologies now enable the creation of realistic avatars from a few or even one image, with the latter being crucial when only one source image is available. Cross-driving synthesis, where avatars are animated with different identities is a key technique for virtual reality, filmmaking, photo animation, etc. However, accurate transfer of facial expressions, especially those that are intense and uneven, remains a substantial challenge in avatar animation, particularly for cross-driving synthesis. While previous research focused mainly on preserving identity and transferring moderate facial motions, our work also seeks to accurately drive high-intensity and asymmetric expressions.

Our research is built on and extends the MegaPortraits [9] model, notable for state-of-the-art results in cross-driving synthesis. An additional advantage using this method is that, unlike many avatar systems that depend on limited [2, 31] predefined motion descriptors, MegaPortraits, learns expression representations from scratch, allowing for greater adaptability to a wider range of expressions. We investigate the latent expression spaces and training methods of MegaPortraits to enhance its ability to depict a broad spectrum of facial expressions. Our comprehensive analysis reveals that, while the original model shows limited effectiveness in intense motion representing, it has significant potential for improvement through targeted architectural modifications, adjustments in the training approach, and the integration of our novel dataset.

## 1.2. Integrating speech driving

We integrate speech in our model to be used either as a complementary or an alternative driver, a crucial aspect for applications like virtual assistance or mixed reality when a primary visual signal is absent.

By improving disentanglement of facial expression latent space, we remove disturbing information from face motion descriptors and emphasize components that solely control lip movements. Based on this, we formulated a novel loss that helps to achieve desirable results. Furthermore, our method can generate plausible head rotations and blinks, which appear natural and enhance its applicability across various tasks. Thus, our final method can drive source identity through image, audio, or a mix of them.

## 1.3. FEED: Facial Extreme Emotions Dataset

Public data scarcity is a common obstacle in deep learning, the topic of human facial expressions and movements is not an exception. In particular, we note a lack of high quality video data capturing a wide range of facial expressions, with existing datasets not going beyond basic facial actions shown at Fig. 2. To address this, we collect a new dataset that

includes basic expressions and also captures complex movements like blinks, winks, and head and tongue movements, along with varied extreme expressions which are difficult or impossible to categorize through basic facial actions. Given this variety in expressions in our dataset, we believe that it will be a valuable resource for research in human emotions and facial reconstruction fields. In summary, our main contributions are as follows:

- We introduce our new **EMOPortraits** model for one-shot head avatars synthesis, that is capable of transferring intense facial expression, showing *state-of-the-art* results in cross-driving synthesis. We achieve this through careful latent facial expression space development as well as novel losses and a minimal amount of domain-specific data ( $\sim 0.1\%$  of the train set).
- We integrate a speech-driving mode in our model, that demonstrates cutting-edge results in speech-driven animations. It functions effectively alongside visual signals or independently, also generating realistic head rotations and eye blinks.
- We present a unique multi-view dataset that spans a broad spectrum of extreme facial expressions, filling the gap of absence of such data in existing datasets.

## 2. Related Work

### 2.1. Visual-driven head avatars

In recent years, the domain of neural avatars has branched into two prominent sub-fields: person-specific and person-agnostic avatars. While person-specific methodologies [11, 12, 36, 45, 46] excel in delivering stunning realism and motion fidelity for a particular individual, they encounter challenges representing an arbitrary person. Moreover, these approach demands multiple frames of training data, distinct training for each avatar, and can struggle to replicate motions not encountered during training.

Person-agnostic methods, on the other hand, don't need train or fine-tuning for each new person and present an alternative approach to talking-head synthesis. Earlier works in this domain generated avatars in a few-shot technique [2, 40], while subsequent studies introduced one-shot capabilities [1, 8, 9, 27, 35, 37, 39, 41]. Many of these works employ predefined motion representations, such as 3DMM's blendshapes [8, 17, 37, 39]. In contrast, some learned latent emotion representation from scratch [2, 9, 27, 41, 42]. This latter approach holds potential for better motion representation as in such settings, expression descriptors become entirely trainable and free from inheriting the limitations typical for blendshapes as was shown in [2, 31]. In our work, we adopt this strategy, remarkably improving upon the methodology used in MegaPortraits[9].

Method	FEED dataset internal structure						
	BFEs	Head rotations	Winks	Eyes move	Asym. em.	Tongue em.	Extreme em.
Mode	Mild/Strong	Axes/random	Mild/Strong	Random	Strong	Strong	Extreme
# of part.	21/23	21	23	21	23	23	23
Avg. length	1:38/1:52	0:58	0:49	0:33	1:41	1:53	3:02

Table 1. The tasks we asked our participants to perform for our FEED dataset. Here BFEs stands for basic facial expressions, shown in Fig. 2

Dataset	#Actors	#Views	Datasets comparison				Resolution
			Basic exp.	Strong exp.	Extreme & Assym. exp.	Tongue exp.	
SAVEES	4	1	✓	✗	✗	✗	1280 × 1024
RAVGESS	24	1	✓	✗	✗	✗	1920 × 1080
CREMA-D	91	1	✓	✗	✗	✗	1280 × 720
MMI	25	1	✓	✗	✗	✗	1920 × 1080
MEAD	48	7	✓	✓	✗	✗	1920 × 1080
Ours (FEED)	23	3	✓	✓	✓	✓	3840 × 2160

Table 2. Comparisons with modern, high-quality audio-visual datasets created in controlled settings.

## 2.2. Face expression datasets

Early image datasets [22?] predominantly offered annotated face expression data corresponded to up to 8 basic emotions. However, these datasets are not well-suited for training head avatars, which necessitate video or multi-view data.

The SAVEE [16] dataset captures facial expressions in speech, but involves just 4 actors. The MMI [23] dataset, more expansive in actor participation, offers spontaneous moderate-intensity expression but is restricted to single-view sequences. RAVDESS [20] distinguishes itself by capturing two motion intensities, but its limited recordings challenge broad applicability. CREMA-D [3] encompasses three expression intensities, and though MEAD [33] provides detailed multi-view data across three intensities, it’s limited to the standard eight expression groups. Our novel FEED dataset aims to address the scientific community’s demand for high-quality multi-view facial expression videos outside the standard categories: Joy, Fear, Sadness, Disgust, Anger Contempt and Surprise.

## 2.3. Speech-driven head avatars

While numerous studies introduce speech-driven avatars [25, 30, 37, 43, 46, 48?], few excel in producing high-quality talking heads with authentic rotations, blinks, and the capability to use both visual and audio inputs. For instance, Wav2Lip [25] aims at re-dubbing videos with accurate lip motions, but often falls short in realism with using single images. Diffused heads [30] struggles to generate long sequences and doesn’t provide access to pose control. MakeItTalk [48] animates facial landmarks in a speaker-specific manner, yet struggles with head pose controls due to its non-reliance on 3D. PC-AVS [46] offers a solution, but demands a driving video for pose modulation. Newly developed models such as SadTalker [43] and StyleHEAT [37] have shown promising outcomes in generating high-quality talking heads. We conducted a quantitative evaluation of our audio-driven mode compared to the aforementioned models.

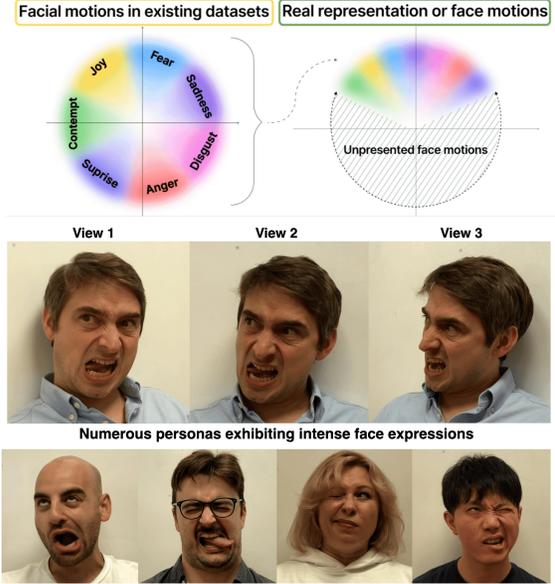


Figure 2. Illustration of the problem in publicly available face expression data and selected examples from our FEED dataset.

## 3. FEED dataset

As previously mentioned, the current publicly available datasets with human face expression videos fall short in capturing a broad range of facial manipulations (see illustration at Fig. 2). In our view, this gap could stem from the complexity of designing clear tasks for participants, coupled with the reality that not everyone is ready for performing extreme expressions on camera. To bridge this gap, our pioneering FEED dataset is designed to meet the scientific community’s need for high-quality, multi-view emotion videos that extend beyond standard emotions. It contains various expressions, including strong asymmetric ones, tongue and cheeks movements, winks, head rotations, eye movements and more nuanced gestures. ;

Our dataset consists of 520 multi-view videos of 23 subjects, captured with 3 cameras. As extreme face motions are complex and their perception can be heavily influenced by subtle differences, we use a high resolution of 4k for all video, capturing the whole face up to the level of individual hair strands and wrinkles, as shown in Fig. 2. For more examples, please refer to the supplementary materials. Our participants were asked to perform 7 tasks (see Tab. 1) to cover as many facial and head movements as possible.

In our comparison with other expression datasets detailed in Tab. 2, the MEAD dataset [33] emerges as the closest to ours regarding viewpoints and maximal expression intensity. However, despite MEAD’s larger number of actors and viewpoints, it offers significantly less variety of facial expressions and have lower resolution than our dataset.

## 4. Expression enhancement

We found that the MegaPortraits [9] model fails to transfer intense expressions correctly, as shown in Fig. 9. Using our FEED dataset of strong and uneven facial expression to fine-tune a pre-trained MegaPortraits model didn’t enhance final results as shown at Tab. 4. Training from scratch on FEED leads to fast overfitting due to the small number of identities presented in the dataset compared to VoxCeleb2 [6] used in [9] (23 vs ~5000). Our method effectively injects this small dataset into training, yielding desirable results. We list our key findings here, for all details, see supplementary materials. Our model’s scheme, displayed in Fig. 4, shares similarities with [9], leading us to occasionally refer to certain elements of our scheme in the context of [9]’s pipeline.

### 4.1. Latent expression space

We begin by exploring MegaPortraits’ latent expression descriptors’ space ( $\mathbf{z}_{s/d}$  in Fig. 4), which is crucial for expression transfer, through PCA analysis. Inspired by [18], which demonstrated that the area under the cumulative explained variance curve of singular values (denoted here as  $\text{AUC}_{\mathbf{z}}$ ) can serve as an effective metric for dimensionality collapse in some latent spaces, we applied this measure to our study.

This metric allows us to forecast model performance: greater collapse suggests reduced entropy in the latent space, which correlates with lower representational quality. High quality of the expression space, produced by  $\mathbf{E}_{motion}$  (as shown in Fig. 4), is crucial to capture and differentiate between the nuances of strong facial expressions. Our ablation study confirms (see Tab. 4) that there is a notable correlation between how broad and isotropic model’s latent space is (expressed by  $\text{AUC}_{\mathbf{z}}$ ) and its final performance. Visual comparison is available in supplementary materials. As shown in Fig. 3, our model’s latent space outperforms MegaPortraits in expression representation ability. This is supported by both visual (Fig. 9) and quantitative analysis (Tab. 3).

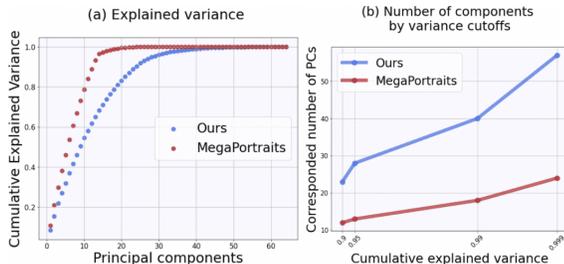


Figure 3. Comparison of latent spaces. Left plot shows that our model’s latent space is wider and exhibits more even variance distribution. Also, as shown on the right plot, a greater number of principal components are involved in capturing variance across various thresholds. This implies a more robust representational capacity of expression space compared to [9]. The VoxCeleb2 test set was used for both plots.

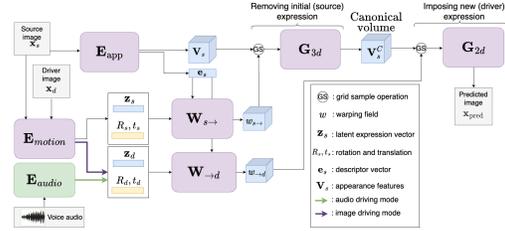


Figure 4. Method Overview. We use  $\mathbf{E}_{app}$  to extract volume features  $\mathbf{V}_s$  and a global descriptor  $\mathbf{e}_s$  from the source image. Then  $\mathbf{E}_{motion}$  or  $\mathbf{E}_{audio}$  generates motion representations from source and driver, including head rotations  $\mathbf{R}_{s/d}$ , translations  $\mathbf{t}_{s/d}$ , and expression descriptors  $\mathbf{z}_{s/d}$ . Using them, we predict warpings  $\mathbf{w}_{s \rightarrow}$  and  $\mathbf{w}_{d \rightarrow}$ . First warping and  $\mathbf{G}_{3D}$  transform  $\mathbf{V}_s$  into a canonical volume  $\mathbf{V}_s^C$  by removing the source motions. Second warping and  $\mathbf{G}_{2D}$  imposes the driver’s motions and renders the final image.

Both metrics shown in Fig. 3 were vital in our model’s development since they highly correlate with the final model’s ability to transfer intense expressions. This was particularly useful for finding adjustments described in Sec. 4.2 and Sec. 4.3, as early visual results during training are not very telling, but changes in the latent space are noticeable after a few epochs. We also found that in MegaPortraits, just a few principal components significantly affect variance, with only 18 components making up 99% (Fig. 3), suggesting their 512  $\mathbf{z}$ ’s dimension is oversized. For our model, we found a 128-dimensional  $\mathbf{z}$  is optimal, offering solid performance and reducing overfitting in imbalanced data, as our ablation study reveals. We calculate the Explained Variance ( $\text{EV}_i$ ) for the  $i^{th}$  component and ( $\text{AUC}_{\mathbf{z}}$ ) using equations shown in Eq. (1), which involves standardizing the vector set  $\mathcal{Z}$  for PCA and sorting the eigenvalues  $\lambda$  in descending order.

$$\mathbf{Z}_{std} = \frac{\mathbf{Z} - \bar{\mathbf{Z}}}{\sigma(\mathcal{Z})}, \quad \lambda = \text{sorted}_{\text{desc}} \left( \text{eig} \left( \frac{1}{n-1} \mathbf{Z}_{std}^T \mathbf{Z}_{std} \right) \right), \quad (1)$$

$$\text{EV}_i = \frac{\lambda_i}{\sum_{j=1}^p \lambda_j}, \quad \text{AUC}_{\mathbf{z}} = \frac{1}{d} \frac{\sum_{i=1}^d \sum_{k=1}^i \lambda_k}{\sum_{j=1}^d \lambda_j},$$

### 4.2. Canonical volume

Despite being designed to exclude face expression details, we found that canonical volume ( $\mathbf{V}^C$  in Fig. 4) in MegaPortraits actually retains significant expression information from the source image, contributing to poor translation of intense expressions, as shown in our experiment detailed in Fig. 5. In this experiment, using portraits with varied expressions, we assessed the expression translation accuracy and how expression intensity affects the  $\mathbf{V}^C$ s. Our findings suggest that the canonical volumes in MegaPortraits are not truly neutral. As confirmed by our ablation study (Tab. 4), creating an expression-free  $\mathbf{V}^C$  is essential for precise translation of intense motions. We believe that a more neutral canonical volume improves tractability and effectiveness in expression translation tasks.

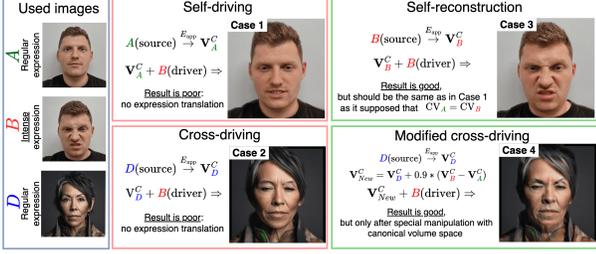


Figure 5. Canonical volumes ( $\mathbf{V}^C$ ) in MegaPortraits [9] are not expression-neutral. To show it we use three portraits: A (regular expression), B (intense expression), D (regular expression, new identity). *Case 1* visualizes poor results in transferring B’s intense expression to A using a self-driving mode, contrasting with the effective reconstruction of B when used as both driver and source in *case 3*. This discrepancy, indicative of expression leakage into the canonical volume, is further quantified by a 43% relative difference in  $\mathbf{V}^C$ s between B and C, contrary to expectations of their similarity for same identity. In cross-driving generation *case 2*, using B as the driver and D as the source yielded poor results. However, in *case 4*, after using additive operations on source canonical volumes, the output expression is much closer to driver than in *case 2*. This manipulation again confirms the significant retention of expression information in canonical volumes.

To overcome this issue, we propose to match canonical volumes from the different images of the same person ( $\mathbf{V}_{s^n}^C, \mathbf{V}_{d^n}^C$ ) during training as following:

$$\mathcal{L}_{CV}^n = \mathcal{L}_{MAE}(\mathbf{V}_{s^n}^C, \mathbf{V}_{d^n}^C), \quad (2)$$

This loss ensures  $\mathbf{V}^C$  remains stable and expression-independent, crucial for translating intense expressions, as shown in our ablation study (Tab. 4).

### 4.3. Source-driver mismatch loss

In addition to maintaining expression-free  $\mathbf{V}^C$ , it’s vital to remove all identity information from the expression vector  $\mathbf{z}$ . While contrastive losses address this in [9], our experiments indicate they are insufficient to prevent overfitting in our imbalanced data scenario, where emotionally intense images represent only 0.1% of our training set but are sampled 25% of the time. We introduce a novel self-supervised loss which mitigates identity information in latent expression vectors:

$$\mathcal{L}_{sdm}(\mathbf{z}_s, \mathbf{z}_d) = w * \max(0, \cos(\mathbf{z}_s, \mathbf{z}_d) - margin),$$

$$w; margin = \begin{cases} 1; 0.5 & \text{if } \mathbf{z}_s, \mathbf{z}_d \text{ are from VC2} \\ 10; 0.25 & \text{if } \mathbf{z}_s, \mathbf{z}_d \text{ are from FEED} \end{cases} \quad (3)$$

Here,  $\mathbf{z}_s, \mathbf{z}_d$  are emotion vectors from source and driver images,  $w, margin$  adjust the loss’s intensity and strictness respectively. We increase  $w$  for FEED due to its higher sampling rate and decrease  $margin$  for more assured expression variance in  $s, d$ . Refer to Fig. 6 for a visual summary of all self-supervised losses enhancing our latent emotion space.

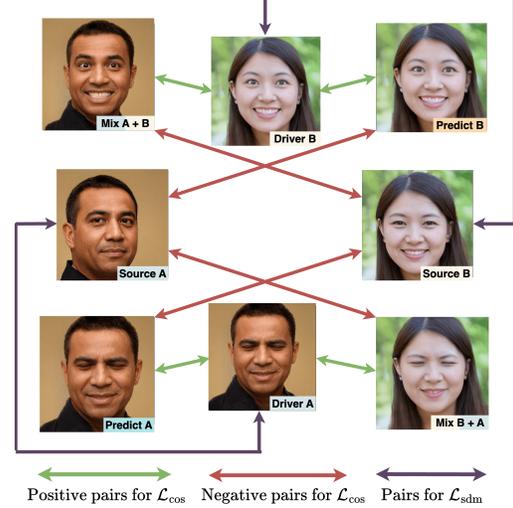


Figure 6. The visualisation of our self-supervised losses.  $\mathcal{L}_{cos}$  is a contrastive loss, similar to one used in [9] but uses different pairs as described in Sec. 5.1. Our novel  $\mathcal{L}_{sdm}$  is designed to prevent overfitting in an imbalanced dataset.

## 5. Incorporating speech-driving mode

### 5.1. Latent space disentanglement

During our latent space analysis described in Sec. 4, we found that combining expression from one driver with the head pose from another resulted in poor performance for MegaPortraits model, as illustrated in Fig. 7 (bottom row). Indeed, the base model lacks a mechanism to intentionally prevent head pose leakage. Disentangled expression latent space, besides expanding the model’s use cases, plays the crucial role for the speech-driving mode. When the latent vector  $\mathbf{z}$  is entangled with head pose data, predicting it from speech becomes challenging because speech lacks head rotation information, unlike images. Moreover, we would like to have a full control over head rotations during speech mode. We were able to make our latent space disentangled by changing the way we sample images for  $\mathcal{L}_{cos}$  Eq. (5).

The MegaPortraits model uses  $\mathcal{L}_{cos}$  - a modified *Large Margin Cosine Loss* [32] to mitigate the transfer of appearance features to expression descriptor (latent vector  $\mathbf{z}$  on Fig. 4), which is crucial for cross-driving mode. Without this, differences in appearance between the source and driver, like hairstyle or skin tone, leak from driver to output image.

For computing this loss, the authors utilize a supplementary source-driving pair ( $\mathbf{x}_{s^*}$  and  $\mathbf{x}_{d^*}$ ) from a different video with another identity, ensuring a distinct appearance from the current  $\mathbf{x}_s, \mathbf{x}_d$  pair. The base model is then employed to produce *cross-reenacted* image ( $\hat{\mathbf{x}}_{s^* \rightarrow d} = \mathbf{G}_{base}(\mathbf{x}_{s^*}, \mathbf{x}_d)$ ). Concurrently, they determine a separate motion descriptor,  $\mathbf{z}_{d^*} = \mathbf{E}_{motion}(\mathbf{x}_{d^*})$ . Descriptors  $\mathbf{z}_{* \rightarrow d}$  from the respective forward passes are also used for positive and negative pairs.



Figure 7. The illustration of disentanglement problem. Images in the right row indicate the expression descriptor  $\mathbf{z}$  entangled with head pose for MegaPortraits, whereas we avoid it using changes proposed in Sec. 5.1.

Our innovation here is a new sampling strategy: for each  $\mathbf{x}_s, \mathbf{x}_d$  pair, we sample one more random additional pair  $\mathbf{x}_{s^m}$  and  $\mathbf{x}_{d^m}$  apart from  $\mathbf{x}_{s^*}$  and  $\mathbf{x}_{d^*}$ . We then apply the model to produce the following *cross-reenacted moved* image using source identity from  $\mathbf{x}_{s^*}$ , desired head pose from  $\mathbf{x}_{s^m}$  and emotions from  $\mathbf{x}_d$ , represented as:  $\hat{\mathbf{x}}_{s^*m \rightarrow d} = \mathbf{G}_{\text{base}}(\mathbf{x}_{s^*m}, \mathbf{x}_d)$ . This creates a positive pair with the same desired emotion but varying head poses.

Motion descriptors are organized into *positive pairs*  $\mathcal{P}$  for alignment, and *negative pairs*  $\mathcal{N}$  for non-alignment:  $\mathcal{P} = \{(\mathbf{z}_{s \rightarrow d}, \mathbf{z}_d), (\mathbf{z}_{s^* \rightarrow d}, \mathbf{z}_d), (\mathbf{z}_{s^*m \rightarrow d}, \mathbf{z}_d)\}$ , and  $\mathcal{N} = \{(\mathbf{z}_{s \rightarrow d}, \mathbf{z}_{d^*}), (\mathbf{z}_{s^* \rightarrow d}, \mathbf{z}_{d^*}), (\mathbf{z}_{s^*m \rightarrow d}, \mathbf{z}_{d^*})\}$  (see Fig. 6). These pairs are used to calculate the following cosine distance:

$$d(\mathbf{z}_i, \mathbf{z}_j) = s \cdot (\langle \mathbf{z}_i, \mathbf{z}_j \rangle - m), \quad (4)$$

where both  $s$  and  $m$  are hyperparameters. This distance is then used to calculate a large margin cosine loss (CosFace) [? ]:

$$\mathcal{L}_{\text{cos}} = - \sum_{(\mathbf{z}_k, \mathbf{z}_l) \in \mathcal{P}} \log \frac{\exp \{d(\mathbf{z}_k, \mathbf{z}_l)\}}{\exp \{d(\mathbf{z}_k, \mathbf{z}_l)\} + \sum_{(\mathbf{z}_i, \mathbf{z}_j) \in \mathcal{N}} \exp \{d(\mathbf{z}_i, \mathbf{z}_j)\}} \quad (5)$$

This loss prevents the head pose from leaking into the embeddings as shown in Fig. 7 (upper row)

By doing this, we achieved the capability to interpret  $\mathbf{z}_i$  as a latent vector encapsulating the emotional content from an image  $\mathbf{x}_i$ . This advancement opened the door for an intriguing possibility during inference: the prediction of this vector could potentially be derived from an audio signal instead of a facial image using the motion encoder  $\mathbf{E}_{\text{motion}}$ , effectively transforming the model into a speech-driven system. The critical challenge here is aligning the audio input with the pretrained latent space. We addressed this by introducing an additional audio encoder,  $\mathbf{E}_{\text{aud}}$ , and employing  $\mathbf{E}_{\text{motion}}$  as a teacher model for this new encoder. A crucial factor in achieving remarkable accuracy in lip synchronization was our ability to isolate specific components within the expression latent space, that are solely responsible for mouth movements. We describe it in the next subsection Sec. 5.2.

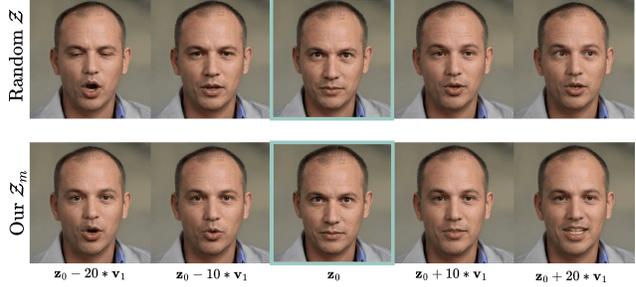


Figure 8. In the upper row, the first principal component from a random subset of expression vectors  $\mathcal{Z}$  affects not just the mouth but also blinks and gaze, making it difficult to interpret and isolate mouth movement components, needed for training  $\mathbf{E}_{\text{aud}}$ . However, as shown in the bottom row, using  $\mathcal{Z}_m$  from Sec. 5.2, the first principal component,  $\mathbf{v}_1$ , focused on mouth movements, is effectively isolated,  $\mathbf{z}_0$  - latent vector, corresponded to the central image.

## 5.2. Mouth movements

To train the audio encoder  $\mathbf{E}_{\text{aud}}$  for predicting  $\mathbf{z}$  from speech, we chose to use the motion encoder  $\mathbf{E}_{\text{motion}}$  as a teacher model. However, direct matching vectors  $\mathbf{z}_i^{\text{aud}}$  from  $\mathbf{E}_{\text{aud}}$  and pseudo-ground truth  $\mathbf{z}_i$  from  $\mathbf{E}_{\text{motion}}$  gives poor results (see Tab. 6). We believe that the reason behind this is that speech-derived  $\mathbf{z}_i^{\text{aud}}$ , while able to capture lip movements, struggles with other facial motions, especially in the upper face. Thus, we focused on isolating mouth movement components in the expression latent space using PCA analysis.

Employing PCA on a broad array of expression vectors  $\mathcal{Z}$ , generated from images with varied facial emotions, yields principal components with their explained variances ( $EV_i$ ). Altering these components, particularly those with high  $EV_i$ , significantly changes facial expressions in an image. However, they typically affect a combination of facial features- Fig. 8 (upper row), not isolated areas like the mouth or eyes.

To focus on mouth movements, we created a unique set of  $\mathcal{Z}_m$  from just one video of a person performing mouth manipulations and manually edited it to minimize upper face movements. This involved using a still upper face from the first frame in all subsequent frames, ensuring that the principal components mainly represented mouth movements. Applying PCA of  $\mathcal{Z}_m$  reveals that most principal components are responsible for solely mouth movements. For illustration of the first component, see Fig. 8 (bottom row). Based on distilled components, we introduce mouth PCA mouth loss:

$$\mathcal{L}_{\text{PCA}}(v_i, v_j, n) = \frac{1}{n} \sum_{k=1}^n |\text{PC}_{\mathbf{v}_i}(k) - \text{PC}_{\mathbf{v}_j}(k)| \quad (6)$$

Where  $n$  is the number of principal components to be considered,  $\text{PC}_{\mathbf{z}_i}(k)$  and  $\text{PC}_{\mathbf{z}_j}(k)$  are the  $k$ -th principal components of vectors  $\mathbf{z}_i$  and  $\mathbf{z}_j$  respectively. Implementing this loss has demonstrably enhanced the quality and accuracy of our results, as our ablation study confirms (see Tab. 6).

## 6. Experiments

In this section, we describe experiments we set for comparison for both our image-driven and speech-driven mode. Please see implementation details such as all used losses, component’s sizes and architecture details, training procedure, data preprocess in supplementary materials.

### 6.1. Image-driven comparison

**Methods.** We evaluated our image-driven model alongside various one-shot models. For this comparison, we included models like NOFA [39] and StyleHEAT [37], which use 3DMM expression blendshapes without disentanglement training, and models with trainable latent face motion representations like FOMM [27], UVA [28], MetaPortrait [42] and MegaPortraits [9]. It’s noteworthy that NOFA and UVA require per-source optimization, unlike the others, including our model.

**Evaluation metrics and data.** For our evaluation, we use 100 random, non-child, single-person images from the FFHQ dataset as sources. From the MEAD dataset, we chose 100 images covering all high-intensity emotions except neutral, and another 100 from the FEED dataset’s extreme emotion task, featuring 5 identities not seen during training of our model. Each source image was paired with 20 driving images from these sets, totaling 2000 pairs.

For assessing reenactment quality, we employ various metrics including the Frechet Inception Distance (FID) [14] to measure the distributional discrepancy between synthetic and real images. Cosine similarity (CSIM) from a face recognition network [4] quantifies the identity preservation in generated images. Additionally, we conducted user studies to determine preferences regarding motion (UMTN) and appearance (UAPP) preservation, presenting participants with triplets—either a driving or source image and two method outputs—and asking them to choose the one with better motion or appearance preservation. Our user study includes 34 participants and around 300 unique questions. Both metrics numerically reflect preference percentages. As our visual results Fig. 9, and quantitative metrics Tab. 3 show, our method demonstrates exceptional ability in translating strong and extreme expressions, notably surpassing other image-driven avatar techniques in user preference for facial expression translation (UMTN) and FID scores.

### 6.2. Speech-driven comparison

**Methods.** To assess our model’s speech-driven capabilities, we compared it with StyleHEAT [37], which also operates using both image and audio inputs, and audio-exclusive models including SadTalker [43], MakeItTalk [48], and PC-AVS [46].

**Evaluation metrics and data.** We employ the evaluation procedure described in [?]. We measure video realism using Frechet Inception Distance (FID) [14]. We also extract

Driver:	MEAD				FEED			
Method	FID↓	CSIM↑	UMTN↑	UAPP↑	FID↓	CSIM↑	UMTN↑	UAPP↑
MetaPortrait [42]	77.5	0.64	*	*	78.8	0.66	*	*
NOFA [39]	69.2	0.65	*	*	69.8	0.66	3.8	8.8
FOMM [27]	84.4	0.56	2.4	0.3	85.9	0.55	1.8	0.5
UVA [19]	78.7	0.68	0.9	6.4	79.5	0.68	0.6	9.4
StyleHEAT [37]	71.9	0.63	3.7	3.3	72.4	0.61	2.6	3.8
MegaPortraits [9]	61.1	0.73	22.7	46.6	62.8	0.73	16.7	42.8
Our (EMOPortraits)	59.6	0.74	70.3	43.4	60.2	0.70	74.6	34.7

Table 3. Our method notably outperforms others in the FID score and strongly leads in the user preference metric for face expression translation (UMTN). It excels in reliably translating strong and extreme expressions from the driving image, distinguishing it as the most capable among the compared methods. Additionally, our approach maintains the source image’s shape and appearance (measured by CSIM and UAPP metrics) on par with MegaPortraits [9] when using strong and regular expressions from the MEAD dataset as the driving signal. However, we observe a minor deviation in identity preservation metrics for extreme emotions (using the FEED dataset). This difference may arise because our method attempts to transfer asymmetric facial expressions, while MegaPortraits might modify the source image less, or not at all, with a challenging driver (as depicted in Fig. 5, case 1), thereby achieving better identity preservation. UVA was excluded from the MEAD-driven user study due to the authors were able to provide only part of the test data. Also, MetaPortrait was excluded from both user studies because of late data provision by authors, offering only the basic model results, without their super-resolution module outputs.

Method	FID↓	CSIM↑	UMTN↓	UAPP↓	AUC <sub>z</sub> ↓
MegaPortraits fine-tuned	62.3	0.63	8.5	12.4	0.89
Ours w/o $\mathcal{L}_{CV}^n$	61.4	0.67	9.2	17.9	0.88
Ours w/o $\mathcal{L}_{sdm}$	61.8	0.59	11.3	6.7	0.85
Ours $dim(\mathbf{z}) = 512$	63.4	0.44	25.6	1.1	0.77
Ours	59.6	0.74	45.4	61.9	0.80

Table 4. Our ablation study examines key aspects of the image-driven mode. After fine-tuning MegaPortraits on the FEED dataset, we observed a decline in identity preservation due to overfitting on FEED identities, without any improvement in emotion translation. The lack of  $\mathcal{L}_{CV}^n$  primarily impacts the ability to translate extreme expression, as indicated by UMTN. Additionally, not using  $\mathcal{L}_{sdm}$  or setting  $dim(\mathbf{z}) = 512$  results in identity leakage, particularly evident in the latter scenario. For this ablation, we use the part of our driven by MEAD dataset and described in Sec. 6.1

facial landmarks from outputs of all methods, frontalized, and normalized to position the mouth edges at (-1, 0) and (1, 0). We compute mean absolute error (MAE) for both predicted mouth landmark positions ( $M_P$ ) and velocities ( $M_V$ ), along with errors in facial geometry ( $F_P$ ) and velocity ( $F_V$ ). For our evaluation, we use 100 randomly selected video sequences from the HDTF dataset [44], each up to 30 seconds. As the results show in Tab. 5, our speech-driven model performs on par with leading methods, excelling in realism and facial dynamics. SyncNet [5] is used to assess lip sync quality, providing offset (temporal misalignment between the audio and video) and confidence (certainty of audio-visual alignment) scores.

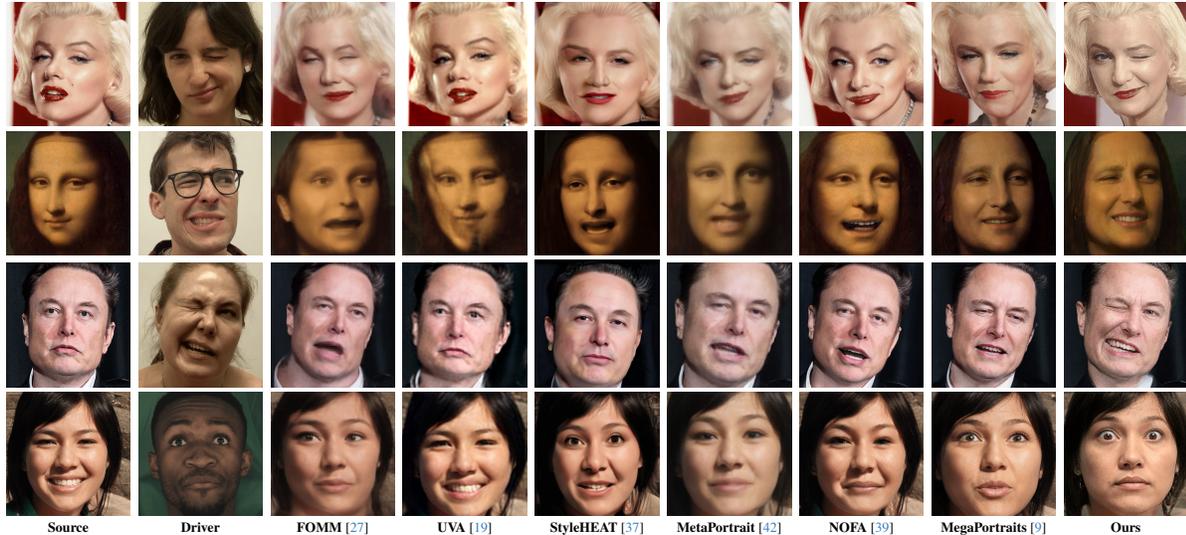


Figure 9. A qualitative comparison of head avatar systems in cross-reenactment scenario.

Method	FID↓	$F_P$ ↓	$F_V$ ↓	$M_P$ ↓	$M_V$ ↓	Sync. offset ↓	Sync. conf. ↑
StyleHEAT [37]	62.4	1.34	0.41	2.98	1.67	0.41	7.22
PC-AVS [46]	71.5	2.78	0.51	2.56	<b>1.04</b>	3.79	<b>8.42</b>
SadTalker [43]	<u>55.6</u>	<u>0.91</u>	<u>0.35</u>	<b>2.12</b>	<u>1.17</u>	<u>0.27</u>	<u>7.58</u>
MakeItTalk [48]	63.3	2.15	0.38	3.11	1.59	2.98	5.17
Ours (EMOPortraits)	<b>28.5</b>	<b>0.82</b>	<b>0.33</b>	<u>2.53</u>	1.38	<b>0.07</b>	5.77

Table 5. Our model matches other leading speech-driven models, excelling in realism (FID), facial geometry and velocity, but is slightly less accurate in mouth landmark positions and velocities. We also score best in audio-video sync offset, but lower in confidence. We believe that while SyncNet metrics are still helpful, they may not fully represent true synchronization quality, as generated videos can sometimes score higher than ground truth without necessarily being better synchronized. In our experimentation, ground truth videos scored an offset of 1.13 (which is worse than our method, SadTalker, and StyleHEAT) and a confidence score of 8.35 (less than PC-AVS). For qualitative analysis, please refer to our supplementary materials.

Method	FID↓	$M_P$ ↓	$M_V$ ↓	$F_P$ ↓	$F_V$ ↓
Ours w/o disentanglement	44.7	0.98	0.45	5.78	2.55
Ours w/o $\mathcal{L}_{PCA}$	29.3	0.84	0.35	3.13	1.88
Ours with $MAE(\mathbf{z}_i, \hat{\mathbf{z}}_i^{\text{speech}})$	<b>28.2</b>	<u>0.84</u>	<u>0.34</u>	<u>2.61</u>	<u>1.45</u>
Ours	<u>28.5</u>	<b>0.82</b>	<b>0.33</b>	<b>2.53</b>	<b>1.38</b>

Table 6. Ablation study of our speech-driven mode, where in the first line, we show that it notably underperforms without head pose - face expression disentanglement of the latent space (see Sec. 5.1). Removing  $\mathcal{L}_{PCA}$  leads to a significant reduction in generating realistic mouth movements. Replacing  $\mathcal{L}_{PCA}$  with  $MAE(\mathbf{z}_i, \hat{\mathbf{z}}_i^{\text{speech}})$  causes a noticeable, decline in mouth movement generation. This suggests that identifying mouth-related PCA components was beneficial in enhancing the outcome.

### 6.3. Ablation study

We conducted an extensive ablation study to evaluate the contributions of individual components within our method. For our main model, we present the evaluation of the importance of the proposed source-driver mismatch loss, canonical volume loss, and change of  $dim(\mathbf{z})$  from 512 in original model to 128 Tab. 4. For our speech-driven mode, we show how vital the disentanglement of the latent space for speech-drive ability, how  $\mathcal{L}_{PCA}$  is better than naive matching of PCA components from random  $\mathcal{Z}$  and how it helps match speech and latent expression space Tab. 6. For more ablation analysis, plots and visual comparison, please refer to our supplementary.

## 7. Conclusion

In this paper, we introduce EmoPortraits, a novel method for creating neural avatars with superior performance in image-driven, cross-identity emotion translation. Our speech-driven mode makes it possible to drive the facial animation through multiple conditions (video, audio, head motion). We collected FEED dataset which, we believe, will be a valuable asset for researchers in diverse human-centered studies.

However, our method has some limitations. It doesn't generate the avatar's body or shoulders, limiting some use cases. We currently integrate our output with a source image body. Additionally, the model sometimes struggles with accurate expression translation and underperforms with extensive head rotation. These challenges are crucial for future enhancements and remain central to our ongoing research efforts.

## References

- [1] Stella Bounareli, Christos Tzelepis, Vasileios Argyriou, Ioannis Patras, and Georgios Tzimiropoulos. Hyperreenact: One-shot reenactment via jointly learning to refine and retarget faces, 2023. [2](#)
- [2] Egor Burkov, Igor Pasechnik, Artur Grigorev, and Victor Lempitsky. Neural head reenactment with latent pose descriptors. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2020. [2](#)
- [3] Houwei Cao, David G Cooper, Michael K Keutmann, Ruben C Gur, Ani Nenkova, and Ragini Verma. Crema-d: Crowd-sourced emotional multimodal actors dataset. *IEEE transactions on affective computing*, 5(4):377–390, 2014. [3](#)
- [4] Qiong Cao, Li Shen, Weidi Xie, Omkar M. Parkhi, and Andrew Zisserman. Vggface2: A dataset for recognising faces across pose and age, 2018. [7](#)
- [5] Joon Son Chung and Andrew Zisserman. Out of time: automated lip sync in the wild. In *Computer Vision—ACCV 2016 Workshops: ACCV 2016 International Workshops, Taipei, Taiwan, November 20–24, 2016, Revised Selected Papers, Part II 13*, pages 251–263. Springer, 2017. [7](#)
- [6] Joon Son Chung, Arsha Nagrani, and Andrew Zisserman. Voxceleb2: Deep speaker recognition. In *INTERSPEECH*, 2018. [4](#), [11](#), [13](#)
- [7] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, K. Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009. [11](#)
- [8] Michail Christos Doukas, Stefanos Zafeiriou, and Viktoriia Sharmanska. Headgan: One-shot neural head synthesis and editing, 2021. [2](#)
- [9] Nikita Drobyshev, Jenya Chelishev, Taras Khakhulin, Aleksei Ivakhnenko, Victor Lempitsky, and Egor Zakharov. Megaportraits: One-shot megapixel neural head avatars, 2023. [2](#), [4](#), [5](#), [7](#), [8](#), [11](#), [12](#), [15](#)
- [10] Tobias Fischer, Hyung Jin Chang, and Y. Demiris. Rt-gene: Real-time eye gaze estimation in natural environments. In *ECCV*, 2018. [11](#)
- [11] Guy Gafni, Justus Thies, Michael Zollhöfer, and Matthias Nießner. Dynamic neural radiance fields for monocular 4d facial avatar reconstruction, 2020. [2](#)
- [12] Philip-William Grassal, Malte Prinzler, Titus Leistner, Carsten Rother, Matthias Nießner, and Justus Thies. Neural head avatars from monocular rgb videos, 2022. [2](#)
- [13] Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron Courville. Improved training of wasserstein gans, 2017. [14](#)
- [14] 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. In *Advances in Neural Information Processing Systems*, 2017. [7](#)
- [15] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A. Efros. Image-to-image translation with conditional adversarial networks, 2018. [14](#)
- [16] Philip Jackson and SJUoSG Haq. Surrey audio-visual expressed emotion (savee) database. *University of Surrey: Guildford, UK*, 2014. [3](#)
- [17] Taras Khakhulin, Vanessa Sklyarova, Victor Lempitsky, and Egor Zakharov. Realistic one-shot mesh-based head avatars, 2022. [2](#)
- [18] Alexander C. Li, Alexei A. Efros, and Deepak Pathak. Understanding collapse in non-contrastive siamese representation learning, 2022. [4](#)
- [19] Xueting Li, Shalini De Mello, Sifei Liu, Koki Nagano, Umar Iqbal, and Jan Kautz. Generalizable one-shot neural head avatar, 2023. [7](#), [8](#), [15](#)
- [20] Steven R Livingstone and Frank A Russo. The ryerson audio-visual database of emotional speech and song (ravdess): A dynamic, multimodal set of facial and vocal expressions in north american english. *PLoS one*, 13(5):e0196391, 2018. [3](#)
- [21] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *ICLR*, 2019. [12](#), [13](#)
- [22] Patrick Lucey, Jeffrey F Cohn, Takeo Kanade, Jason Saragih, Zara Ambadar, and Iain Matthews. The extended cohn-kanade dataset (ck+): A complete dataset for action unit and emotion-specified expression. In *2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition-Workshops*, pages 94–101. IEEE, 2010. [3](#)
- [23] Maja Pantic, Michel Valstar, Ron Rademaker, and Ludo Maat. Web-based database for facial expression analysis. In *2005 IEEE International Conference on Multimedia and Expo*, pages 5–pp. IEEE, 2005. [3](#)
- [24] Omkar M. Parkhi, Andrea Vedaldi, and Andrew Zisserman. Deep face recognition. In *BMVC*, 2015. [11](#), [13](#)
- [25] KR Prajwal, Rudrabha Mukhopadhyay, Vinay P Namboodiri, and CV Jawahar. A lip sync expert is all you need for speech to lip generation in the wild. In *Proceedings of the 28th ACM International Conference on Multimedia*, pages 484–492, 2020. [3](#)
- [26] Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. Robust speech recognition via large-scale weak supervision, 2022. [13](#)
- [27] Aliaksandr Siarohin, Stéphane Lathuilière, S. Tulyakov, Elisa Ricci, and N. Sebe. First order motion model for image animation. *ArXiv*, abs/2003.00196, 2019. [2](#), [7](#), [8](#), [15](#)
- [28] Aliaksandr Siarohin, Willi Menapace, Ivan Skorokhodov, Kyle Olszewski, Hsin-Ying Lee, Jian Ren, Menglei Chai, and Sergey Tulyakov. Unsupervised volumetric animation. *arXiv preprint arXiv:2301.11326*, 2023. [7](#)
- [29] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *CoRR*, abs/1409.1556, 2015. [11](#)
- [30] Michał Stypułkowski, Konstantinos Vougioukas, Sen He, Maciej Zieba, Stavros Petridis, and Maja Pantic. Diffused heads: Diffusion models beat gans on talking-face generation, 2023. [3](#)
- [31] Mani Kumar Tellamekala, Ömer Sümer, Björn W. Schuller, Elisabeth André, Timo Giesbrecht, and Michel Valstar. Are 3d face shapes expressive enough for recognising continuous emotions and action unit intensities?, 2023. [2](#)
- [32] Hao Wang, Yitong Wang, Zheng Zhou, Xing Ji, Dihong Gong, Jingchao Zhou, Zhifeng Li, and Wei Liu. Cosface: Large margin cosine loss for deep face recognition, 2018. [5](#)

- [33] Kaisiyuan Wang, Qianyi Wu, Linsen Song, Zhuoqian Yang, Wayne Wu, Chen Qian, Ran He, Yu Qiao, and Chen Change Loy. Mead: A large-scale audio-visual dataset for emotional talking-face generation. In *European Conference on Computer Vision*, pages 700–717. Springer, 2020. [3](#)
- [34] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. High-resolution image synthesis and semantic manipulation with conditional gans. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018. [12](#)
- [35] Ting-Chun Wang, Arun Mallya, and Ming-Yu Liu. One-shot free-view neural talking-head synthesis for video conferencing. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2021. [2](#)
- [36] Yuelang Xu, Lizhen Wang, Xiaochen Zhao, Hongwen Zhang, and Yebin Liu. AvatarMAV: Fast 3d head avatar reconstruction using motion-aware neural voxels. In *Special Interest Group on Computer Graphics and Interactive Techniques Conference Conference Proceedings*. ACM, 2023. [2](#)
- [37] Fei Yin, Yong Zhang, Xiaodong Cun, Mingdeng Cao, Yanbo Fan, Xuan Wang, Qingyan Bai, Baoyuan Wu, Jue Wang, and Yujiu Yang. Styleheat: One-shot high-resolution editable talking face generation via pre-trained stylegan, 2022. [2](#), [3](#), [7](#), [8](#), [15](#)
- [38] Changqian Yu, Jingbo Wang, Chao Peng, Changxin Gao, Gang Yu, and Nong Sang. Bisenet: Bilateral segmentation network for real-time semantic segmentation, 2018. [11](#), [13](#)
- [39] Wangbo Yu, Yanbo Fan, Yong Zhang, Xuan Wang, Fei Yin, Yunpeng Bai, Yan-Pei Cao, Ying Shan, Yang Wu, Zhongqian Sun, and Baoyuan Wu. Nofa: Nerf-based one-shot facial avatar reconstruction, 2023. [2](#), [7](#), [8](#), [12](#), [15](#)
- [40] Egor Zakharov, Aliaksandra Shysheya, Egor Burkov, and Victor Lempitsky. Few-shot adversarial learning of realistic neural talking head models, 2019. [2](#)
- [41] Egor Zakharov, Aleksei Ivakhnenko, Aliaksandra Shysheya, and Victor Lempitsky. Fast bi-layer neural synthesis of one-shot realistic head avatars, 2020. [2](#)
- [42] Bowen Zhang, Chenyang Qi, Pan Zhang, Bo Zhang, Hsiang-Tao Wu, Dong Chen, Qifeng Chen, Yong Wang, and Fang Wen. Metaportrait: Identity-preserving talking head generation with fast personalized adaptation, 2023. [2](#), [7](#), [8](#), [15](#)
- [43] Wenxuan Zhang, Xiaodong Cun, Xuan Wang, Yong Zhang, Xi Shen, Yu Guo, Ying Shan, and Fei Wang. Sadtalker: Learning realistic 3d motion coefficients for stylized audio-driven single image talking face animation, 2023. [3](#), [7](#), [8](#)
- [44] Zhimeng Zhang, Lincheng Li, Yu Ding, and Changjie Fan. Flow-guided one-shot talking face generation with a high-resolution audio-visual dataset. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3661–3670, 2021. [7](#)
- [45] Yufeng Zheng, Victoria Fernández Abrevaya, Marcel C. Bühler, Xu Chen, Michael J. Black, and Otmar Hilliges. I m avatar: Implicit morphable head avatars from videos, 2022. [2](#)
- [46] Hang Zhou, Yasheng Sun, Wayne Wu, Chen Change Loy, Xiaogang Wang, and Ziwei Liu. Pose-controllable talking face generation by implicitly modularized audio-visual representation, 2021. [2](#), [3](#), [7](#), [8](#)
- [47] Yi Zhou, Connelly Barnes, Jingwan Lu, Jimei Yang, and Hao Li. On the continuity of rotation representations in neural networks. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5738–5746, 2019. [13](#)
- [48] Yang Zhou, Xintong Han, Eli Shechtman, Jose Echevarria, Evangelos Kalogerakis, and Dingzeyu Li. MakeltTalk. *ACM Transactions on Graphics*, 39(6):1–15, 2020. [3](#), [7](#), [8](#)
- [49] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A. Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. *2017 IEEE International Conference on Computer Vision (ICCV)*, pages 2242–2251, 2017. [12](#)

## 8. Main model implementation details

### 8.1. Full List of Findings

Beyond the core findings highlighted in the main text, this section outlines additional key distinctions that set our work apart from the MegaPortraits model [9].

**Model Architecture.** Our main model’s structure, depicted in Figure 4 of the main text, excluding the speech-driven component, shares a conceptual resemblance with that of MegaPortraits [9]. However, we have implemented several architecture alterations. First is the reduction in the dimensionality of the latent expression descriptors from 512 to 128, a change detailed in Section 4.1. This reduction enhances the efficiency of the model without compromising the quality of expression representation. Additionally, we have made comprehensive modifications to the architecture and size of the model’s main components to optimize model’s performance. These modifications are visually represented in Fig. 10.

**Dropout.** We have incorporated dropout as a last layer of  $E_{motion}$  that predicts the expression descriptors ( $z$ ) in our model. This implementation serves to improve  $E_{motion}$ ’s capability to construct a more nuanced and robust latent representation of facial expressions, and also aids in preventing overfitting. By ensuring that the model does not become overly reliant on any element of the latent vector, we achieve a more generalized and versatile expression representation capability, essential for dealing with a wide range of facial motions.

**Enhanced Loss Functions.** Beyond new loss functions introduced in the main text, such as the canonical volume loss (outlined in Section 4.2) and the source-driver mismatch loss (described in Section 4.3), we have also integrated the  $\mathcal{L}_{head}$  loss, as mentioned in Section 8.3. This specific loss function plays a noticeable role in the precise predicting of facial regions critical for emotional expression, particularly the eyes and mouth. Additionally, it addresses the previously noted challenges in accurately generating ears. The integration of this loss underscores our model’s attention to detail and commitment to achieving a high degree of realism in facial expression synthesis.

### 8.2. Implementation Details

**Data Preparation.** Our data preparation approach for the VoxCeleb2 (VC2) dataset [6] follows the protocol established in the original model [9]. For our novel FEED dataset, we cropped frames around the face region and resized them. The dataset subset used in our experiments, detailed in Table 1, includes “Winks,” “Tongue Emotion,” and “Extreme & Asymmetric Emotion.” Notably, during training, we did not exploit the multi-view nature of the FEED dataset. Specifically, in each iteration, both the source and driver images were selected from the same single-camera video. Both datasets were employed for training and evaluating our

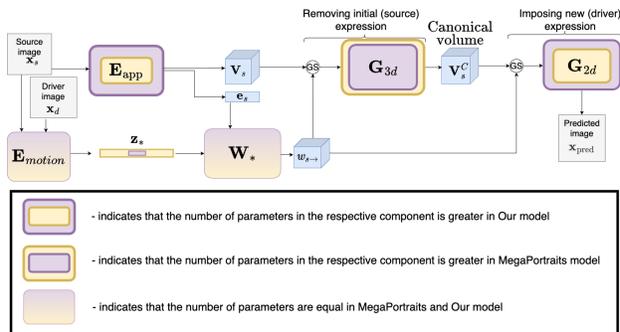


Figure 10. Comparison of our model’s scheme with the MegaPortraits scheme, showing the relative sizes of individual components.

model at a resolution of  $512 \times 512$ .

**Training Details.** Our framework was trained on 8 Nvidia Tesla A100 GPUs for 250,000 iterations, using a batch size of 2 per GPU (16 in total). The training data consisted of a mixture of 75% VC2 examples and 25% FEED examples. Every second iteration involved a batch comprising one image pair from VC2 and one from the FEED dataset. This sampling strategy, integrating image pairs from both datasets, proved more effective than using separate batches from each dataset. By employing contrastive losses where positive and negative pairs spanned both datasets, we mitigated overfitting risks associated with the limited identity variety in the FEED dataset. This approach also facilitated the asymmetric emotion translation to unseen identities.

### 8.3. Used losses

**Photometrical losses.** These are key to aligning the motion and appearance of the predicted image ( $\hat{x}_{s \rightarrow d}$ ) with the ground truth image ( $x_d$ ). To achieve this, we use three distinct pre-trained networks:

- *VGG19* [29] (ILSVRC/ImageNet [7] Trained): This helps in matching the overall content of the images.
- *VGGFace* [24] (Face Recognition Focused): Essential for aligning facial features accurately.
- *Gaze Direction Based on VGG16* [10]: Specifically trained to emulate a top-notch gaze detection system, ensuring precise gaze direction matching.

We measure the similarity by calculating the L1 distance between the feature maps of both the predicted and ground truth images, utilizing all these networks. Additionally, we employ face masks for the eyes, mouth, and ears (sourced from FaceParsing network [38]), focusing our model on these critical head areas. Then we use these masks to match mentioned head regions on the predicted and the ground truth images using L1 loss between pixel corresponded to a specific region. The final photometric loss combines these individual perceptual losses, formulated as:

$$\mathcal{L}_{\text{pho}} = w_{\text{IN}}\mathcal{L}_{\text{IN}} + w_{\text{face}}\mathcal{L}_{\text{face}} + w_{\text{gaze}}\mathcal{L}_{\text{gaze}} + \mathcal{L}_{\text{head}}. \quad (7)$$

Here,  $\mathcal{L}_{\text{head}}$  further breaks down into:

$$\mathcal{L}_{\text{head}} = w_{\text{eyes}}\mathcal{L}_{\text{eyes}} + w_{\text{mouth}}\mathcal{L}_{\text{mouth}} + w_{\text{ears}}\mathcal{L}_{\text{ears}}. \quad (8)$$

**Self-supervised losses.** As detailed in Section 5, we trained our expression descriptors using a modified large margin cosine loss (CosFace) [?], denoted as  $\mathcal{L}_{\text{cos}}$  and presented in Equation 5. This approach is similar to the one employed by the authors of MegaPortraits [9]. Additionally, we introduced two more losses. The source-driver mismatch loss  $\mathcal{L}_{\text{sdm}}$  Equation 3 (described in Section 4.3), which directly influences the expression’s latent space. This loss is pivotal in eliminating identity information from the expression descriptor  $\mathbf{z}_i$  and in preventing overfitting, especially in the context of our extremely imbalanced dataset. The combination of these two losses forms our latent space loss:

$$\mathcal{L}_{\text{lat}} = w_{\text{cos}}\mathcal{L}_{\text{cos}} + w_{\text{sdm}}\mathcal{L}_{\text{sdm}}. \quad (9)$$

The second additional self-supervised loss that enhances the disentanglement of identity and expression is the canonical volume loss  $\mathcal{L}_{\text{CV}}$  (described in Section 4.2 and Equation 2). This loss functions to extract expression information from the canonical volume, thereby reducing the overlap of information contained in  $\mathbf{V}_i^C$  and  $\mathbf{z}_i$ .

**Adversarial losses.** To ensure the predicted images look realistic, adversarial losses are computed using the same predicted and reference images. We follow [9] and train a multi-scale patch discriminator [49] with a hinge adversarial loss. To boost training stability, a standard feature-matching loss is also included ([34]). The GAN loss for the generator is expressed as:

$$\mathcal{L}_{\text{GAN}} = w_{\text{adv}}\mathcal{L}_{\text{adv}} + w_{\text{FM}}\mathcal{L}_{\text{FM}}. \quad (10)$$

To conclude, the total loss which is used to train our model is the sum of individual losses:

$$\mathcal{L}_{\text{main}} = \mathcal{L}_{\text{pho}} + \mathcal{L}_{\text{lat}} + w_{\text{CV}}\mathcal{L}_{\text{CV}} + \mathcal{L}_{\text{GAN}}. \quad (11)$$

We utilized the AdamW optimizer [21] with a cosine learning rate schedule. The initial learning rate was gradually reduced from  $2 \times 10^{-4}$  to  $1 \times 10^{-6}$  over the training iterations. The hyperparameters for the losses were set as follows:  $w_{\text{IN}} = 20$ ,  $w_{\text{face}} = 10$ ,  $w_{\text{gaze}} = 10$ ,  $w_{\text{adv}} = 1$ ,  $w_{\text{FM}} = 40$ ,  $w_{\text{cos}} = 2$ ,  $w_{\text{std}} = 1$  (increased to 10 for pairs from the FEED dataset), and  $w_{\text{CV}} = 1$ . Additionally, we set  $s = 5$  and  $m = 0.2$  in the cosine loss.



Figure 11. More selected examples from our FEED dataset.

## 8.4. Visual Comparison

Our choice of baseline methods, as outlined in our experiment section in the main text, was driven by two key factors. Firstly, these methods are prominent in the field of talking-head video generation using arbitrary identities, making them relevant benchmarks for our study. Secondly, the accessibility of their source code and pretrained model weights, either through public availability or provided by the authors for use with our test set, was a crucial consideration.

The setup for our visual comparison, as described in Section 6.1, was chosen based on specific criteria. We choose FFHQ images as source images due to the consistent clarity of facial features across the dataset, which is essential for accurate comparison. Additionally, it was critical to select identities that were not part of the training datasets for any of the methods used in comparison. This ensures that our comparisons are based on novel identities, providing a fair assessment of each method’s generalization capabilities. This criterion was also applied in selecting driving identities from the MEAD and FEED datasets, which were not used in training by any of the compared methods.

To supplement the comparisons described in the main paper, we provide additional examples for each method (refer to Figure 14). As, for NOFA [39], the range of examples is limited due to the restricted number of inferred identities provided by the authors, we provide a second set of additional examples excluding NOFA [39] (see Figure 15). Furthermore, we include a visual comparison Figure 16 for our ablation study (see Table 4 in the main text).

## 9. Speech-driven mode

### 9.1. Implementation details

In this section, we detail the workings of our audio encoder,  $\mathbf{E}_{\text{aud}}$ , as depicted in Fig. 4, for its application in speech-driven scenarios. The encoder,  $\mathbf{E}_{\text{aud}}$ , is designed for generating latent expression vectors, denoted as  $\mathbf{z}$ , from speech inputs. The scheme of the encoder presented on Fig. 12.

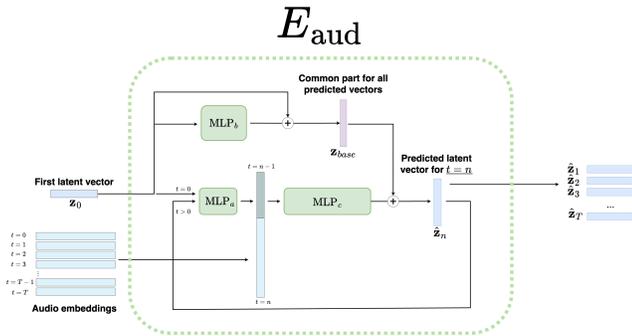


Figure 12. Comparison of our audio encoder used to predict latent expression vectors during speech-driven mode.

First, we use Whisper [26] model to retrieve audio embeddings from a raw audio clip containing speech. This step yields a series of  $T$  audio embedding vectors, where each vector is linked to a specific frame in the video clip. Next, we employ a multilayer perceptron (MLP), designated as  $MLP_b$ , to compute the base component of our latent expression vectors,  $z_{base}$ . This base component encapsulates common facial features such as initial gaze direction and the pose of the upper facial region, as observed in the first frame of the input. Then, another MLP,  $MLP_a$ , is employed. It uses the previously calculated latent vector  $z$ —initially which is  $z_0$  for  $t = 0$  and  $\hat{z}_n$  for  $t > 0$ —to align the latent features of  $z$  with those of the audio embeddings and derive features that useful for final prediction. In the final step, after merging the relevant audio embedding with the output from  $MLP_a$ , the next network,  $MLP_c$  is utilized and responsible for computing a part of the vector which, when added to  $z_{base}$ , forms the final latent expression vector  $\hat{z}_n$  for each frame.

## 9.2. Data Preparation

Similar to our primary model, we employ the VoxCeleb2 dataset [6] for training our audio encoder,  $E_{aud}$ . During each training iteration, we randomly select an audio clip, varying in length from 50 to 200 frames. The corresponding audio segment and the initial frame of this clip are fed into  $E_{aud}$  as inputs. All subsequent frames from the clip are utilized as reference frames (ground truth) for training purposes. Both the training and evaluation processes are conducted using a resolution of  $512 \times 512$ .

## 9.3. Utilized Loss Functions

Training of  $E_{aud}$  incorporates three distinct types of loss functions:

**Photometrical Identity Preservation Losses:** To ensure the identity of the individual in the video is preserved, we apply an L1 loss between the predicted image (using output expression vectors from  $E_{aud}$ , denoted as  $\hat{x}_n^{aud}$ ) and the

ground truth image  $x_n$ . Additionally, we implement a perceptual loss comparing facial features in these images using the *VGGFace* model [24]. The identity preservation loss is represented as:

$$\mathcal{L}_{idt} = w_{L1}\mathcal{L}_{L1} + w_{face}\mathcal{L}_{face} \quad (12)$$

**Latent Mouth Movement Losses:** For accurately translating mouth movements, we use an L1 loss focusing on the principal components related to mouth movements in  $\hat{z}_n$  and  $z_n$ . This loss, detailed in Sec. 5.2, is denoted as  $\mathcal{L}_{PCA}(z_i, z_j, n)$ , with  $n = 8$  in our experiments. A secondary L1 loss,  $\mathcal{L}_{vtr}$ , with a reduced weight, ensures a closer match between  $\hat{z}_n$  and  $z_n$  vectors.

$$\mathcal{L}_{latent} = w_{PCA}\mathcal{L}_{PCA} + w_{vtr}\mathcal{L}_{vtr} \quad (13)$$

**Photometrical Lip Movement Losses:** For enhanced translation of mouth movements, we employ the FaceParsing network [38] to generate masks for the upper lip, lower lip, and inner mouth regions in both the predicted ( $\hat{x}_n^{aud}$ ) and ground truth ( $x_n$ ) images. These masks are compared using the Binary Cross Entropy loss,  $\mathcal{L}_{BCE}$ . Additionally, we extract pixels corresponding to the mouth in both predicted and actual images, comparing them using an L1 loss,  $\mathcal{L}_{lips}$ :

$$\mathcal{L}_{mouth} = w_{BCE}\mathcal{L}_{BCE} + w_{lips}\mathcal{L}_{lips} \quad (14)$$

The overall loss function used to train our model is the cumulative sum of these individual losses:

$$\mathcal{L}_{speech} = \mathcal{L}_{idt} + \mathcal{L}_{latent} + \mathcal{L}_{mouth} \quad (15)$$

We utilized the AdamW optimizer [21] with a cosine learning rate schedule. The initial learning rate was gradually reduced from  $1 \times 10^{-4}$  to zero over the training iterations. The hyperparameters for the losses were set as follows:  $w_{L1} = 10$ ,  $w_{face} = 100$ ,  $w_{PCA} = 200$ ,  $w_{vtr} = 5$ ,  $w_{BCE} = 5 * 10^3$ ,  $w_{lips} = 5 * 10^5$ .

## 10. Generating head rotations

To control the main model using head rotations generated from speech, we have developed a generative adversarial model. This model takes a speech recording as input and outputs a series of rotations, as depicted in Fig 13.

### 10.1. Rotation Representation

We represent the 3D rotations with six dimensions as described in [47]. This ensures the continuity of the representation, which is more suitable for learning. We also add an extra three parameters to predict the translation. Therefore, we can formulate a head transformation sequence  $X$  as a sequence of rotations and translation across  $T$  consecutive frames,  $X \in R^{T \times 9}$  where each  $X_t \in R^9$  is a vector representing the transformation from the reference frame. To map

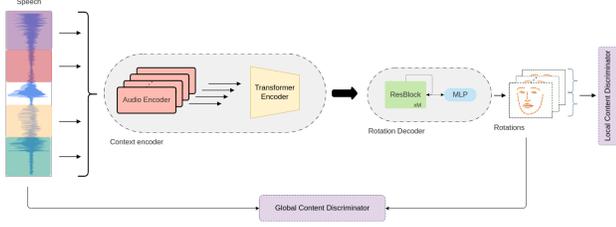


Figure 13. **Audio-to-rotations model.** The input signal is split into overlapping segments for the audio encoder, processed through a Transformer and a ResNet decoder. The resulting rotation sequence and audio input are then given to the Global Content Discriminator, while the Local Temporal Discriminator receives smaller segments of these rotations.

the 6D representation again to the 3D rotation group, we can use the following formula:

$$f_{GS} \left( \begin{bmatrix} | & | \\ a_1 & a_2 \\ | & | \end{bmatrix} \right) = \begin{bmatrix} | & | & | \\ b_1 & b_2 & b_3 \\ | & | & | \end{bmatrix} \quad (16)$$

$$b_i \left[ \begin{cases} N(a_1) & \text{if } i = 1 \\ N(a_2 - (b_1 \cdot a_2)b_1) & \text{if } i = 2 \\ b_1 \times b_2 & \text{if } i = 3 \end{cases} \right]^T \quad (17)$$

Here  $N(\cdot)$  denotes a normalization function and  $f_{GS}$  a Gram-Schmidt process. The model produces a  $3 \times 2$  matrix with  $a_1, a_2$  being its columns. The Gram-Schmidt process in Equation (17) produces the third column  $b_3$  by taking the cross product of the two first columns  $b_1$  and  $b_2$ , making it normal to the plane containing them. This process ensures that the resulting  $3 \times 3$  matrix is orthogonal. The remaining three dimensions map directly to the translations by construction.

## 10.2. Generator

Our generator is split into two components: a context encoder and a head pose decoder. The context encoder merges an audio encoder, Whisper, with a transformer encoder for temporal analysis. This setup efficiently leverages existing audio embeddings from other parts of our system. The head pose is decoded using the encoder’s hidden states, processed through ResNet layers and an MLP. This converts  $R^{T \times H}$  to  $R^{T \times 9}$ , where  $H$  is the Transformer’s hidden size.

## 10.3. Discriminators

To evaluate our generated rotations, we use two discriminators assessing local and global coherence.

**Local Content Discriminator.** Inspired by Isola et al. [15], we use a 1D temporal PatchGAN variant. This discriminator targets patch-level structures, classifying patches of  $N$  frames as real or fake. With the discriminator convolution spanning the entire sequence, averaging all responses yields

the final judgment. We found  $N=8$  optimizes frame-to-frame coherence.

**Global Content Discriminator.** This discriminator assesses the full sequence’s coherence with the audio input. We encode the rotation sequence using 1D ResNet blocks and a global pooling layer, then concatenate it with the audio embeddings from Whisper. A final MLP layer determines if the sequence is authentic or generated.

## 10.4. Losses

To train the model, we use a weighted combination of different losses. The resulting loss is described as follows:

$$\mathcal{L}_{tot} = \lambda_{recons} \mathcal{L}_{recons} + \lambda_{adv} \mathcal{L}_{adv} + \lambda_{smooth} \mathcal{L}_{smooth} \quad (18)$$

**Reconstruction loss  $\mathcal{L}_{recons}$ .** Is a  $L_1$  loss between predicted  $X$  and ground-truth  $Y$  head poses. The head poses can be separated into rotation  $r$  and translation  $t$ .

$$\mathcal{L}_{recons} = \sum_{i=1}^T |r_i - \hat{r}_i| + \sum_{i=1}^T |t_i - \hat{t}_i| \quad (19)$$

We then have two different coefficient factors  $\lambda_{rot}$  and  $\lambda_{trans}$  for the reconstruction of rotations and translations.

**Smoothing loss  $\mathcal{L}_{smooth}$ .** Acts as a regularization loss and ensures smoothness over consecutive frames.

$$\mathcal{L}_{smooth} = \sum_{i=2}^T |X_i - X_{i-1}| \quad (20)$$

**Adversarial loss  $\mathcal{L}_{adv}$ .** We adopt WGAN-GP [13] for improved stability of the training and for avoiding mode collapse. The discriminator and generator losses are as follows:

$$\mathcal{L}_D = \mathbb{E}_{\hat{x} \sim \mathbb{P}_g} [D(\hat{x})] - \mathbb{E}_{x \sim \mathbb{P}_r} [D(x)] + \lambda \mathbb{E}_{\hat{x} \sim \mathbb{P}_{\hat{x}}} [(\|\nabla_{\hat{x}} D(\hat{x})\|_2 - 1)^2] \quad (21)$$

$$\mathcal{L}_G = -\mathbb{E}_{\hat{x} \sim \mathbb{P}_g} [D(\hat{x})] \quad (22)$$

$\mathbb{P}_{\hat{x}}$  is defined as a uniform sampling along straight lines between pairs of points sampled from the data distribution  $\mathbb{P}_r$  and the generator distribution  $\mathbb{P}_g$ . We set  $\lambda$  to 10 and update five times the discriminator for every single update of the generator as suggested in the original paper.

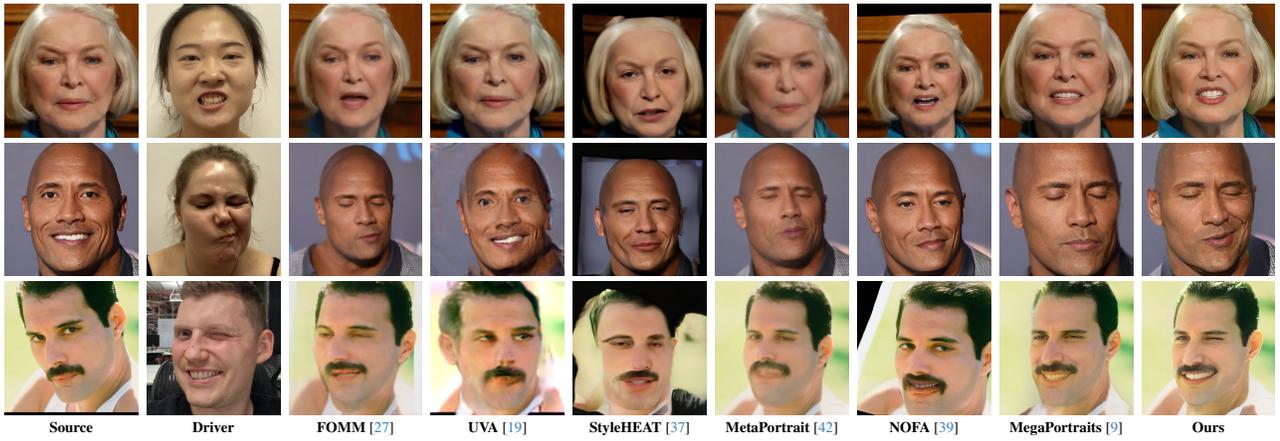


Figure 14. An additional qualitative comparison of head avatar systems in cross-reenactment scenario.

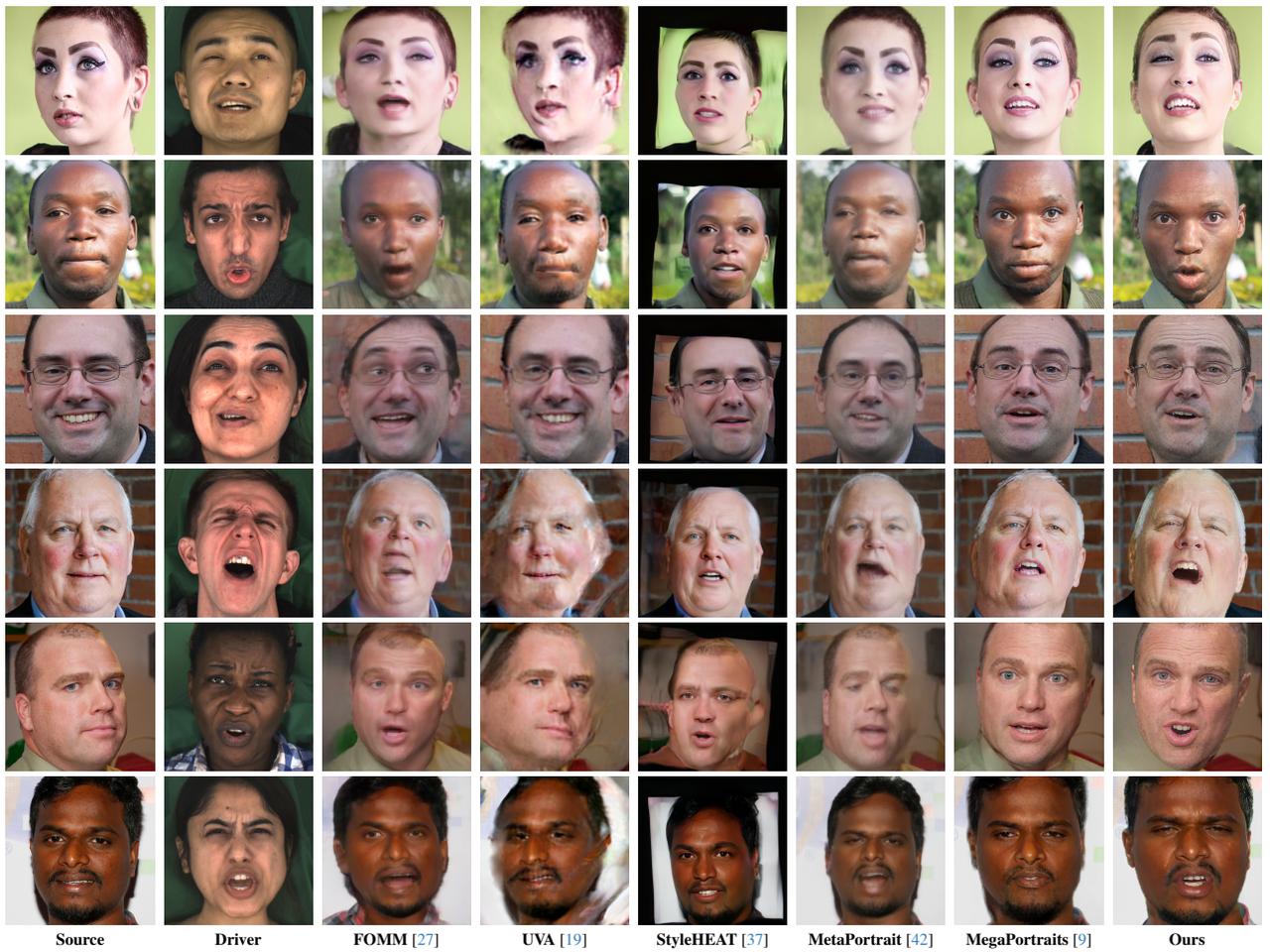


Figure 15. An additional qualitative comparison of head avatar systems in cross-reenactment scenario.

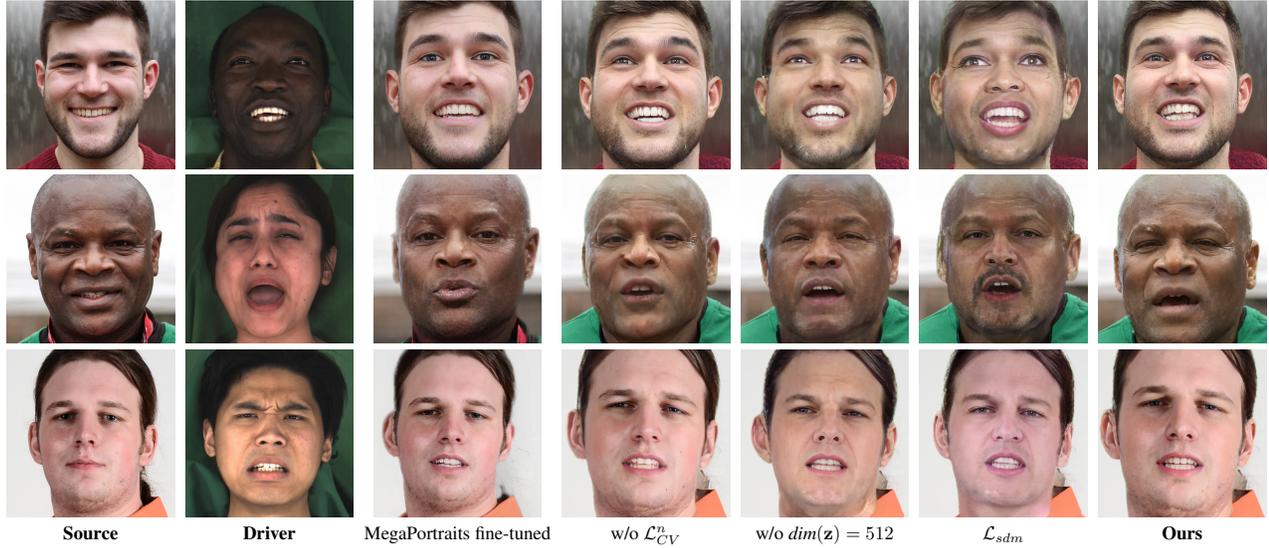


Figure 16. Visual comparison for our ablation study Fig. 16.

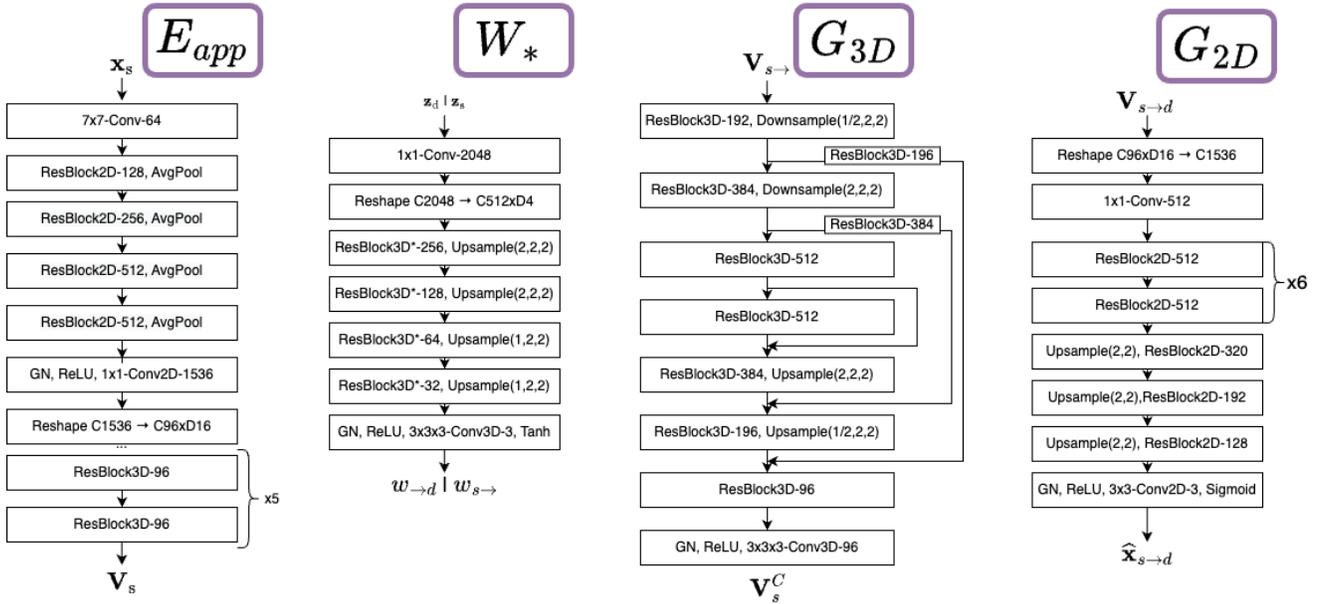


Figure 17. Architectures of main components of our main model