

Efficient Region-Aware Neural Radiance Fields for High-Fidelity Talking Portrait Synthesis

Jiahe Li¹, Jiawei Zhang¹, Xiao Bai^{1*}, Jun Zhou², Lin Gu^{3,4}

¹School of Computer Science and Engineering, State Key Laboratory of Software Development Environment, Jiangxi Research Institute, Beihang University

²School of Information and Communication Technology, Griffith University

³RIKEN AIP ⁴The University of Tokyo

Abstract

This paper presents ER-NeRF, a novel conditional Neural Radiance Fields (NeRF) based architecture for talking portrait synthesis that can concurrently achieve fast convergence, real-time rendering, and state-of-the-art performance with small model size. Our idea is to explicitly exploit the unequal contribution of spatial regions to guide talking portrait modeling. Specifically, to improve the accuracy of dynamic head reconstruction, a compact and expressive NeRF-based Tri-Plane Hash Representation is introduced by pruning empty spatial regions with three planar hash encoders. For speech audio, we propose a Region Attention Module to generate region-aware condition feature via an attention mechanism. Different from existing methods that utilize an MLP-based encoder to learn the cross-modal relation implicitly, the attention mechanism builds an explicit connection between audio features and spatial regions to capture the priors of local motions. Moreover, a direct and fast Adaptive Pose Encoding is introduced to optimize the head-torso separation problem by mapping the complex transformation of the head pose into spatial coordinates. Extensive experiments demonstrate that our method renders better high-fidelity and audio-lips synchronized talking portrait videos, with realistic details and high efficiency compared to previous methods. Code is available at <https://github.com/Fictionarry/ER-NeRF>.

1. Introduction

Audio-driven talking portrait synthesis is an important and challenging issue with several applications such as digital humans, virtual avatars, film-making, and video conferencing. Over the past few years, many researchers have tackled the task with deep generative models [10, 32, 41, 58,

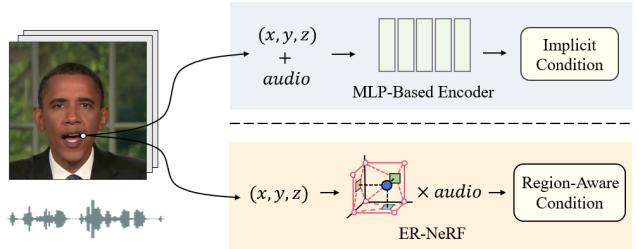


Figure 1. Instead of learning the implicit audiovisual relation by an MLP-based encoder, we explicitly attend to the cross-modal interaction between speech audio and spatial regions. Region awareness enables ER-NeRF to render more accurate facial motions.

57, 29, 50]. Recently, Neural Radiance Fields (NeRF) [30] is introduced into audio-driven talking portrait synthesis. It provides a new way to learn a direct mapping from the audio feature to the corresponding visual appearance by a deep multi-layer perceptron (MLP). Since then, several studies condition NeRF on audio signals in an end-to-end way [23, 28, 33, 47] or by some intermediate representations [48, 7] to reconstruct a specific talking portrait. Though these vanilla NeRF-based methods have shown great success in the synthesis quality, the inference speed is far from meeting real-time requirements, which seriously limits their practical applications.

Several recent works on efficient neural representation have demonstrated tremendous speedups over vanilla NeRF by replacing part of the MLP network with sparse feature grids [38, 31, 6, 8, 18, 5, 19]. Instant-NGP [31] introduces a hash-encoded voxel grid for static scene modeling, allowing fast speed and high-quality rendering with a compact model. RAD-NeRF [40] first applies this technique to talking portrait synthesis and builds a real-time framework with state-of-the-art performance. However, this approach requires a complex MLP-based grid encoder to learn the regional audio-motion mapping implicitly, which limits its convergence and reconstruction quality.

*Corresponding author: Xiao Bai (baixiao@buaa.edu.cn).

This paper aims to explore a more effective solution for efficient and high-fidelity talking portrait synthesis. Based on previous studies, we notice that different spatial regions contribute unequally to representing talking portraits: (1) In volumetric rendering, since only the surface regions contribute to representing the dynamic head, most other spatial regions are empty and can be pruned with some efficient NeRF techniques to reduce the training difficulty; (2) As the fact that different facial areas have varying associations with speech audio [28], different spatial regions are inherently related to the audio signal in their own distinct manners and lead to unique audio-driven local motions. Inspired by these observations, we *explicitly exploit the unequal contribution of spatial regions to guide the talking portrait modeling*, and present a novel Efficient Region-aware talking portrait NeRF (**ER-NeRF**) framework for realistic and efficient talking portrait synthesis, which achieves high-quality rendering, fast convergence, and real-time inference with small model size.

Our first improvement focuses on the dynamic head representation. Although RAD-NeRF [40] has leveraged Instant-NGP to represent the talking portrait and achieves a fast inference, its rendering quality and convergence are hampered by hash collisions when modeling the 3D dynamic talking head. To address this problem, we introduce a *Tri-Plane Hash Representation* that factorizes the 3D space into three orthogonal planes via a NeRF-based tri-plane decomposition [6]. During the factorization, all spatial regions are squeezed onto 2D planes, with the corresponding feature grids pruned. Hence, hash collisions only occur in low-dimensional subspaces and are reduced in number. With fewer noises, the network can pay more attention to processing audio features, leading to the capability of reconstructing more accurate head structures and finer dynamic motions.

To capture the regional impact of audio signals, we further explore the relevance between the audio feature and position encoding of the proposed Tri-Plane Hash Representation. Instead of concatenating the raw features and learning the audiovisual correlation by a large MLP-based encoder, we propose a *Region Attention Module* that adjusts the audio feature to best fit certain spatial regions via a cross-modal attention mechanism. Hence, the dynamic parts of the portrait can acquire more appropriate features to model accurate facial movements, while other static portions remain unaffected by the changing signals. By gaining regional awareness, high-quality and efficient modeling for local motions can be achieved.

Moreover, a simple but effective *Adaptive Pose Encoding* is proposed in our framework to solve the head-torso separation problem. It maps the complex pose transformation onto spatial coordinates and provides a clearer position relation for torso-NeRF to learn its own pose implicitly.

The main contributions of our work are summarized as

follows:

- We introduce an efficient *Tri-Plane Hash Representation* to facilitate dynamic head reconstruction, which also achieves high-quality rendering, real-time inference and fast convergence with a compact model size.
- We propose a novel *Region Attention Module* to capture the correlation between the audio condition and spatial regions for accurate facial motion modeling.
- Extensive experiments show that the proposed ER-NeRF renders realistic talking portraits with high efficiency and visual quality, which outperforms state-of-the-art methods on both objective evaluation and human studies.

2. Related Work

2D-Based Talking Portrait Synthesis. Driving talking portraits by arbitrary speech audio is an active research topic in computer vision and computer graphics. This task aims to reenact the specific person with high image quality and audio-visual consistency. Conventional methods [4, 3] define phoneme-mouth correspondence rules and stitch the mouth shapes. Early deep learning-based methods focus on synthesizing the audio-synchronized lip motions for a given facial image [32, 17, 27, 10, 46]. Later, to enhance controllability, intermediate representations like facial landmarks and 3D facial models are utilized in several multi-stage methods [41, 44, 50, 29]. However, extra errors and information losses would occur in the estimation of these intermediate representations. More recently, diffusion models have been used to improve lip-sync and image quality [49, 34, 37], but they are slow in inference. Due to the lack of an explicit 3D structure representation, 2D-based methods have drawbacks in the naturalness and consistency of head pose control.

NeRF-based Talking Portrait Synthesis. 3D vision techniques aim to learn the 3D structure from images and videos relying on multi-view correspondence, and have been widely developed in many areas [45, 53, 42, 52, 55]. Recently, Neural Radiance Fields (NeRF) [30] has been applied to tackle 3D head structure problems in audio-driven talking portrait synthesis. Earlier works [23, 47, 33, 28] are mainly built on a vanilla NeRF renderer, making them slow and costly for memory. Among them, SSP-NeRF [28] is the first to consider the different impacts of audio on facial areas and adopts a semantic sampling strategy to encourage local motion modeling. By applying Instant-NGP [31], RAD-NeRF [40] has made huge improvements in visual quality and efficiency. Nevertheless, it requires a complex module to handle audio signals. These end-to-end methods take the whole or part of a large MLP network as the encoder to learn the connection between audio and regions, increasing their complexity and training difficulty. Some multi-stage methods [48, 7] pre-train a model to learn the audio-

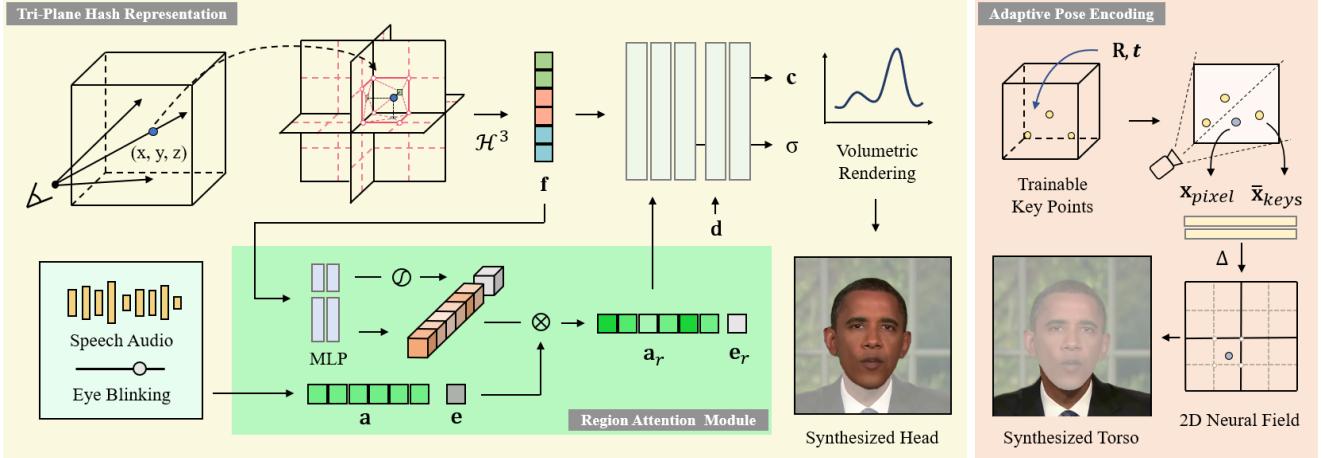


Figure 2. Overview of ER-NeRF framework. The head part of the talking portrait is modeled by the Tri-Plane Hash Representation. A tri-plane hash encoder \mathcal{H}^3 is used to encode the 3D coordinate \mathbf{x} into its spatial geometry feature \mathbf{f} . The input condition features of speech audio \mathbf{a} and eye blinking \mathbf{e} are reweighted in channel-level with the Region Attention Module and converted to region-aware condition features \mathbf{a}_r and \mathbf{e}_r . Then the region-aware features combined with spatial geometry feature \mathbf{f} and the view direction \mathbf{d} are input into an MLP decoder to predict the color \mathbf{c} and density σ of the head. The torso part is rendered by another torso-NeRF with the Adaptive Pose Encoding. The corresponding head pose $\mathbf{P} = (\mathbf{R}, \mathbf{t})$ is applied to transform the trainable key points to get their normalized 2D coordinates $\bar{\mathbf{X}}_{keys}$, which conditions a certain 2D Neural Field to predict the torso image.

visual relation by intermediate representations, and utilize a NeRF-based renderer for image generation. However, they are inefficient due to the complex architecture. This paper proposes an efficient NeRF-based method that significantly improves visual quality and audio-lips synchronization.

Efficient Neural Representation. Many reported works focus on the efficiency of NeRF. Recently, several hybrid explicit-implicit representations [6, 8, 31, 38, 20] are proposed for static scene reconstruction and strike a balance between speed and memory cost. In these methods, a high-dimensional scene would be separated and stored into sparse feature grids. Plane-based approaches [6, 8] factorize the space into multiple low-dimensional planes and vectors to get a compact representation. Instant-NGP [31] employs multiple hash tables to store the sparse details in multiresolution, assuming most empty regions have been pruned, which hugely improves memory utilization and rendering quality as well. Although the size of each hash map is usually insufficient for representing all positions in the space, the method does not handle the hash collision explicitly but leaves it to the MLP decoder. These methods are mainly designed for static scenes and are incapable of generating dynamic representation. In the field of dynamic NeRFs, current efficient methods are either focused on how to rebuild the scene along the timeline [5, 19, 36, 8, 18, 43] or can only control some simple deformations [51], both of which are unsuitable for modeling audio-driven talking portrait. By leveraging the advantages of the hash and plan-based methods, we introduce an efficient representation for high-quality dynamic head modeling that achieves fast training and inference with small model size.

3. Method

3.1. Preliminaries and Problem Setting

Given a set of multi-view images and camera poses, NeRF [30] represents a static 3D scene with an implicit function $\mathcal{F} : (\mathbf{x}, \mathbf{d}) \rightarrow (\mathbf{c}, \sigma)$, where $\mathbf{x} = (x, y, z)$ is the 3D spatial coordinate and $\mathbf{d} = (\theta, \phi)$ is the viewing direction. The output $\mathbf{c} = (r, g, b)$ denotes the emitted color and σ is the volume density. The color $C(\mathbf{r})$ of one pixel crossed by the ray $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$ from camera center \mathbf{o} can be calculated by aggregating the color \mathbf{c} along the ray:

$$\hat{C}(\mathbf{r}) = \int_{t_n}^{t_f} \sigma(\mathbf{r}(t)) \cdot \mathbf{c}(\mathbf{r}(t), \mathbf{d}) \cdot T(t) dt, \quad (1)$$

where t_n and t_f are the near and far bounds. $T(t)$ is the accumulated transmittance from t_n to t :

$$T(t) = \exp(- \int_{t_n}^t \sigma(\mathbf{r}(s)) ds). \quad (2)$$

In hash grid-based NeRF [31], a multiresolution hash encoder \mathcal{H} is utilized to encode the spacial point by its coordinate \mathbf{x} . Therefore, conditioned with the audio feature \mathbf{a} , the basic implicit function of hash NeRF-based audio-driven talking portrait synthesis can be formulated as:

$$\mathcal{F}^A : (\mathbf{x}, \mathbf{d}, \mathbf{a}; \mathcal{H}) \rightarrow (\mathbf{c}, \sigma). \quad (3)$$

In this paper, we adopt the same basic setting as previous NeRF-based works [23, 28, 40]. Specifically, we use a few minutes of single-person video as the training data, which is captured from the front view by a static camera. The camera’s intrinsic and extrinsic parameters for each frame are calculated from the head poses, which are estimated by

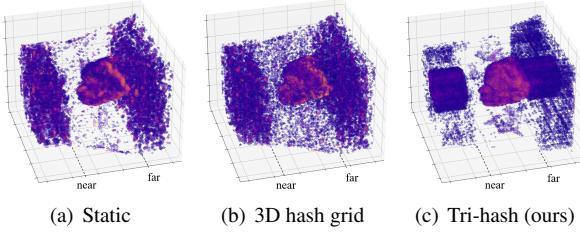


Figure 3. **The visualized occupancy grids.** We show the predicted head surfaces according to σ . (a) 3D hash grid without audio condition. (b, c) 3D hash grid and our tri-plane hash representation conditioned with audio. The MLP decoder of the 3D hash grid becomes overloaded after being required to handle audio features and learn the dynamic motions at the same time, while our representation can still reconstruct the fine surface.

a 3DMM model. Audio features are extracted from a pre-trained DeepSpeech [24] model. We also employ an off-the-shelf semantic parsing method to separate the head, torso, and background for various usages. Moreover, we train and render the head and torso separately for acceleration.

3.2. Tri-Plane Hash Representation

Instant-NGP [31] utilizes a set of hash tables to reduce the number of feature grids for efficient neural representation. Inspired by this idea, RAD-NeRF [40] is developed as a real-time and high-quality talking portrait synthesis framework, which leveraged the hash map to represent the small number of surface regions for the portrait head in multiresolution. However, a general 3D hash grid representation is not natively suitable for our task.

A particular problem is the hash collision. Hashing in Instant-NGP treats every position in 3D space equally, which enhances its expressive ability for complex scenes. Nevertheless, the number of hash collisions linearly increases with the number of sampling points, which makes it a burden for the MLP decoder to solve the conflicting gradients. This problem has little effect when reconstructing static scenes, but for talking portrait synthesis, it gets serious when the MLP decoder needs to handle multiple audio features at the same time, as illustrated in Fig. 3.

Factorization for Hash Grid. Since fewer sampling points always mean lower quality, it's hard to solve this problem by directly reducing the sampling number per ray. Another thinking is to avoid hash collisions from high dimensions. As previous works have proved that a static 3D space of the head can be represented by three 2D tensors [6], it's possible to squeeze the dynamic talking head into several low-dimensional subspaces with little information loss. From this perspective, we factorize the 3D spatial feature volume into three orthogonal 2D hash grids.

For a given coordinate $\mathbf{x} = (x, y, z) \in \mathbb{R}^{XYZ}$, we separately encode its projected coordinates by three 2D-

multiresolution hash encoders [31]:

$$\mathcal{H}^{\mathbf{AB}} : (a, b) \rightarrow \mathbf{f}_{ab}^{\mathbf{AB}} \quad (4)$$

where the output $\mathbf{f}_{ab}^{\mathbf{AB}} \in \mathbb{R}^{LF}$ is the plane-level geometry feature for the projected coordinate (a, b) and $\mathcal{H}^{\mathbf{AB}}$ is the multiresolution hash encoder for plane $\mathbb{R}^{\mathbf{AB}}$, with the number of levels L , feature dimensions per entry F . Then we concatenate the results to get the final geometry feature $\mathbf{f}_g \in \mathbb{R}^{3 \times LF}$:

$$\mathbf{f}_x = \mathcal{H}^{XY}(x, y) \oplus \mathcal{H}^{YZ}(y, z) \oplus \mathcal{H}^{XZ}(x, z). \quad (5)$$

The symbol \oplus denotes the concatenation operator that concatenates features into a $3 \times LF$ -channel vector.

Our proposed factorization significantly reduces hash collision, as now the collision only occurs in 2D planes. Assuming a common situation that the query rays are almost perpendicular to the frontal plane, the collision can be reduced from $O(R^2N)$ to $O(R^2 + 2RN)$, where R^2 is the number of target pixels and N is the sampling number. With a usual setting of $N = 16$ and $R \approx 256$ in RAD-NeRF [40], our representation can ideally achieve a $5\times$ reduction in hash collision with the same model size. This reduction enables the MLP decoder to focus more on processing audio features, leading to improved convergence and dynamic rendering quality.¹

Overall Head Representation. The input to the MLP decoder consists of \mathbf{f}_x , the view direction \mathbf{d} and a dynamic condition feature set \mathcal{D} including audio feature. The implicit function of the tri-plane hash representation can be formulated as:

$$\mathcal{F}^{\mathcal{H}} : (\mathbf{x}, \mathbf{d}, \mathcal{D}; \mathcal{H}^3) \rightarrow (\mathbf{c}, \sigma), \quad (6)$$

where $\mathcal{H}^3 : \mathbf{x} \rightarrow \mathbf{f}_x$ denotes a tri-plane hash encoder consisting of all of three planar hash encoders in Eq. 4.

3.3. Region Attention Module

Dynamic conditions like audio seldom influence the whole portrait equally. Hence, learning how these conditions affect different regions of the portrait is essential for generating natural facial movements. Many previous works [23, 28, 47] ignore this point at the feature level and use some costly approaches to learn the correlation implicitly. By leveraging the multi-resolution regional information stored in the hash encoder, we introduce a lightweight region attention mechanism to explicitly fetch the relations between the dynamic feature and different spatial regions.

Region Attention Mechanism. The region attention mechanism involves an external attention step to calculate the attention vector and a cross-modal channel attention step for reweighting. We aim to connect the dynamic condition feature with the multiresolution geometry feature $\mathbf{f}_x \in \mathbb{R}^N$, which is encoded by the hash encoder \mathcal{H} for a spatial point \mathbf{x} . However, since this hierarchical feature is

constructed by concatenation, no explicit information flow exists during encoding.

To improve the regional information exchange between different levels of \mathbf{f}_x efficiently, and further discriminate the importance of audio for each region via the norm of the attention vector, we use a two-layer MLP to capture the global context of the space. Hence it can be explained as the form of external attention mechanism [22] with two external memory units M_k and M_v for individual levels connection and self-condition query:

$$A = \text{ReLU}(FM_k^T), \quad (7)$$

$$V_{out} = AM_v.$$

where vector \mathbf{f}_x is viewed as an matrix $F \in \mathbb{R}^{N \times 1}$.

Then, similar to the channel attention mechanism proposed by Hu et al. [26], we treat the resulting feature $V_{out} \in \mathbb{R}^{O \times 1}$ as the region attention vector $\mathbf{v} \in \mathbb{R}^O$ to reweight each channel of the dynamic condition feature $\mathbf{q} \in \mathbb{R}^O$. Finally, the output feature vector is:

$$\mathbf{q}_{out} = \mathbf{v} \odot \mathbf{q} \quad (8)$$

where \odot denotes the Hadamard product. The resulting region-aware feature \mathbf{q}_{out} at each channel is related to hierarchical regions where x is located, since the region attention vector \mathbf{v} includes an informative multi-resolution representation of the space. Therefore, the multi-resolution spatial region can decide which part of the information in \mathbf{q} should be kept or enhanced.

Speech Audio. For audio signals, given a query coordinate x and an audio feature $\mathbf{a} \in \mathbb{R}^A$, we calculate the geometry feature of x by the tri-plane hash encoder \mathcal{H}^3 of our tri-plane hash representation. Then we feed it into a two-layer MLP to generate the region attention vector $\mathbf{v}_{a,x} \in \mathbb{R}^A$ for audio with the same number of channels A . After that, channel-wise attention is applied to \mathbf{a} by $\mathbf{v}_{a,x}$:

$$\mathbf{v}_{a,x} = \text{MLP}_a(\mathcal{H}^3(x)), \quad (9)$$

$$\mathbf{a}_{r,x} = \mathbf{v}_{a,x} \odot \mathbf{a}.$$

During training, in regions that vary with the audio, the attention vector $\mathbf{v}_{a,x}$ is optimized for better utilization of the audio feature \mathbf{a} . Instead, for the static parts, the audio conditions are considered noises and $\mathbf{v}_{a,x}$ is going to be a zero vector to help denoising the useless information.

Eye Blinking. We also apply the mechanism for explicit eye blinking control. We use a scalar to describe the action of eye blinking and regard it as a vector e with one dimension. Differently, the region attention vector $\mathbf{v}_e \in \mathbb{R}^1$ for eye blinking is output by a sigmoid layer:

$$\mathbf{v}_{e,x} = \text{MLP}_e(\mathcal{H}^3(x)), \quad (10)$$

$$\mathbf{e}_{r,x} = e \cdot \text{Sigmoid}(\mathbf{v}_{e,x}).$$

The result $\mathbf{e}_{r,x}$ is scaled by $\mathbf{v}_{e,x}$ according to its geometry position. In the region of the eyes, $\mathbf{e}_{r,x}$ conditions the ap-

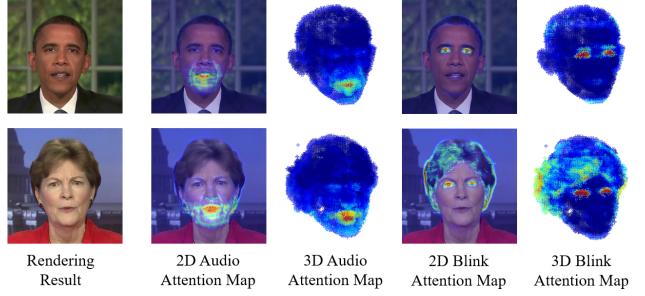


Figure 4. **Visualization of Region Attention Module.** Even if influenced by some uncertain details like fluffy hair, our region attention module successfully captures the relation between dynamic conditions and spatial regions without explicit annotation.

pearance significantly and is close to e for maximizing its effect. Otherwise, it tends to become 0 to reduce the negative interference.

3.4. Training Details

Adaptive Pose Encoding. To solve the head-torso separation problem, we make an improvement based on previous works [40, 48]. Instead of directly using the whole image or pose matrix as the condition, we map the complex transformation of the head pose into the coordinates of several key points that have clearer position information, and lead the torso-NeRF to learn an implicit torso pose from these coordinates.

The encoding process is very simple. We initialize N points in the 3D canonical space with trainable homogeneous coordinates $\mathbf{X}_{keys} \in \mathbb{R}^{4 \times N}$ and apply the head pose $\mathbf{P} = (\mathbf{R}, t)$ to transform the key points $\hat{\mathbf{X}}_{keys} = \mathbf{P}^{-1}\mathbf{X}_{keys}$. Then we project $\hat{\mathbf{X}}_{keys}$ onto the image plane and get the 2D coordinates $\bar{\mathbf{X}}_{keys} \in \mathbb{R}^{2 \times N}$ which are the final encoding results to condition the torso-NeRF. In this work, we use $N = 3$ and a 2D deformable neural field [40] to render the pixel-wise color of the torso.¹

Coarse-to-Fine Optimization. We apply a two-staged coarse-to-fine training process for better image quality. At the coarse stage, we follow the vanilla NeRF to use the MSE loss for the predicted color $\hat{C}(\mathbf{r})$ of the image \mathcal{I} :

$$\mathcal{L}_{coarse} = \sum_{i \in \mathcal{I}} \left\| C(i) - \hat{C}(i) \right\|_2^2. \quad (11)$$

Since MSE loss has a weakness in optimizing sharp details, we then apply an overall finetune with LPIPS loss [54]. Similar to RAD-NeRF [40], we randomly sample a set of patches \mathcal{P} from the whole image and combine the LPIPS loss by a weight λ to enhance details:

$$\mathcal{L}_{fine} = \sum_{i \in \mathcal{P}} \left\| C(i) - \hat{C}(i) \right\|_2^2 + \lambda \text{LPIPS}(\hat{\mathcal{P}}, \mathcal{P}). \quad (12)$$

¹Additional descriptions and detailed discussions can be found in the supplementary material.

Methods	PSNR \uparrow	LPIPS \downarrow	FID \downarrow	LMD \downarrow	AUE \downarrow	Sync \uparrow	Time	FPS	Size (MB)
Ground Truth	N/A	0	0	0	0	7.584	-	-	-
Wav2Lip [32]	-	-	31.08	5.124	3.861	8.576	-	19	>400
PC-AVS [57]	18.25	0.2440	101.97	4.816	3.142	8.397	-	32	>500
AD-NeRF [23]	30.75	0.1034	18.60	3.345	2.201	5.205	18h	0.13	5.21
RAD-NeRF [40]	33.13	0.0519	12.05	2.812	2.102	5.052	5h	32	11.8
RAD-NeRF \dagger	33.26	0.0486	12.20	2.802	1.750	5.197	-	-	-
ER-NeRF (Ours)	33.10	0.0291	10.42	2.740	1.629	5.708	2h	34	2.51

\dagger using AU45 and overall LPIPS finetune.

Table 1. **The quantitative results of the head reconstruction setting.** The best results are in **bold**. Since Wav2Lip can see the ground truth during the self-driven evaluation, we provide another clip of video as the image input. Hence PSNR and LPIPS are not valid. The inference FPS of NeRF-based methods is tested on the Obama dataset [23] under the resolution of 450×450 .

4. Experiments

4.1. Experimental Settings

Dataset. For a fair comparison, the dataset for our experiments is obtained from publicly-released video sets [23, 28, 33]. We collect four high-definition speaking video clips with an average length of about 6500 frames in 25 FPS. Each raw video is cropped and resized to 512×512 with a center portrait, except the one from AD-NeRF [23] with the size of 450×450 . A pre-trained DeepSpeech model is used to extract the basic audio feature from the speech audio.

Comparison Baselines. We compare our method with recent representative one-shot and person-specific models, including Wav2Lip [32], PC-AVS [57], NVP [41], LSP [29] and SynObama [39]. In addition, we also compare our method with the three end-to-end NeRF-based models: AD-NeRF [23], SSP-NeRF, and RAD-NeRF [28]. Furthermore, we evaluate our method directly on the Ground Truth to provide a clearer comparison.

Implementation Details. We implement our method on PyTorch. For a specific portrait, we train the head part for 100,000 and 25,000 iterations at the coarse and the fine stage, respectively. In each iteration, we randomly sample a batch of 256^2 rays from one image. Each 2D hash encoder is set with $L = 14$, $F = 1$, and with resolutions from 64 to 512. The torso part is trained separately for another 100,000 iterations. We use AdamW optimizer for both networks with a learning rate of 0.01 for hash encoders and 0.001 for other modules. For the control of eye blinking, we choose AU45 [16] to describe the degree of the action. All experiments are performed on a single RTX 3080Ti GPU. Both the training for the head and torso take about 2 hours.

4.2. Quantitative Evaluation

Metrics. We employ Peak Signal-to-Noise Ratio (**PSNR**) to measure the overall image quality and Learned Perceptual Image Patch Similarity (**LPIPS**) [54] to measure the details. As we have already used the LPIPS during training, for a fair comparison, an additional feature-based loss Fréchet Inception Distance (**FID**) [25] is involved for evaluating image quality. We also utilize the landmark distance

Methods	Testset A		Testset B	
	LMD \downarrow	Sync \uparrow	LMD \downarrow	Sync \uparrow
Ground Truth	0	6.701	0	7.309
Wav2Lip [32]	6.221	8.378	7.393	8.966
PC-AVS [57]	7.112	8.087	7.722	8.565
SynObama [39]	6.540	6.802	-	-
NVP [41]	-	-	7.954	4.313
LSP [29]	5.905	4.287	8.122	5.843
AD-NeRF [23]	<u>6.192</u>	5.195	<u>8.006</u>	4.316
SSP-NeRF [28]	6.332	5.422	-	-
RAD-NeRF [40]	6.357	6.186	8.332	6.680
RAD-NeRF \dagger	6.339	6.119	8.355	6.392
Ours	6.254	<u>6.242</u>	8.150	6.830

\dagger using AU45 and overall LPIPS finetune.

Table 2. **The quantitative results of lip synchronization setting.** The best overall results and the best NeRF-based methods are in **bold** and underline, respectively.

(**LMD**) [9] and SyncNet confidence score (**Sync**) [12, 13] for lip synchronization and action units error (**AUE**) [2, 1] to evaluate face motion accuracy.

Comparison Settings. In quantitative evaluation, we focus on the synthesized quality of the head. Our comparisons are divided into two settings: 1) The *head reconstruction setting*, where we split each video into training and test dataset to evaluate the reconstruction quality of the head for a specific portrait. 2) The *lip synchronization setting*, where we use the audio track of unseen videos to drive all methods for comparisons in lip synchronization.

For the first setting, we use all videos in the collected dataset described in Sec. 4.1 and split each video for both training and evaluation. For the second setting, we extract two audio clips from the public demos of NVP and SynObama, named **Testset A** and **Testset B**. Due to the lack of pre-trained models and codes for NVP, SynObama, and SSP-NeRF, we also get their generated videos from released demos for evaluation. Following previous works [23, 28, 40], we train our method and other baselines on the Obama dataset released with AD-NeRF [23]. For each generated result, we crop and rescaled the facial area into the same size for a fair comparison.



Figure 5. **The comparison of the key frames and details of generated portraits.** We show the generated results of the baselines [32, 57, 23, 40] under the head reconstruction setting and the ground truth. For NeRF-based methods, we also synthesize the torso part for evaluation. Please **zoom in** for better visualization.

Methods	Wav2Lip [32]	PC-AVS [57]	SynObama [39]	LSP [29]	NVP [41]	AD-NeRF [23]	RAD-NeRF [40]	ER-NeRF (Ours)
Lip-sync Accuracy	2.67	2.50	3.56	2.67	2.83	3.25	<u>3.81</u>	4.14
Image Quality	1.92	1.83	4.22	3.83	3.75	3.33	<u>3.69</u>	4.08
Video Realness	1.89	1.83	3.33	2.92	<u>3.50</u>	3.02	3.47	3.86

Table 3. **User Study.** The rating is of scale 1-5, the higher the better. We highlight the **best** and second best results.

Evaluation Results. The results of the *head reconstruction setting* and *lip synchronization setting* are shown in Table 1 and Table 2, respectively. It can be observed that: (1) In the head reconstruction setting, our method achieves the best reconstruction quality in vision and lip synchronization. Although the one-shot methods (Wav2Lip and PC-AVS) perform best in Sync and can synthesize talking heads without per-scene training, they get poor scores in other metrics, which shows that they cannot accurately reconstruct the specific portrait. For a fair comparison, we also apply the overall LPIPS finetune and AU45 [16] to RAD-NeRF to enhance its image quality and eye blinking but cause no obvious improvement in image details. Our ER-NeRF performs the best in most metrics while reaching a higher score than other baselines in Sync. The results show that our method can synthesize realistic portraits with high lip-sync accuracy. (2) In the lip synchronization setting, our method shows an excellent generalization ability to synthesize lip-sync talking portraits. AD-NeRF and SSP-NeRF encounter an over-smoothing lip movement, leading to a high LMD score but low SyncNet confidence. While getting the high-

est Sync score among NeRF-based methods, our method exceeds some representative baselines in lip synchronization. (3) Our method reaches real-time inference, with a faster training time and smaller model size. In Table 1, we report the inference FPS, model size and the time cost for training person-specific models. In comparison, our ER-NeRF achieves the best performance in all three metrics, which demonstrates its high efficiency.

4.3. Qualitative Evaluation

Evaluation Results. For an intuitive comparison of the whole portrait, we show the key frames of a clip and details of four portraits in Figure 5. For NeRF-based methods, we synthesize the torso part to evaluate the whole portrait. The result shows that our ER-NeRF renders more details and has the highest personalized lip-sync accuracy. Although Wav2Lip and PC-AVS achieve a high score in Sync, their generated results have an obvious gap from the ground truth. To evaluate the torso part, all three NeRF-based methods render the torso and head separately. AD-NeRF severely suffers from head-torso separation (yellow

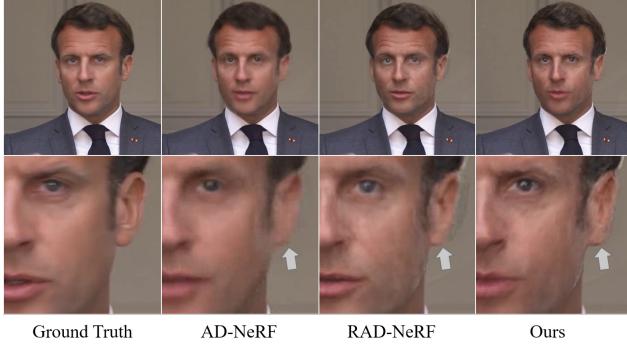


Figure 6. **Evaluation of the out-of-range pose.** Even with a more compact representation, our method can still render accurate structure at a large rotation angle which is rare in the training video.

arrow), and the torso of RAD-NeRF also fails to align with the head sometimes (red arrow). With the same basic representation for the torso as RAD-NeRF, our method demonstrates higher robustness and quality thanks to the design of *Adaptive Pose Encoding*.

In addition, we also compare results with some out-of-range poses, as shown in Figure 6. Despite having pruned most feature grids, our method performs the best in image quality and structure accuracy, which means the robustness and efficiency of our Tri-Plane Hash Representation.

User Study We conducted a user study to better judge the visual quality of the generated heads. We sample 28 generated video clips from the quantitative evaluation, and invite 18 attendees to join the study. The Mean Opinion Scores (MOS) rating protocol is adopted for evaluation and the attendees are required to rate the generated videos from three aspects: (1) Lip-sync Accuracy; (2) Video Realness; (3) Image Quality. The average scores of each method are shown in Table 3. Our ER-NeRF performs the best in lip-sync accuracy and realness. For image quality, only SynObama [39] gets a higher score than our method, which however relies on a large number of training videos and cannot render in real-time. The results show the excellent visual quality of our method for high-fidelity talking portrait synthesis.

4.4. Ablation Study

In this section, we report the ablation study under the head reconstruction setting to prove the effectiveness of our two main contributions. We test settings of different backbones, dynamic feature integration methods and attention targets. The results are shown in Table 4 and Table 5.

Representation. We evaluate three representation backbones on the quality of head reconstruction. The first is an MLP-based network, which is the same as AD-NeRF [23]. For grid-based backbones, we compare our tri-hash representation with pure tri-plane in EG3D [6] and the Instant-NGP [31] 3D hash grid that is used in RAD-NeRF [40]. Due to our specialized architecture, the proposed tri-hash representation achieves the best image quality and makes a significant improvement in lip synchronization.

Backbone	Concat	Att.	PSNR↑	LPIPS↓	LMD↓	AUE↓	Sync↑
MLP	✓		30.75	0.103	3.345	2.201	5.205
Pure	✓		32.11	0.033	2.960	1.812	4.441
Tri-Plane		✓	<u>33.14</u>	<u>0.030</u>	2.825	1.677	5.233
iNGP [31]	✓		33.05	0.031	2.919	1.729	4.664
		✓	33.12	<u>0.030</u>	<u>2.810</u>	1.689	<u>5.257</u>
Tri-Hash	✓		33.25	0.029	2.881	1.634	5.123
		✓	33.10	0.029	2.740	<u>1.646</u>	5.708

Table 4. **Ablation Study** on Tri-Plane Hash Representation and Region Attention Module.

Type	PSNR ↑	LPIPS ↓	LMD ↓	AUE ↓	Sync ↑
Feature-Wise	33.14	0.030	2.781	1.650	5.465
Channel-Wise	33.10	0.029	2.740	1.646	5.708

Table 5. **Ablation Study** on types of attention.

Region Attention Module. We evaluate the region attention mechanism on two backbones compared with directly concatenating. The results show the enormous impact of our method on modeling accurate motions. Note that by only using our attention mechanism with existing backbones, we can get better scores in both image quality and lip synchronization than current state-of-the-art methods with half of the training time and fewer parameters, which shows the high efficiency of our attention mechanism.

Attention Type. In Table 5, we compare three types of attention for the region attention mechanism: feature-wise and channel-wise. Feature-wise attention scales the entire audio feature with a one-dimensional attention vector, while channel-wise reweights each channel, as described in Section 3.3. The outperforming of channel-wise attention indicates that the proposed region attention mechanism successfully captures the distinct impacts of spatial regions and significantly improves lip motion quality.

5. Ethical Consideration

We hope our ER-NeRF can enhance interactive experiences and benefits human beings. However, it could be misused for some malicious purposes. As part of our responsibility, we will share our generated results to help develop stronger deepfake detectors. We believe that the responsible use of this technique can promote the healthy growth of both machine learning research and the digital industry.

6. Conclusion

In this paper, We propose an efficient and effective framework ER-NeRF for high-quality talking portrait synthesis, mainly consisting of a Tri-Plane Hash Representation and a Region Attention Module. Our framework achieves significant improvement in realistic talking portrait synthesis with higher efficiency. Due to the space limitation, we have put the discussion in the supplementary material. We hope our work can benefit human beings and also inspire more novel conditional NeRF techniques.

Acknowledgments. In this work, Jiahe Li, Jiawei Zhang and Xiao Bai are supported by the National Natural Science Foundation of China (No. 62276016), Lin Gu is supported by JST Moonshot R&D Grant Number JPMJMS2011, Japan.

References

- [1] Tadas Baltrušaitis, Marwa Mahmoud, and Peter 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. [6](#)
- [2] Tadas Baltrušaitis, Amir Zadeh, Yao Chong Lim, and Louis-Philippe 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. [6](#)
- [3] Matthew Brand. Voice puppetry. In *Proceedings of the 26th annual conference on Computer graphics and interactive techniques*, pages 21–28, 1999. [2](#)
- [4] Christoph Bregler, Michele Covell, and Malcolm Slaney. Video rewrite: Driving visual speech with audio. In *Proceedings of the 24th annual conference on Computer graphics and interactive techniques*, pages 353–360, 1997. [2](#)
- [5] Ang Cao and Justin Johnson. Hexplane: A fast representation for dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 130–141, 2023. [1, 3](#)
- [6] Eric R Chan, Connor Z Lin, Matthew A Chan, Koki Nagano, Boxiao Pan, Shalini De Mello, Orazio Gallo, Leonidas J Guibas, Jonathan Tremblay, Sameh Khamis, et al. Efficient geometry-aware 3d generative adversarial networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16123–16133, 2022. [1, 2, 3, 4, 8](#)
- [7] Aggelina Chatziagapi, ShahRukh Athar, Abhinav Jain, Rothith Mysore Vijaya Kumar, Vimal Bhat, and Dimitris Samaras. Lipnerf: What is the right feature space to lip-sync a nerf. In *International Conference on Automatic Face and Gesture Recognition 2023*, 2023. [1, 2, 14](#)
- [8] Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. Tensorf: Tensorial radiance fields. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXII*, pages 333–350. Springer, 2022. [1, 3](#)
- [9] Lele Chen, Zhiheng Li, Ross K Maddox, Zhiyao Duan, and Chenliang Xu. Lip movements generation at a glance. In *Computer Vision–ECCV 2018: 15th European Conference, Munich, Germany, September 8–14, 2018, Proceedings, Part VII 15*, pages 538–553. Springer, 2018. [6](#)
- [10] Lele Chen, Ross K Maddox, Zhiyao Duan, and Chenliang Xu. Hierarchical cross-modal talking face generation with dynamic pixel-wise loss. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 7832–7841, 2019. [1, 2](#)
- [11] Liang Chen, Yong Zhang, Yibing Song, Lingqiao Liu, and Jue Wang. Self-supervised learning of adversarial example: Towards good generalizations for deepfake detection. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 18710–18719, 2022. [14](#)
- [12] Joon Son Chung and Andrew Zisserman. Lip reading in the wild. In *Computer Vision–ACCV 2016: 13th Asian Conference on Computer Vision, Taipei, Taiwan, November 20–24, 2016, Revised Selected Papers, Part II 13*, pages 87–103. Springer, 2017. [6](#)
- [13] 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. [6](#)
- [14] Brian Dolhansky, Joanna Bitton, Ben Pflaum, Jikuo Lu, Russ Howes, Menglin Wang, and Cristian Canton Ferrer. The deepfake detection challenge (dfdc) dataset. *arXiv preprint arXiv:2006.07397*, 2020. [14](#)
- [15] Xiaoyi Dong, Jianmin Bao, Dongdong Chen, Ting Zhang, Weiming Zhang, Nenghai Yu, Dong Chen, Fang Wen, and Baining Guo. Protecting celebrities from deepfake with identity consistency transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9468–9478, 2022. [14](#)
- [16] Paul Ekman and Wallace V. Friesen. *Facial Action Coding System: Manual*. Palo Alto: Consulting Psychologists Press, 1978. [6, 7](#)
- [17] Tony Ezzat, Gadi Geiger, and Tomaso Poggio. Trainable videorealistic speech animation. *ACM Transactions on Graphics (TOG)*, 21(3):388–398, 2002. [2](#)
- [18] Jiemin Fang, Taoran Yi, Xinggang Wang, Lingxi Xie, Xiaopeng Zhang, Wenyu Liu, Matthias Nießner, and Qi Tian. Fast dynamic radiance fields with time-aware neural voxels. In *SIGGRAPH Asia 2022 Conference Papers*, pages 1–9, 2022. [1, 3](#)
- [19] Sara Fridovich-Keil, Giacomo Meanti, Frederik Rahbæk Warburg, Benjamin Recht, and Angjoo Kanazawa. K-planes: Explicit radiance fields in space, time, and appearance. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12479–12488, 2023. [1, 3](#)
- [20] Sara Fridovich-Keil, Alex Yu, Matthew Tancik, QinHong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5501–5510, 2022. [3](#)
- [21] Luca Guarnera, Oliver Giudice, and Sebastiano Battiato. Deepfake detection by analyzing convolutional traces. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, pages 666–667, 2020. [14](#)
- [22] Meng-Hao Guo, Zheng-Ning Liu, Tai-Jiang Mu, and Shi-Min Hu. Beyond self-attention: External attention using two linear layers for visual tasks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022. [5](#)
- [23] Yudong Guo, Keyu Chen, Sen Liang, Yong-Jin Liu, Hujun Bao, and Juyong Zhang. Ad-nerf: Audio driven neural ra

- diance fields for talking head synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5784–5794, 2021. 1, 2, 3, 4, 6, 7, 8, 12, 15
- [24] Awni Hannun, Carl Case, Jared Casper, Bryan Catanzaro, Greg Diamos, Erich Elsen, Ryan Prenger, Sanjeev Satheesh, Shubho Sengupta, Adam Coates, et al. Deep speech: Scaling up end-to-end speech recognition. *arXiv preprint arXiv:1412.5567*, 2014. 4, 12
- [25] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017. 6
- [26] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7132–7141, 2018. 5
- [27] Amir Jamaludin, Joon Son Chung, and Andrew Zisserman. You said that?: Synthesising talking faces from audio. *International Journal of Computer Vision*, 127:1767–1779, 2019. 2
- [28] Xian Liu, Yinghao Xu, Qianyi Wu, Hang Zhou, Wayne Wu, and Bolei Zhou. Semantic-aware implicit neural audio-driven video portrait generation. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXVII*, pages 106–125. Springer, 2022. 1, 2, 3, 4, 6, 15
- [29] Yuanxun Lu, Jinxiang Chai, and Xun Cao. Live speech portraits: Real-time photorealistic talking-head animation. *ACM Trans. Graph.*, 40(6), dec 2021. 1, 2, 6, 7, 15
- [30] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *European Conference on Computer Vision*, pages 405–421. Springer, 2020. 1, 2, 3
- [31] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Transactions on Graphics (ToG)*, 41(4):1–15, 2022. 1, 2, 3, 4, 8
- [32] 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. 1, 2, 6, 7, 14, 15
- [33] Shuai Shen, Wanhua Li, Zheng Zhu, Yueqi Duan, Jie Zhou, and Jiwen Lu. Learning dynamic facial radiance fields for few-shot talking head synthesis. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XII*, pages 666–682. Springer, 2022. 1, 2, 6, 12, 13
- [34] Shuai Shen, Wenliang Zhao, Zibin Meng, Wanhua Li, Zheng Zhu, Jie Zhou, and Jiwen Lu. DiffTalk: Crafting diffusion models for generalized audio-driven portraits animation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1982–1991, 2023. 2
- [35] Kaede Shiohara and Toshihiko Yamasaki. Detecting deep-fakes with self-blended images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18720–18729, 2022. 14
- [36] Liangchen Song, Anpei Chen, Zhong Li, Zhang Chen, Lele Chen, Junsong Yuan, Yi Xu, and Andreas Geiger. Nerfplayer: A streamable dynamic scene representation with decomposed neural radiance fields. *arXiv preprint arXiv:2210.15947*, 2022. 3
- [37] Michał Stypulkowski, Konstantinos Vougioukas, Sen He, Maciej Zięba, Stavros Petridis, and Maja Pantic. Diffused heads: Diffusion models beat gans on talking-face generation. *arXiv preprint arXiv:2301.03396*, 2023. 2
- [38] Cheng Sun, Min Sun, and Hwann-Tzong Chen. Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5459–5469, 2022. 1, 3
- [39] Supasorn Suwajanakorn, Steven M Seitz, and Ira Kemelmacher-Shlizerman. Synthesizing obama: learning lip sync from audio. *ACM Transactions on Graphics (ToG)*, 36(4):1–13, 2017. 6, 7, 8, 15
- [40] Jiaxiang Tang, Kaisiyuan Wang, Hang Zhou, Xiaokang Chen, Dongliang He, Tianshu Hu, Jingtuo Liu, Gang Zeng, and Jingdong Wang. Real-time neural radiance talking portrait synthesis via audio-spatial decomposition. *arXiv preprint arXiv:2211.12368*, 2022. 1, 2, 3, 4, 5, 6, 7, 8, 12, 13, 15
- [41] Justus Thies, Mohamed Elgharib, Ayush Tewari, Christian Theobalt, and Matthias Nießner. Neural voice puppetry: Audio-driven facial reenactment. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XVI* 16, pages 716–731. Springer, 2020. 1, 2, 6, 7
- [42] Chen Wang, Xiang Wang, Jiawei Zhang, Liang Zhang, Xiao Bai, Xin Ning, Jun Zhou, and Edwin Hancock. Uncertainty estimation for stereo matching based on evidential deep learning. *Pattern Recognition*, 124:108498, 2022. 2
- [43] Feng Wang, Sinan Tan, Xinghang Li, Zeyue Tian, and Huaping Liu. Mixed neural voxels for fast multi-view video synthesis. *arXiv preprint arXiv:2212.00190*, 2022. 3
- [44] 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 *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXI*, pages 700–717. Springer, 2020. 2
- [45] Xiang Wang, Chen Wang, Bing Liu, Xiaoqing Zhou, Liang Zhang, Jin Zheng, and Xiao Bai. Multi-view stereo in the deep learning era: A comprehensive review. *Displays*, 70:102102, 2021. 2
- [46] Olivia Wiles, A Sophia Koepke, and Andrew Zisserman. X2face: A network for controlling face generation using images, audio, and pose codes. In *Computer Vision–ECCV 2018: 15th European Conference, Munich, Germany, September 8–14, 2018, Proceedings, Part XIII* 15, pages 690–706. Springer, 2018. 2
- [47] Shunyu Yao, RuiZhe Zhong, Yichao Yan, Guangtao Zhai, and Xiaokang Yang. Dfa-nerf: Personalized talking head generation via disentangled face attributes neural rendering. *arXiv preprint arXiv:2201.00791*, 2022. 1, 2, 4

- [48] Zhenhui Ye, Ziyue Jiang, Yi Ren, Jinglin Liu, Jinzheng He, and Zhou Zhao. Geneface: Generalized and high-fidelity audio-driven 3d talking face synthesis. In *The Eleventh International Conference on Learning Representations*, 2022. [1](#), [2](#), [5](#), [13](#), [14](#)
- [49] Zhentao Yu, Zixin Yin, Deyu Zhou, Duomin Wang, Finn Wong, and Baoyuan Wang. Talking head generation with probabilistic audio-to-visual diffusion priors. *arXiv preprint arXiv:2212.04248*, 2022. [2](#)
- [50] Chenxu Zhang, Yifan Zhao, Yifei Huang, Ming Zeng, Saifeng Ni, Madhukar Budagavi, and Xiaohu Guo. Facial: Synthesizing dynamic talking face with implicit attribute learning. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 3867–3876, 2021. [1](#), [2](#)
- [51] He Zhang, Fan Li, Jianhui Zhao, Chao Tan, Dongming Shen, Yebin Liu, and Tao Yu. Controllable free viewpoint video reconstruction based on neural radiance fields and motion graphs. *IEEE Transactions on Visualization and Computer Graphics*, 2022. [3](#)
- [52] Jiawei Zhang, Xiang Wang, Xiao Bai, Chen Wang, Lei Huang, Yimin Chen, Lin Gu, Jun Zhou, Tatsuya Harada, and Edwin R Hancock. Revisiting domain generalized stereo matching networks from a feature consistency perspective. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13001–13011, 2022. [2](#)
- [53] Pengcheng Zhang, Lei Zhou, Xiao Bai, Chen Wang, Jun Zhou, Liang Zhang, and Jin Zheng. Learning multi-view visual correspondences with self-supervision. *Displays*, 72:102160, 2022. [2](#)
- [54] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 586–595, 2018. [5](#), [6](#), [13](#)
- [55] Youmin Zhang, Yimin Chen, Xiao Bai, Suihanjin Yu, Kun Yu, Zhiwei Li, and Kuiyuan Yang. Adaptive unimodal cost volume filtering for deep stereo matching. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 12926–12934, 2020. [2](#)
- [56] Hanqing Zhao, Wenbo Zhou, Dongdong Chen, Tianyi Wei, Weiming Zhang, and Nenghai Yu. Multi-attentional deepfake detection. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2185–2194, 2021. [14](#)
- [57] 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. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4176–4186, 2021. [1](#), [6](#), [7](#), [15](#)
- [58] Yang Zhou, Xintong Han, Eli Shechtman, Jose Echevarria, Evangelos Kalogerakis, and Dingzeyu Li. Makeltalk: speaker-aware talking-head animation. *ACM Transactions On Graphics (TOG)*, 39(6):1–15, 2020. [1](#)

A. Supplementary Material

A.1. Overview

In the supplemental document, we introduce the details of our torso-nerf with Adaptive Head Encoding, model architecture details, user study details, additional experiments and analysis, ethical considerations, and the discussion of this work.

A.2. Torso-NeRF Details

We combine the proposed Adaptive Pose Encoding and the 2D deformable neural field from RAD-NeRF [40] to render the torso part. As described in Section 3.4 of the main paper, we init three points in the 3D canonical space with trainable homogeneous coordinates:

$$\mathbf{X}_{keys} = (\mathbf{x}_{keys}, \mathbf{y}_{keys}, \mathbf{z}_{keys}, \mathbf{1})^T \in \mathbb{R}^{4 \times 3}. \quad (13)$$

For each frame, we form the pose of head \mathbf{P} as:

$$\mathbf{P} = \begin{pmatrix} \mathbf{R} & \mathbf{t} \\ 0 & 1 \end{pmatrix} \quad (14)$$

and apply it to transform the key points:

$$\hat{\mathbf{X}}_{keys} = \mathbf{P}^{-1} \mathbf{X}_{keys}. \quad (15)$$

where $\hat{\mathbf{X}}_{keys}$ is the transformed coordinates. Then we convert $\hat{\mathbf{X}}_{keys}$ to the ordinary coordinates and project them onto the plane $Z = 1$ to calculate their 2D coordinates $\bar{\mathbf{X}}_{keys} \in \mathbb{R}^{2 \times 3}$ on the imaging plane, where

$$\bar{\mathbf{X}}_{keys}(i, j) = \gamma \cdot \hat{\mathbf{X}}_{keys}(i, j) / \hat{\mathbf{z}}_{keys}(j), \quad (16)$$

and γ is the coefficient learned by the network.

The overview of the torso-NeRF is shown in Figure 7. We use $\bar{\mathbf{X}}_{keys}$ to condition the 2D deformable neural field [40] for rendering the pixel-wise color and alpha of the torso at the image pixel coordinate \mathbf{x}_{pixel} . Specifically, to render the pixel at $\mathbf{x}_{pixel} \in \mathbb{R}^2$ on the image, we firstly feed $\bar{\mathbf{X}}_{keys}$ and the pixel coordinate \mathbf{x}_{pixel} into an MLP, and add the output $\Delta\mathbf{x}$ to \mathbf{x}_{pixel} for a 2D deformation. The deformed coordinate is then encoded by the 2D multiresolution hash encoder \mathcal{H}^t . Finally, another MLP is used to calculate the pixel-wise transparency α and color \mathbf{c}_t .

The implicit function of the torso-NeRF can be formulated as:

$$\mathcal{F}^T : (\mathbf{x}_{pixel}, \bar{\mathbf{X}}_{keys}; \mathcal{H}^t) \rightarrow (\mathbf{c}_t, \alpha) \quad (17)$$

During training, the coordinates \mathbf{X}_{key} can be optimized to gain the ability in representing the implicit relationship between the poses of the head and torso. And due to only linear transformations involved during forwarding, the torso quality is improved without a significant increase in the amount of calculation.

Methods	AD-NeRF	RAD-NeRF	ER-NeRF
Stability	1.33	2.89	3.89
Image Quality	2.67	3.33	4.00

Table 6. **User Study of Torso Quality.** The rating is of scale 1-5, the higher the better.

User Study. We also conduct a user study to evaluate the synthesized torso part. We invite the attendees to rate the stability and image quality of generated torsos in the *head reconstruction setting*. To compare our method, we selected AD-NeRF [23] and RAD-NeRF [40] as the baselines since they are the only two NeRF-based methods that can synthesize the torso part and have released their codes. The results are reported in Table 6. We can observe that our ER-NeRF achieves the best both on Stability and Image Quality by just adding a straightforward encoding step *without any deep neural network*, which demonstrates the high efficiency of our Adaptive Pose Encoding.

A.3. Architecture Details

Audio Feature Extractor. In the experiments, we use the pretrained *DeepSpeech* [24] model to extract raw audio features. We then process these features with the same audio attention module as previous NeRF-based works [23, 33, 40], except for changing the output dimension from 64 to 32.

Region Attention Module. The speech audio branch utilizes an attention vector MLP with 2 layers and 64 hidden dimensions. Conversely, the eye-blinking branch employs a 2-layer MLP with only 16 hidden dimensions.

Tri-plane Hash Representation. The 2D hash encoders are configured to have 14 resolution levels and a single entry assigned to each level, with a range of multiple resolutions from 64 to 512. The density MLP decoder contains 3 layers, and the color MLP decoder contains 2 layers, both of which have 64 hidden dimensions.

A.4. User Study Details

The study involves 18 participants with an age range of 20-30 years old. To facilitate more accurate judgments, we combine all generated videos and the ground truth into a single high-resolution video. This allows participants to observe all motions simultaneously. To ensure fairness in the comparison process, we assign a number to each generated result instead of identifying them by their method. Participants are asked to evaluate the three perspectives of the generated portraits: (1) Lip-sync Accuracy; (2) Video Realness; (3) Image Quality. To evaluate the torso-NeRF, we additionally invite the attendees to judge two aspects of the synthesized torso: (1) Stability; (2) Image Quality.

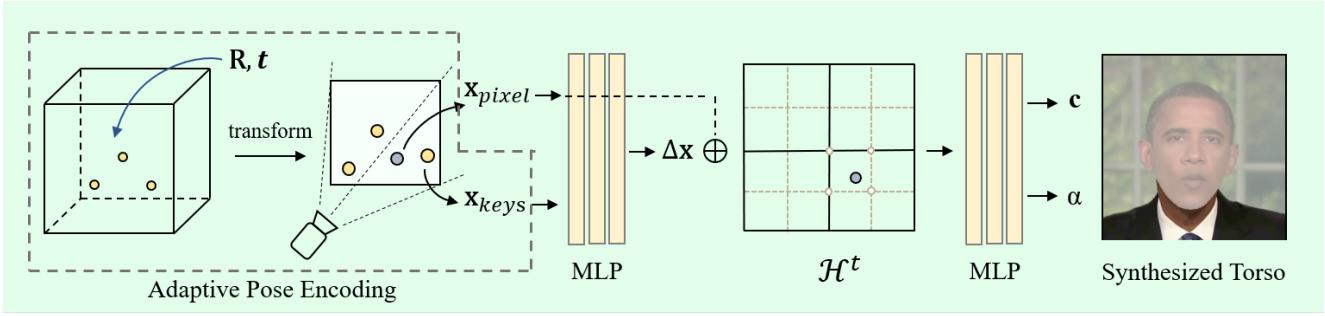


Figure 7. Overview of the Torso-NeRF.

Grid	Instant-NGP	Tri-Hash			
		Frontal	Side 1	Side 2	Total
Collision	835186	138345	31041	26048	195434

Table 7. The number of hash collisions occurring in one feature lookup step on a single grid resolution.

A.5. Tri-Plane Hash Representation

Complexity of Hash Collision Here we give the proof of the complexity $O(R^2 + 2RN)$ in Section 3.2 for our Tri-Hash Representation: 1) For the frontal plane, the projected area is linearly correlated to R^2 , thus the collision is $O(R^2)$; 2) The ideal projected area for the other two side planes is $(\lambda R)R$, where λ is an adjustment. But notice only the nearest N points can be sampled at some side areas due to occlusion, so λR is partly correlated to N , and the collision is $O(\lambda R^2 + RN)$. Overall, $O(R^2 + 2RN)$ is given.

The Number of Hash Collisions. Here we give the evaluation during one lookup step to directly verify our effect on hash collision reduction. The hashtable size is set to 2^{14} and divided by 3 for each planar grid in our Tri-Hash, with the grid resolution of 512, the max in the experiment. Adjustments of 1/8 and 1/4 are applied due to bilinear interpolation. The point coordinates are scaled up to encourage uniform hashing. In practice, the benefit of our method would be more obvious, since indeed the coordinates cannot be uniformly separated among the hash table and so the overlapping of grids becomes more serious.

A.6. Additional Experiments

LPIPS Finetune. It may seem counter-intuitive that the overall LPIPS [54] finetuning is less effective for RAD-NeRF [40] but has a significant impact on the high-frequency details of our ER-NeRF despite having a smaller model size. This phenomenon is likely due to differences in training difficulty. Our ablation study shows that even a simplified architecture with only a 3D hash grid backbone and an audio feature dimension of 32 can reproduce fine details. On the other hand, RAD-NeRF uses a

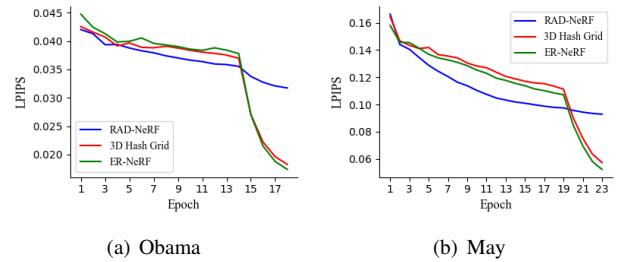


Figure 8. The validation LPIPS loss on our Obama dataset and May dataset with different architectures. A complex network is much harder to be optimized by the LPIPS finetune and reproduce fine details.

more complex architecture with an additional hash grid and higher-dimensional audio features to improve lip-sync performance, which increases the training difficulty and makes the network harder to optimize. As a result, the LPIPS finetuning has a weaker impact on its rendering quality. The variations in LPIPS loss during training are illustrated in Figure 8.

Region Attention for Eye Blinking. We perform an ablation study on the eye-blinking branch of the Region Attention Module in isolation. When we skip the region attention mechanism and directly concatenate the AU45 with the input of the MLP decoder, some unnatural facial movements appear, like jittering and unreasonable lip movements with eye blinking (Figure 9). This might be due to the module’s inability to accurately identify the regional impact of eye blinking and thus learns an incorrect motion mapping with other facial regions. The results indicate that our Region Attention Module can help decouple different semantic motions and improve robustness.

A.7. Comparison with GeneFace and DFRF

In table 8 and 9, we have also compared our ER-NeRF with two current SOTA methods GeneFace [48] and DFRF [33], both of which are designed for different settings, notably. Meanwhile, since the code of GeneFace is released

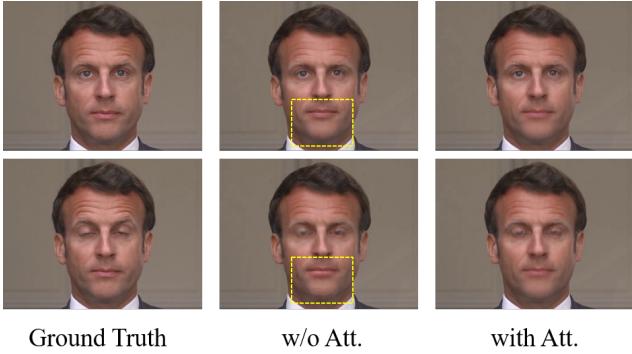


Figure 9. **Ablation on Region Attention for Eye Blinking.** Some unnatural facial movements appear when directly concatenating the AU45 with the input to control eye blinking. After applying the proposed region attention mechanism, the robustness has been improved.

Methods	PSNR \uparrow	LPIPS \downarrow	FID \downarrow	LMD \downarrow	AUE \downarrow	Sync \uparrow
DFRF	30.74	0.0881	13.32	3.553	2.538	4.385
GeneFace	30.24	0.0817	11.16	3.496	2.854	5.403
ER-NeRF (Ours)	33.10	0.0291	10.42	2.740	1.629	5.708

Table 8. DFRF and GeneFace at the *head reconstruction setting*.

Methods	A: LMD \downarrow	A: Sync \uparrow	B: LMD \downarrow	B: Sync \uparrow
DFRF	6.551	4.854	8.126	4.127
GeneFace	5.465	5.849	7.237	6.275
ER-NeRF (Ours)	6.254	6.242	8.150	6.830

Table 9. DFRF and GeneFace at the *lip synchronization setting*.

too close to the submission deadline, it was not taken into the baselines in the main paper. We consider the comparisons not entirely fair for them, and the results are just for reference.

A.8. Additional Qualitative Comparison

We show some additional generated key frames on the Testset A under the *lip synchronization setting* with high resolution in Fig. 10. In this setting, we only synthesize the head part. The results show that our ER-NeRF can outperform most baselines in image quality while retaining a high lip-sync accuracy. We strongly recommend watching our [supplemental video](#) for better visualization and more results.

A.9. Ethics Considerations

Our proposed ER-NeRF synthesizes high-fidelity talking portraits with accurate lip-audio synchronization. The generated portrait video is highly realistic and difficult for people to distinguish fake from real. We hope it can facilitate a wide range of applications, such as digital humans, video production, and human-computer interaction assis-

tance. On the other hand, however, such techniques may be misused for malicious purposes and make harm. It's significant to tell the users whether a video is real or fake. Recent studies have already achieved success in deepfake detection for face swapping, reenactment and other generating videos [21, 56, 14, 11, 35, 15], but it remains a challenge to discriminate synthesized high-fidelity portraits from recent NeRF-based methods. Besides sharing our generated results to the deepfake detection communication and to help develop more powerful deepfake detectors, we also provide some possible perspectives to fight against the malicious use of talking portrait synthesis:

- **Protect real portrait speech videos.** Since current NeRF-based techniques rely heavily on specific training videos, protection for these real videos is valid to prevent misuse. For example, we can add digital watermarks to the portrait part which can be easily detected even in the generated fake videos.
- **Limit the use of deepfake techniques.** Nowadays, little cost of deepfakes leads to an unconstrained use of these techniques. The negative impact of the malicious use of deepfakes can be amplified when they are unintentionally created and shared by the public on social media platforms. Even though the creators may have no malicious intent, the spread of these deepfakes can still have harmful consequences. We suggest the laws should state how to properly make use of these face-generation techniques. On the other hand, the public should also be aware of the potential harm of deepfakes and treat them cautiously.

A.10. Limitation and Future Work

Compared to the one-shot methods like Wav2Lip [32], our method has some advantages in results quality and resolution, however, needs per-scene training when generating new target portraits. Enabling the generative ability may be the target we work for.

Besides, the proposed method has two main limitations. Firstly, our method still encounters a challenge with the small scale of a single training video, leading to a weak lip-audio synchronization with out-of-domain audio, such as some cross-lingual speech or singing voice. Currently, we rely on a pretrained speech recognition model to extract audio features. We have noticed that some recent works [48, 7] employed a pretrained model to enhance their generalizability. In future work, we will consider incorporating priors from large audiovisual datasets to address this limitation. Secondly, although our method has improved the robustness and image quality of the torso part, there remain some blurry regions. We analyze this may be caused by uncertain movements and the form of representation itself. In future work, we will focus on addressing this issue.



Figure 10. Additional Qualitative Comparisons. We show the synthesized head results of the *lip synchronization setting* on Testset A. (a) Ground truth; (b) AD-NeRF [23]; (c) SynObama [39]; (d) RAD-NeRF [40]; (e) **ER-NeRF (ours)**; (f) LSP [29]; (g) SSP-NeRF [28]; (h) Wav2Lip [32]; (i) PC-AVS [57].