# AvatarStudio: Text-driven Editing of 3D Dynamic Human Head Avatars

MOHIT MENDIRATTA, Max Planck Institute for Informatics and Saarland University, Germany

XINGANG PAN*, Max Planck Institute for Informatics, SIC, Germany

MOHAMED ELGHARIB*, Max Planck Institute for Informatics, SIC, Germany

KARTIK TEOTIA, Max Planck Institute for Informatics and Saarland University, Germany

MALLIKARJUN B R, Max Planck Institute for Informatics and Saarland University, Germany

AYUSH TEWARI, MIT CSAIL, United States of America

VLADISLAV GOLYANIK, Max Planck Institute for Informatics, SIC, Germany

ADAM KORTYLEWSKI, University of Freiburg and Max Planck Institute for Informatics, SIC, Germany

CHRISTIAN THEOBALT, Max Planck Institute for Informatics, SIC, Germany

Fig. 1. Our method AvatarStudio takes as input a 3D NeRF volume of a dynamic head (top) and produces visual edits that correspond to a target text prompt (second and third rows). Our method is the first designed specifically to handle text-based editing of videos. It also produces 3D-consistent results that can be viewed from an arbitrary camera viewpoint.

*Indicates equal contribution

Authors' addresses: Mohit Mendiratta, Max Planck Institute for Informatics and Saarland University, Germany, mmendira@mpi-inf.mpg.de; Xingang Pan, Max Planck Institute for Informatics, SIC, Germany, xpan@mpi-inf.mpg.de; Mohamed Elgharib, Max Planck Institute for Informatics, SIC, Germany, elgharib@mpi-inf.mpg.de; Kartik Teotia, Max Planck Institute for Informatics and Saarland University, Germany, ktoetia@mpi-inf.mpg.de; Mallikarjun B R, Max Planck Institute for Informatics and Saarland University, Germany, mbr@mpi-inf.mpg.de; Ayush Tewari, MIT CSAIL, United States of America, ayusht@mit.edu; Vladislav Golyanik, Max Planck Institute for Informatics, SIC, Germany, golyanik@mpi-inf.mpg.de; Adam Kortylewski, University of Freiburg and Max Planck Institute for Informatics, SIC, Germany, akortyle@mpi-inf.mpg.de; Christian Theobalt, Max Planck Institute for Informatics, SIC, Germany, theobalt@mpi-inf.mpg.de.

Capturing and editing full head performances enables the creation of virtual characters with various applications such as extended reality and media production. The past few years witnessed a steep rise in the photorealism of human head avatars. Such avatars can be controlled through different input data modalities, including RGB, audio, depth, IMUs and others. While these data modalities provide effective means of control, they mostly focus on editing the head movements such as the facial expressions, head pose and/or camera viewpoint. In this paper, we propose AvatarStudio, a text-based method for editing the appearance of a dynamic full head avatar. Our approach builds on existing work to capture dynamic performances of human heads using neural radiance field (NeRF) and edits this representation with a text-to-image diffusion model. Specifically, we introduce an optimization

strategy for incorporating multiple keyframes representing different camera viewpoints and time stamps of a video performance into a single diffusion model. Using this personalized diffusion model, we edit the dynamic NeRF by introducing view-and-time-aware Score Distillation Sampling (VT-SDS) following a model-based guidance approach. Our method edits the full head in a canonical space, and then propagates these edits to remaining time steps via a pretrained deformation network. We evaluate our method visually and numerically via a user study, and results show that our method outperforms existing approaches.Our experiments validate the design choices of our method and highlight that our edits are genuine, personalized, as well as 3D- and time-consistent.

CCS Concepts: • **Computing methodologies** → **Computer vision**; *Image manipulation.*

Additional Key Words and Phrases: Text-driven editing, neural rendering, 3D dynamic human head avatar, diffusion model

## 1 INTRODUCTION

The human face is at the center of our visual communications and hence its digitization is of utmost importance. The past few years have witnessed a sharp rise in the photorealism of digital faces. To achieve this, several methods were proposed, such as generative adversarial networks for 2D images [Karras et al. 2017, 2018]. Other methods build such high-quality in 3D using either explicit [Gecer et al. 2021; Lombardi et al. 2018], or more recently, learnable implicit [Lombardi et al. 2021; Zheng et al. 2022] scene representations. In addition to photorealism, controlling and rigging digital faces received a lot of attention. This includes for instance methods [Athar et al. 2022; Gao et al. 2022; Grassal et al. 2021] that utilize a low-dimensional parametric representation in form of 3D morphable model [Blanz and Vetter 1999; Egger et al. 2020] or some other latent space [Abdal et al. 2019; Teotia et al. 2023; Wang et al. 2021a]. Moreover, several data modalities have been explored as a control signal, such as RGB images [Bansal et al. 2018; Lombardi et al. 2021; Siarohin et al. 2019], audio [Suwajanakorn et al. 2017; Thies et al. 2020], sparse image representations such as contours and keypoints [Mihajlovic et al. 2022; Zakharov et al. 2019] and even input from sensors such as IMUs and IR cameras [Li et al. 2015; Wei et al. 2019].

The vast majority of existing methods for controlling digital faces, however, focus on editing the motion of the face. That is, controlling the facial expressions, head pose and/or the camera viewpoint [Kirschstein et al. 2023; Raj et al. 2021]. Controlling the facial appearance has been mostly studied in the context of facial relighting [Ranjan et al. 2023; Rao et al. 2022; Tan et al. 2022]. Here, methods were developed that edit the facial appearance as a function of the scene illumination. For this, the target illumination is commonly described via HDRI maps [Mallikarjun et al. 2021; Sun et al. 2019] or via a low dimensional representation such as spherical harmonics [Ranjan et al. 2023; Tewari et al. 2020; Zhou et al. 2019]. There are also methods for editing the facial appearance in a non-photorealistic manner [Fišer et al. 2017; Selim et al. 2016; Yang et al. 2022]. These methods usually take a target painting as input and can handle moving heads. One input modality that has not been fully explored yet for facial edits is text. Text is one of the most user-friendly data modalities that can be easily defined without any expert knowledge.

In the past few years, text-driven image synthesis attracted the attention of the research community. Thanks to the wide development and adaptation of transformers [Radford and Narasimhan 2018; Vaswani et al. 2017] and diffusion models [Rombach et al. 2022], several works have shown the ability of editing images in 2D [Brooks et al. 2023; Ruiz et al. 2022] and 3D [Haque et al. 2023; Jain et al. 2022; Poole et al. 2022], given text prompt as input. While 2D-based methods produce interesting results, they can not produce edits that are 3D-consistent. In contrast, 3D-based methods [Haque et al. 2023; Jain et al. 2022; Poole et al. 2022] show results on a 3D volume that can be rendered faithfully from an arbitrary camera viewpoint. However, even such methods lack in several ways. First, many of them [Aneja et al. 2022; Jain et al. 2022; Wang et al. 2021b, 2022] optimize their solution in a joint image-text embedding known as CLIP [Radford and Narasimhan 2018]. While such CLIP-based objective function leads to interesting results, it usually tends to generate limited edits. Second, none of the existing methods are designed to handle dynamic scenes and hence cannot process image sequences properly. This usually leads to clear artifacts and limited edibility (as shown in our experiments).

In this paper, we propose AvatarStudio, a text-based method for editing the appearance of a dynamic full head avatar.

Our approach assumes a digital head avatar as input, that can be trained from a multi-view performance capture of a human head. In particular, we follow the approach presented in HQ3DAvatar[Teotia et al. 2023] to learn a volumetric head avatar as being one of the latest in literature. Here, the head is represented as a canonical neural radiance field (NeRF) [Mildenhall et al. 2021] and a deformation network propagates the canonical representation across time. Our approach enables the text-based editing of such dynamic volumetric avatars in a view- and time-coherent manner.

Specifically, we perform editing through text-based conditional image generation with a diffusion model. We make several technical contributions to ensure that the editing of AvatarStudio are genuine, personalized, as well as 3D- and time-consistent. First, we sample several keyframes from the multi-view video that represent different camera viewpoints and time stamps of the performance capture. We introduce an optimization strategy to incorporate these keyframes into a single diffusion model, by fine-tuning a pre-trained model with a unique text identifier as proposed in [Ruiz et al. 2022]. Importantly, to prevent the leakage of information between keyframes we keep the sampled noise constant for every batch during each fine-tuning iteration. Based on this personalized diffusion model we can generate and edit each keyframe individually. We leverage this property to edit the dynamic NeRF by introducing a novel view- and time-aware Score Distillation Sampling (VT-SDS) approach, that iteratively edits the dynamic NeRF by sampling a random set of keyframes across the view and time domain. VT-SDS follows a model-based classifier-free guidance approach [Zhang et al. 2022], where we take advantage of the step-by-step generation process in diffusion models to guide the early stages of the image generation towards the content of the respective keyframe, while performing the editing throughout the later stages of the generation process with a large-scale pre-tained diffusion model. To ensure that the edited dynamic neural radiance field remains faithful and free from overfitting artifacts we use an annealing strategy that gradually lowers the

effect of the personalized diffusion model to enable high-frequency edits. Aspects of novelty of this work include:

- We present the first method for text-driven editing of dynamic 3D human head avatars. Our approaches leverages the state-of-the-art in neural volumetric scene representations together with recent advances in text-driven diffusion models to achieve high-quality editing of dynamic digital avatars.
- A new optimization strategy for incorporating multiple keyframes that represent different camera viewpoints and different time stamps, into a single diffusion model.
- A view- and time-aware Score Distillation Sampling (VT-SDS) that enables high-quality personalized editing in a coherent manner across the view and time domain.

We evaluate our method subjectively and numerically through a user study and compare against related methods. Results show that our approach produces a wide variety of text-based edits, while maintaining the integrity of the input identity (see Fig. 1). It generates temporally coherent results and clearly outperforms related methods.

## 2 RELATED WORK

This section provides an overview of generative models for image synthesis, with emphasize on diffusion models. First, we introduce Generative Adverserial Networks (GANs) and outline some of the landmark works in the literature [Karras et al. 2017, 2018]. We then discuss diffusion models in details, and focus on methods that can edit images using text as input. Here, diffusion models are divided into two main categories; 2D [Brooks et al. 2023; Ruiz et al. 2022] and 3D [Haque et al. 2023; Jain et al. 2022; Poole et al. 2022; Wang et al. 2021b, 2022]. One main difference between both approaches is that 3D methods can produce edits that are multi-view consistent, while 2D methods do not focus on changing the camera viewpoint. We outline important milestones in diffusion models for image synthesis. This includes enabling object-specific edits as in DreamBooth [Ruiz et al. 2022], and the introduction of the probability density distillation in DreamFusion [Poole et al. 2022]. We also discuss the CLIP image-text embedding [Radford et al. 2021], and how it is commonly used in the literature in formulating the objective function [Aneja et al. 2022; Jain et al. 2022; Wang et al. 2021b, 2022]. Our work differs from related works in several ways. To start with, it is the first that enables text-driven editing of image sequences. This is done by the introduction of a novel optimization that allows incorporating temporal frames in a pre-trained diffusion model. Second, our method is 3D by design, thus enables edits that are multi-view consistent. Last by not least, we do not use the CLIP embedding and thus enable stronger and more faithful edits. We now discuss related methods in more details.

### 2.1 Generative Models for Image Synthesis

The use of Generative Adverserial Networks, or GANs, for image synthesis have been an active research topic for the past years. This goes back to the revolutionary work of Goodfellow *et al.* [Goodfellow et al. 2014], where it was shown that a generator and a discriminator can be trained in an adversarial manner until the discriminator is no longer capable of telling whether the generator's output is real or synthesized. Results showed the ability of GANs to synthesis low resolution images of faces and other objects. This sparked a plethora of follow up work in generative models for image synthesis, including two notable works, Progressive GANs [Karras et al. 2017] and StyleGAN [Karras et al. 2018]. Here, GANs were able to synthesize high resolution images (1K) with extreme photorealism.

It was not late until diffusion models found their way into image synthesis. To this end, Ho *et al.* [Ho et al. 2020] have shown that diffusion probabilistic models (DPM) can be represented as a Markov Chain process using autoencoders. Briefly after in Dhariwal *et al.* [Dhariwal and Nichol 2021], it was shown that diffusion models can beat GANs in image synthesis quality. However, one main concern still remained; the computational complexity of such models. Rombach *et al.* [Rombach et al. 2022] addressed this concern by training diffusion models on latent spaces of autoencoders. Furthermore, several means of conditioning the diffusion model were shown, including text. This sparked a greater interest in text-driven image synthesis, especially with the rising popularity of transformers [Radford and Narasimhan 2018; Vaswani et al. 2017]. While earlier works of text-driven synthes such as DALL-E [Ramesh et al. 2021] relied primarily on transformers for language modeling and image synthesis, the follow-up version DALL-E 2 [Ramesh et al. 2022] utilized diffusion models. Other text-driven synthesis methods were also proposed using other generative models such as Style-GAN [Ramesh et al. 2021]. However, with the availability of public implementations of diffusion models such as Stable Diffusion [CompVis 2022; Rombach et al. 2022], text-driven image synthesis have witnessed great progress in the past couple of years.

### 2.2 Text-driven Diffusion Models for Image Synthesis

Current methods for text-driven image synthesis using diffusion models can be classified into 2D [Brooks et al. 2023; Ruiz et al. 2022] and 3D [Aneja et al. 2022; Haque et al. 2023; Jain et al. 2022; Lin et al. 2022; Poole et al. 2022; Wang et al. 2021b, 2022] approaches. While the former produces a wide range of visual edits in terms of content and style, the latter produces results that are 3D-consistent and thus can be viewed from an arbitrary camera angle. DreamBooth [Ruiz et al. 2022] is a 2D-based approach that handles the problem of fine-tuning large text-to-image diffusion models to a specific examined object. Here, multiple images (typically 3-5) of the same object is provided as input, while DreamBooth learns to associate a unique identifier to this object. This embeds the examined object in the output domain of the text-to-image diffusion model, thus allowing a wide variety of text-driven edits. Instruct-Pix2Pix [Brooks et al. 2023] takes a different approach for the same problem of text-driven image synthesis. Their idea is to generate paired synthetic data by utilizing the large language model of GPT-3 [Brown et al. 2020] together with the text-to-image model of Stable Diffusion. This strategy generalizes well to real user-instructions and real input images during test. Unlike our method, none of these 2D-based methods can generate results that are multi-view consistent [Brooks et al. 2023; Ruiz et al. 2022].

3D text-driven image synthesis methods can be classified as ones that utilize a CLIP embedding [Aneja et al. 2022; Jain et al. 2022;

Wang et al. 2021b, 2022] and others that use other means for optimizing their solution [Haque et al. 2023; Lin et al. 2022; Poole et al. 2022]. CLIP [Radford et al. 2021], short for "Contrastive Language-Image Pre-Training", is a joint text and image embedding trained in a way to predict the correct (text,image) pairing. Using CLIP and a Neural Radiance Field (NeRF) [Mildenhall et al. 2021] formulation, Dream Fields [Jain et al. 2022] extended 2D text-to-diffusion models to 3D. Along similar lines, DreamFusion [Poole et al. 2022] also extended 2D models to 3D, however not using CLIP embedding. Instead, a new loss is proposed based on probability density distillation. Here, a solution is initialized with a random 3D NeRF model, which is then optimized in a way so that 2D renderings from an arbitrary viewpoint minimizes the loss. Magic3D [Lin et al. 2022] improves upon the computational efficiency of DreamFusion using a coarse-to-fine manner. Following the 2D text-to-image model of Instruct-Pix2Pix [Brooks et al. 2023], Instruct-NeRF2NeRF [Haque et al. 2023] proposes a 3D text-to-image solution. Here, 2D images are iteratively edited using Instruct-Pix2Pix and the resulting 3D NeRF model is iteratively optimized. Results show 3D-consistent edits of various forms without the need of any additional training data.

ClipFace [Aneja et al. 2022] is a self-supervised approach for text-driven editing of human faces. Here, the face is modelled through a 3D morphable model (3DMM), where each facial component is trained separately. The solution uses a combination of a CLIP-based loss together with an adversarial training. While CLIP based methods achieve good results, they still can be limited in terms of their edibility. In addition, using a 3DMM lacks the ability of editing the full head [Aneja et al. 2022], and thus discard important regions such as the mouth interior. Finally, none of the remaining 3D-based methods are designed to handle image sequences and thus generate clear artifacts with limited editing, in contrast to our approach. Last but not least, we would like to point out that many of these methods are concurrent work [Aneja et al. 2022; Haque et al. 2023]. Despite that, we still compare against the most related methods to us [Haque et al. 2023; Jain et al. 2022; Poole et al. 2022].

## 3 METHOD

Our objective is to edit dynamic 3D full human heads using a text prompt expressing the desired edit (see Fig. 2). We assume the human heads are represented by dynamic NeRF-based models reconstructed with an existing method. In this work, we use HQ3DAvatar [Teotia et al. 2023] to obtain the dynamic human heads due to its high quality. The editing is achieved by leveraging the prior knowledge of a large text-guided latent diffusion model (LDM) [CompVis 2022; Rombach et al. 2022]. We first describe preliminaries on HQ3DAvatar and LDM (Sec. 3.1). We then introduce a new optimization strategy that adapts the LDM to capture the identity of the given head from different viewpoints and time stamps (Sec. 2.2). This optimization step is essential for preserving the original head identity and details during editing, as we will show in experiments. Next, we discuss how we use the adapted LDM to edit the dynamic head with a text prompt (Sec. 3.3). Our approach enables personalized and targeted edits from textual inputs while maintaining the examined identity. It also allows a wide range of text-driven edits of dynamic full heads

and produces temporally consistent results. For example, given a dynamic head and an exemplary textual input "Turn her into a zombie" or "Make him look like Van Gogh", our model can produce multi-view consistent NeRF edits that transfer the zombie or van Gogh styles to the target identity. We evaluate our approach visually and numerically through a user study, and results show that we clearly outperform existing methods.

### 3.1 Preliminaries

*3.1.1 Neural Radiance Fields.* The Neural Radiance Fields (NeRF) algorithm, as described in the paper by Mildenhall *et al.* [Mildenhall et al. 2021], employs a fully-connected deep neural network to represent a scene. The network is fed a continuous 5D coordinate that comprises a spatial location $\mathbf{p}$ with coordinates $(p_x, p_y, p_z)$ and a viewing direction $\mathbf{v}$ $(v_\theta, v_\phi)$. The output of the network is the volume density $\sigma$ and the view-dependent emitted radiance at that location. By utilizing classic volume rendering techniques and querying 5D coordinates along camera rays, NeRF can project output colors and densities onto novel views to synthesize images. The volume rendering process is differentiable, which allows for optimization of the representation using a set of input images with known camera parameters (extrinsic and intrinsic). These camera parameters are used to extract a per-pixel world-space ray parameterization that describes the 3D center $\mathbf{o}$ and direction $\mathbf{d}$ of the camera ray $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$ corresponding to each pixel in each image. The expected color $C(\mathbf{r})$ of camera ray $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$ with near and far bounds $t_n$ and $t_f$ can be calculated as follows:

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

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

*3.1.2 Scene Representation.* To capture complete head performances, we make use of the dynamic volumetric representation of HQ3DAvatar [Teotia et al. 2023], which learns a volumetric representation of the human head using multi-view RGB videos. We use this representation due to its high quality video results. While our method in principle could work with other dynamic NeRF models, examining this is outside the scope of this work.

HQ3DAvatar consists of two main stages. The first contains a deformation network $D$ that maps the input coordinates $\mathbf{p}$ of the world space to deformed positions in the canonical space as follows

$$\mathbf{p}_c = D(\mathbf{p}, \mathbf{e}) + \mathbf{p} \quad (3)$$

Here, $\mathbf{p}_c$ denotes the deformed coordinates in the canonical space, while $\mathbf{e}$ is a time embedding of the input RGB frames. This embedding is extracted using a pre-trained VGG-Face encoder [Parkhi et al. 2015]. The second stage contains an appearance network $A$ that predicts the radiance $\mathbf{c}$ and volume density $\sigma$ for each deformed coordinate. This is written as:

$$A : (\mathbf{p}_c, \mathbf{v}, \mathbf{e}) \rightarrow (\mathbf{c}, \sigma), \quad (4)$$

where $\mathbf{v}$ is the viewing direction. To make the method computationally efficient, a multi-resolution hash-grid based representation is used along the lines of Mueller et al. [Müller et al. 2022]. To this end, the appearance network $A$ consists of two main parts. The first
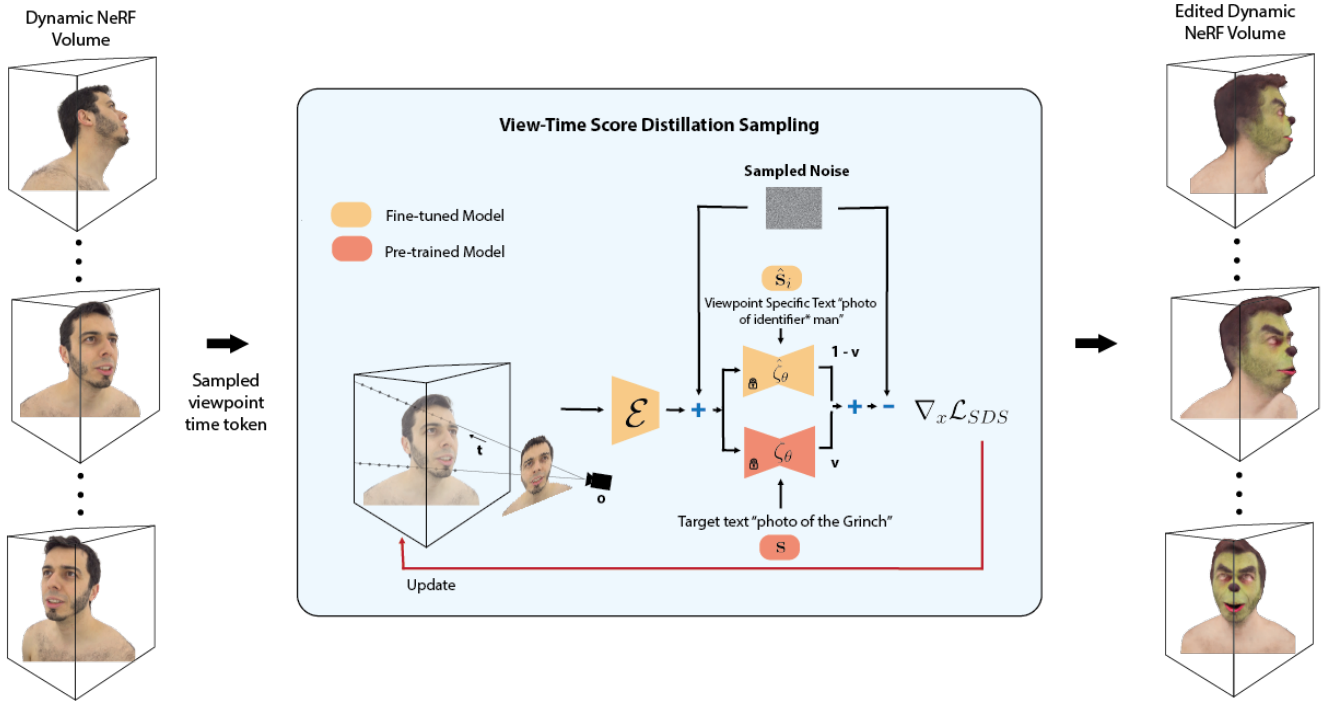
Fig. 2. **An overview of our method of our approach for text-driven editing of dynamic head avatars.** Our method takes as input a reconstructed dynamic NeRF volume (left) and a text prompt **S**, and produces corresponding visual edits (right). These edits can be viewed from an arbitrary camera viewpoint in a 3D-consistent manner. To this end, we propose a novel optimization that fine-tunes pre-trained latent diffusion models on multiple keyframes representing different viewpoints and time stamps (see Sec. 3.2). Furthermore, we employ a new view- and time-aware Score Distillation Sampling (VT-SDS) that combines a pre-trained latent diffusion model with our fine-tuned model (see Sec. 3.3 and Eq. 10).
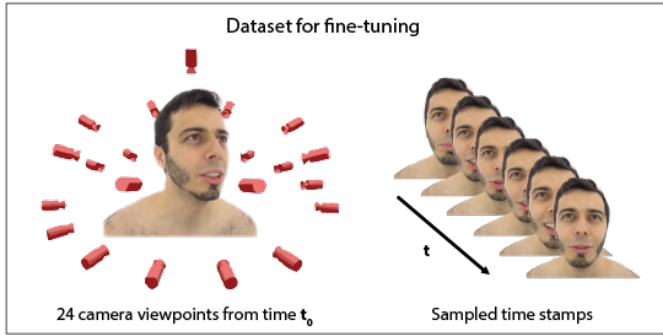


Fig. 3. Sample viewpoints and time stamps used in our diffusion model fine-tuning (see Sec. 3.2.1)

.

is the multi-resolution hash grid, while the second is an MLP-based network that outputs the radiance **c** and volume density $\sigma$. This MLP is conditioned on the time embedding **e** as well as the viewing direction **d**. Once we have the radiance field representation of the scene, we use standard volumetric integration to synthesize color **C** for each ray $r(t)$ using Eqn. (1). For more details on the method including the network design and the efficient rendering, please refer to the manuscript of [Teotia et al. 2023].

*3.1.3 Latent Diffusion Models.* Latent diffusion models (LDMs) [Rombach et al. 2022] are a class of Denoising Diffusion Probabilistic Models (DDPMs) [Ho et al. 2020] that use an vector-quantized auto-encoder [Van Den Oord et al. 2017] to translate an input image into a latent space in which a text-conditioned DDPM is trained. The encoder $\mathcal{E}$ processes a given image $I \in \mathbb{R}^{H \times W \times 3}$ to a latent representation **z**, such that $\mathbf{z} = \mathcal{E}(I)$. The decoder $\mathcal{D}$ reconstructs the estimated image $\tilde{I}$ from the latent, such that $\tilde{I} = \mathcal{D}(\mathbf{z})$ and $\tilde{I} \approx I$. The diffusion model is trained to generate images in the latent space of the encoder. Similar to other types of generative models [Mirza and Osindero 2014], diffusion models are in principle capable of modeling conditional distributions of the form $p(x|s)$, where $s$ is the conditioning variable. Conditional latent diffusion model are learned to optimize the following loss:

$$\mathbb{E}_{\mathcal{E}(I), s, \boldsymbol{\epsilon} \sim \mathcal{N}(0,1), t} \left[ \|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta \left( \mathbf{z}_t, t, s \right)\|_2^2 \right], \quad\quad (5)$$

where $t$ is the diffusion time step and $\mathbf{z}_t$ is the noisy latent code at time $t$. $\boldsymbol{\epsilon}$ is the noise sample, $\boldsymbol{\epsilon}_\theta$ is the denoising model with parameters $\theta$ and $s$ is the conditioning input. During training, $\boldsymbol{\epsilon}_\theta$ is optimized. At inference, a latent code is generated by randomly sampling a noise tensor and denoising it iteratively based on a conditioning input.

## 3.2 Fine-Tuning Text-to-Image Latent Diffusion Model

A key challenge in dynamic full head editing is to preserve the original characteristics of the head (*e.g.*, its identity, details, motions, *etc*) rather than creating a completely different head. Editing using the original LDM would soon lead to drifts of these characteristics due to information leakage. A potential way to alleviate this problem is to fine-tune the LDM on images of the given head using DreamBooth [Ruiz et al. 2022]. However, unlike DreamBooth which aims to sample new 2D images and thus only needs to capture the object identity, we want to edit a 3D head across different viewpoints and time stamps. Thus, ideally, the LDM should not only be identity-aware but also viewpoint-aware and time-aware. Our investigations revealed that implementing DreamBooth for multiple concepts (*i.e.*, multiple viewpoints and time stamps) is also prone to concept leakages, producing suboptimal editing. Hence, below we detail our new optimization strategy which is designed specifically for associating multiple concepts for different viewpoints and time stamps.

*3.2.1 Optimization.* In our work, we fine-tune the LDM of Stable Diffsion [CompVis 2022; Rombach et al. 2022]. To this end, we use images of the given dynamic 3D head from different viewpoints and time stamps (see Fig. 3), denoted as $\{\mathbf{x}_i; i \in \{1, ..., n\}\}$. How we select the images will be discussed in Sec. 3.2.2. We then assign a label $\mathbf{P}_i$ to each of these images, using the format 'photo of a [identifier] [class noun]'. The identifier is a one-of-a-kind code of 10 characters, unique for each image, while the class noun is 'man' or 'woman' depending on the gender of the given head. Each identifier is initialized via a random word generator. Our aim is to fine-tune a pretrained text-to-image LDM $\zeta_\theta$ so that, given an initial noise map $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and a conditioning vector $\mathbf{s}_i = \Gamma(\mathbf{P}_i)$ produced using a text encoder $\Gamma$, the fine-tuned LDM $\hat{\zeta}_\theta$ will reconstruct the image $\mathbf{x}_i$. The model is fine-tuned using a squared error loss to denoise a variably-noised image or latent code $\mathbf{z}_{t,i} := \alpha_t \mathcal{E}(\mathbf{x}_i) + \beta_t \epsilon$ as follows:

$$\mathbb{E}_{\mathbf{x}_i, \mathbf{s}_i, \epsilon, t} \left[ w_t \| \zeta_\theta(\alpha_t \mathbf{x}_i + \beta_t \epsilon, \mathbf{s}_i) - \mathcal{E}(\mathbf{x}_i) \|_2^2 \right] \tag{6}$$

where $\alpha_t, \beta_t, w_t$ are terms that control the noise schedule and sample quality. These terms are a function of the diffusion process time $t \sim \mathcal{U}([0, 1])$ [Ruiz et al. 2022]. To overcome potential language drift of language models [Lee et al. 2019; Lu et al. 2020], we incorporate Class-specific Prior Preservation Loss [Ruiz et al. 2022]. In practice, we use a batch size of 3 during fine-tuning, which corresponds to 3 randomly sampled $\mathbf{x}_i$ with different viewpoints or time stamps. We observe that using different noise $\epsilon$ within each batch may lead to concept leakage between different identifiers. Thus, we use a shared noise $\epsilon$ within each batch, which helps the identifiers to capture the variations in the image avoiding any leakage.

*3.2.2 Time Embedding Sampling.* Here we introduce how we select the images $\{\mathbf{x}_i\}$ used to fine-tune the LDM. As discussed above, $\{\mathbf{x}_i\}$ should include multiple viewpoints and time stamps. Thus, we use all camera views of the first frame as the multiview images. In our experiments we use a multiview camera rig equipped with 24 RGB cameras around the head (see Sec. 4.2). Hence, we use all these 24 camera viewpoints. We also include 6 other frames from the frontal camera view, which should be as diverse as possible.

To achieve this, we employ an empirical strategy. We generate embeddings $\mathbf{e}_j$ in HQ3DAvatar (see Eq. (3))

for each of the $m$ frames in a given dynamic head. We calculate the mean of these $m$ embeddings and then select 6 embeddings and their corresponding frames that exhibit the greatest variation from the mean based on their absolute difference. We make sure to discard similar neighbouring frames with similar deformations to avoid redundancy.

## 3.3 Text-guided Dynamic NeRF Editing

With a pre-trained HQ3DAvatar on a specific identity, we perform text-driven editing by optimizing the appearance network $A$ (see Eq. 4), while keeping the deformation network $D$ fixed. In other words, we edit the appearance in HQ3DAvatar's canonical space. Please refer to Fig. 2 for an overview of this editing process. Our text-driven editing is developed based on Score Distillation Sampling (SDS) loss [Poole et al. 2022], which supervises an image to follow the text prompt using a text-to-image LDM. Here, we render an image $\mathbf{x}$ at each optimization step by randomly sampling from the camera viewpoints and time stamps corresponding to $\{\mathbf{x}_i\}$ (Sec. 3.2.2).

At every step of the optimization process, a random diffusion time instant $t$ is sampled and noise is injected into the rendered image $\mathbf{x}$:

$$\mathbf{x}_t = \mathbf{x} + \epsilon_t, \tag{7}$$

where $\epsilon_t$ is the noise map generated via a noising function $Q(t)$.

With our modified Score Distillation Sampling loss (see Sec. 3.3.1), the gradients for score distillation are calculated on a per-pixel basis as follows:

$$\nabla_x \mathcal{L}_{SDS} = w(t) \left( \epsilon_t - \Psi(\mathbf{x}_t, t, \mathbf{s}, \mathbf{s}_i) \right). \tag{8}$$

Here, $w(t)$ is a weighting function following [Poole et al. 2022], $\mathbf{s}$ is the text embedding of the user-input text for editing the dynamic NeRF, and $\mathbf{s}_i$ is the same as in Eq. 6. Furthermore, $\Psi(\mathbf{x}_t, t, \mathbf{s}, \mathbf{s}_i)$ is the noise predicted by a combination of our fine-tuned and a pre-trained diffusion model given $\mathbf{x}_t$, $t$, $\mathbf{s}$ and $\mathbf{s}_i$ as we will introduce later (Eq. 10).

Moreover, we utilize a regularizer for the density field generated by A [Melas-Kyriazi et al. 2023] as follows:

$$\mathcal{L}_{\text{entropy}} = \omega \cdot \log_2(\omega) - (1 - \omega) \cdot \log_2(1 - \omega). \tag{9}$$

Here, $\omega$ is the cumulative sum of density weights computed along each ray in the scene. This regularizer is an entropy loss that promotes the points to be either completely transparent or completely opaque.

*3.3.1 Modified Score Distillation Sampling.* Inspired by recent work [Zhang et al. 2022] on reducing overfitting and severe language drift in fine-tuned text-driven diffusion models, we utilize our fine-tuned model to provide content features. These features are combined with scores from the pre-trained model in a manner similar to classifier-free guidance [Ho and Salimans 2022].

To this end, we use the notation $\hat{\zeta}_\theta$ to refer to the fine-tuned denoising model, and $\zeta_\theta$ to refer to the pre-trained text-to-image model. During Score Distillation Sampling (SDS), we guide the pre-trained model with our fine-tuned model by using a linear combination of the noise estimated from each model, for a specified range

of optimization steps. Thus, the noise estimation in a SDS step can be defined as:

$$\Psi\left(\mathbf{x}_t, t, \mathbf{s}, \mathbf{s}_i\right) = w\left(v\zeta_\theta\left(\mathbf{x}_t, \mathbf{s}\right) + (1-v)\hat{\zeta}_\theta\left(\mathbf{x}_t, \mathbf{s}_i\right)\right)$$
$$+ (1-w)\zeta_\theta\left(\mathbf{x}_t\right), \tag{10}$$

where $w$ is the overall guidance weight and $v$ stands for the model guidance weight, which depends on $t$ as we will discuss later. $\hat{\mathbf{s}}_i$ is the same as in Eq. 6 and $\mathbf{s}$ is the target language conditioning obtained from the target prompt (see Sec. (3.2.1)).

To ensure that the edited dynamic neural radiance field remains faithful and free from overfitting artifacts, we use Eq. (10) with $0.5 \leq v \leq 0.7$ for sampling when $t > K$, and $v = 1$ for sampling when $t \leq K$. Unless stated otherwise, we use $K = 600$. We follow [Lin et al. 2022] and use an annealed SDS loss function that gradually lowers the maximum time-step used to sample $t$. This enables SDS to emphasize high-frequency information once the edit's outline has been established. In the ablation study of Sec. 4.6, we show that using Eq. (10) for Score Distillation Sampling is essential for generating good edits in both the the spatial and temporal domains.

## 4  EXPERIMENTS

In this section, we evaluate the performance of our method subjectively and numerically.

**Performance measures.** We asses three main aspects of the generated results. First, their ability to respect the target text prompt. Here, it is important to maintain the integrity of the input identity. Second, we asses the ability of generating edits that are 3D-consistent, and thus can be rendered from an arbitrary camera viewpoint. Third, we assess the temporal coherency of the generated edits. Visual results are shown throughout figures and the supplemental video. Numerical results are extracted from a user study. Results show that our method is capable of producing a wide variety of text-driven edits that are 3D-consistent and temporally coherent.

**Baselines.** We compare against two text-driven image synthesis baselines: Dream Fields [Jain et al. 2022] and Instruct-NeRF2NeRF [Haque et al. 2023]. We also compare against an implementation that combines the diffusion model fine-tuning method of DreamBooth [Ruiz et al. 2022] together with the 3D text-based editing method of Dream-Fusion [Poole et al. 2022]. Results show that our method clearly outperforms the state of the art visually and numerically.

In the next two sections we discuss implementation details and the multi-view data that is used in our experiments. Subsequently, we discuss the user study that we perform in assessing our experiments (Sec. 4.3). We show in Sec. 4.4 extensive evaluation of our method using varying prompts. In Sec. 4.5, we compare against related methods subjectively and objectively. Finally, we investigate the various design choices of our method in an ablation study (Sec. 4.6). Here we investigate our two main contributions. First, we investigate the importance of our fine-tuning strategy (Sec. 3.2) which incorporates multiple camera viewpoints and different time stamps. We also investigate the importance of incorporating Eq. 10 in the Score Distillation Sampling as discussed in Sec. 3.3. Results show that all our design choices contribute positively to the final output.

### 4.1  Implementation Details

We employ a consistent set of parameters for all experiments, without optimizing them specifically for each scenario. We utilize the open-source Stable Diffusion model [Rombach et al. 2022] as our prior for the diffusion model. This model was trained on the LAION dataset [Schuhmann et al. 2022], which consists of pairs of text and images. We render our images at a resolution of 256px. Since the Stable Diffusion model is specifically designed for images with a resolution of 512px, we first upsample our renders to 512px before passing them to the Stable Diffusion's latent space encoder, which is a Variational Autoencoder (VAE). We use classifier-free guidance of strength 10. Additionally, we set the VT SDS ratio to v=0.6 and K=600. We optimize the head avatar for a single prompt using the Adam [Kingma and Ba 2014] optimizer with learning rate 1e-3 for 10000 iterations. The optimization process takes approximately 60 minutes on a single A100 GPU. For our fine-tuning step, we optimize the diffusion model for a total of 28000 steps using the Adam optimizer [Kingma and Ba 2014] with image size 512px, batch size 3 and a learning rate 5e-5.

### 4.2  Data Capture

Our method is trained on multi-view data. We use a 360-degree camera rig equipped with 24 Sony RXO II cameras that are hardware-synced and capable of recording 25 frames per second at a 4K resolution. These cameras are positioned in a way to capture the entire human head, including the scalp's hair. They are also accompanied by LED strips to provide uniform illumination. The cameras are calibrated using a static structure with distinctive features. The intrinsic and extrinsic parameters are estimated using Metashape [Agisoft 2020]. Background subtraction is carried out using the matting approach of Lin *et al.* [Lin et al. 2021] to eliminate static elements like wires, cameras, and other objects. To simplify the process of background subtraction, a diffused white sheet is placed inside the rig, which contains holes for each camera lens. Please refer to Fig. 4 for an overview of data captured and used in our experiments. We use data collected by this camera rig for all our experiments, including evaluating related methods. We show results on 5 identities performing a variety of expressions. Some identities are also reading a set of sentences known as Pangrams [1]. Here, each sentence contains all the 26 Latin letters. Our videos are recorded at 25 frames per second, and are between 300-500 frames long.

### 4.3  User Study

It is challenging to evaluate text-driven visual edits numerically due to the absence of ground-truth data. In fact, one text prompt could have several different possible visual edits. Current methods used two main strategies for numerical evaluations. One approach is to measure the alignment of the produced edits with the input text prompt using the CLIP space [Haque et al. 2023; Jain et al. 2022; Poole et al. 2022]. This method, however, does not evaluate the temporal coherency of the solution. In addition, it is expected to favour methods that optimize their solution in the CLIP space. Another strategy for numerical evaluation is to perform a user study, as adapted by [Lin et al. 2022; Wang et al. 2022]. We believe

---

[1]https://callibeth.com/downloads/pangrams111.pdf

Fig. 4. We trained our method using data captured from a multi-view video camera rig. The rig contains 24 video cameras positioned around the head. The figure shows an example capture from each camera viewpoint

this strategy is more suitable as it is not tied to a specific text-image embedding, and due to the subjective nature of the examined problem.

Motivated by this, we designed a user study that assess several important aspects of text-driven video edits. For a given identity and a given text prompt, our user study shows four videos side-by-side. The first is the original input as produced by HQ3DAvatar [Teotia et al. 2023], while the remaining videos are the output of three different text-driven editing methods. The order of these three videos is randomly shuffled. Each video is around 19 seconds, featuring either a talking person or different facial expressions captured by a rotating camera. Participants were asked to watch the video and were given the option to replay it as desired. They were then asked to answer the following set of questions.

- Q1: Which method better retains the identity of the input sequence (identity preservation)?
- Q2: Which method better follows the given input textual prompt (prompt preservation)?
- Q3: Which method better maintains temporal consistency (temporal consistency)?
- Q4: Which method is better overall considering the above 3 aspects (identity preservation, textual preservation, temporal consistency)?

As shown, the first three questions are designed to assess (from the top) the identity preservation, the prompt preservation and the temporal consistency of the output. The user gives an answer to each of these three questions, and hence a specific method output could be chosen as the best one in just a subset of these questions. Finally, the fourth question is a measure of the overall quality considering all these three aspects (identity preservation, textual preservation, temporal consistency).

### 4.4 Text-driven Full Head Editing

Fig. 5-8 show various edits generated by our method. Here, we show results for several identities with various input text prompts. Subjects are talking and performing various expressions while we show results from different camera viewpoints. Results show that our method can handle a wide variety of text-driven visual edits. This includes both photorealistic (e.g. Fig. 5 and Fig. 6 "Photo of an old person"), and non-photorealistic (Fig. 5 and Fig. 9 "Photo of a Panda") edits. We can edit specific regions of the face according to the input text prompt. For instance, Fig 6 (second column) and Fig. 8 (fourth column) show that we can edit the color of hair in isolation, which respects the text prompt of "Photo of a face with blue hair". Our method can also edit the geometry so that it is inline with the text prompt. For instance, see "Photo of the Grinch" (Fig. 6 and Fig. 9) and "Photo of panda" (Fig. 9). Our method handles a variety of facial expressions and head movements. This includes normal speech (Fig. 6), smiling (Fig. 7, middle row), and extreme expressions (Fig. 7, last row). Results also show that our method produces edits that are 3D consistent as well as temporally coherent. This is best observed in the supplemental video. Last but not least, our method achieves this wide range of edits while maintaining the integrity of the input identity. Please refer to the supplemental video for more results.

### 4.5 Comparison Against Related Methods

We compare against two recent 3D text-based editing methods. That is InstructNeRF2NeRF [Haque et al. 2023] and Dream Fields [Jain et al. 2022]. In both methods, we edit only the reconstructed HQ3DAvatar's appearance in the canonical space, while keeping the deformation network fixed. We then use the deformation network to warp the edited appearance to the remaining time stamps. Since the original Dream Fields method [Jain et al. 2022] generates the entire image purely from text, we implemented a version which is initialized
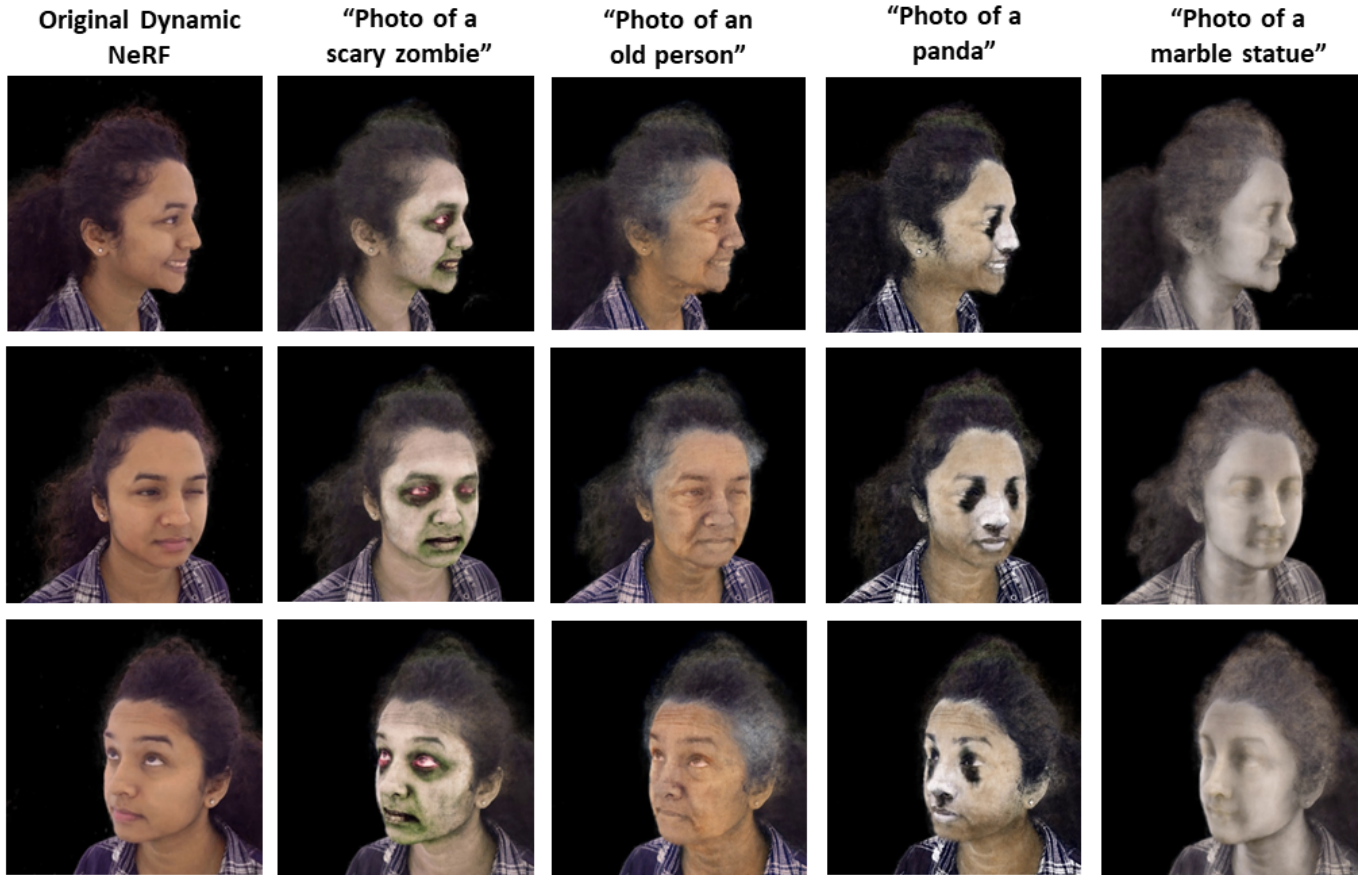
Fig. 5. Our method produces compelling text-driven visual edits for different text prompts. Note the good sharp edits in the eyes (see the second column) and how the fourth column show edits in the facial geometry. Our results are 3D- and temporally consistent as can be seen in the supplemental video.

by our reconstructed HQ3DAvatar model. This is done by optimizing the HQ3DAvatar's appearance via the CLIP loss [Radford et al. 2021]. We call this implementation Dream Fields++ in the rest of the paper. We also call InstructNeRF2NeRF in our figures InstructN2N for brevity.

Fig. 10-11 show results of different methods for different text prompts. Results show that Dream Fields++ generates significant artifacts that destroys the integrity of the input image. These edits are also clearly not inline with the target text prompt. Furthermore, Fig. 10 (third column) shows that Dream Fields++ generates edits in unwanted regions. Here, despite the text prompt says "Photo of a face with blue hair", Dream Fields++ turns the lips blue. Similarly, InstructNeRF2NeRF generates edits that are not well aligned with the target text prompts. For instance, see Fig. 10 and Fig. 11 with prompt "Photo of an old person". InstructNeRF2NeRF could also destroy the original identity (see Fig. 11, last column) and some times it generates edits in wrong regions (see Fig. 10, last column, lips). In contrast, our method generates edits that are more inline with the target text prompt. In addition, it maintains the integrity of the input image and produces clearly more temporally consistent

results. For this, please refer to the supplemental video and the user study (discussed next).

Following Sec. 4.3, we perform a user study to compare our method against Dream Fields++ and InstructNeRF2NeRF. Here, we examine a total of 3 identities, each processed with 3 different prompts. Thus, the users were presented with a total of 9 videos, each consisting of four videos playing side-by-side as discussed in Sec. 4.3. Each time, the users were asked to asses various spatial and temporal aspects of the different methods by answering the four questions listed in Sec. 4.3. Hence, in total they were asked to answer these questions 9 times, where in each time the order of different methods was randomly shuffled. Tab. 1 summarizes the findings of this user study. In total, 48 users participated in the study. Our technique is rated clearly the best in all questions. For instance, users rated our method the best in the overall quality 85.4% of the time. This compares favorably to 3.9% and 10.6% for Dream Fields++ and InstructeNeRF2NeRF respectively. One can create another straightforward baseline by combining the fine-tuning approach of DreamBooth [Ruiz et al. 2022] together with the 3D text-based editing method of DreamFusion [Poole et al. 2022]. We noticed, however, that this approach suffers from overfitting and

Fig. 6. Our method produces pleasing text-driven visual edits for different prompts. Note for instance the geometrical edits in the fourth column and how our method can handle interesting prompts such as "photo of an old person" as shown in the last column. Our method can also edit specific regions as instructed by the prompt (see blue hair, second column). Our results are 3D- and temporally consistent and maintains the original identity.

|  | Dream Fields++ | InstructNeRF2NeRF | StudioAvatar |
|---|---|---|---|
| Q1: Identity preservation | 4.6 | 7.2 | 88.2 |
| Q2: Prompt preservation | 4.4 | 22.2 | 74.4 |
| Q3: Temporal consistency | 6 | 8.3 | 85.6 |
| Q4: Overall | 3.9 | 10.6 | 85.4 |

Table 1. Reporting the results of our user study which included responses from 48 participants. Each participant was presented with the outputs of different methods and was asked to pick his/her preference with respect to the four questions listed in Sec. 4.3. The table reports the percentages at which a method was rated the best with respect to a specific question. Hence the sum of each row should add up to 100. Our method was rated the best in overall quality (Q4) 85.4% of the time. This compares favorably to 3.9% and 10.6% for Dream Fields++ and InstructeNeRF2NeRF respectively. Our method was also clearly rated the best in the remaining questions.

usually gives poor edits that do not follow the prompt well, as shown in Fig. 12. Here, it is quite clear that our method is clearly superior.

## 4.6 Ablation Study

We evaluate the various design choices of our method in an ablation study. First, we investigate the importance of our optimization strategy which accounts for multiple camera viewpoints and different time stamps during model fine-tuning (Sec. 3.2), To achieve this,

we perform two experiments. First, we replace our optimization strategy with DreamBooth's [Ruiz et al. 2022]. While we still embed multiple camera viewpoints and different time stamps, however, DreamBooth's fine-tuning strategy only estimates one token for all images. Fig. 13 shows the output of this process. It is clear that this strategy leads to significant artifacts. However, since our approach estimates multiple tokens for each of the embedded frames, we can generates significantly better results (see top row). We also

Fig. 7. Our method produces compelling text-driven visual edits for different text prompts. Note how the bronze bust editing prompt introduces the effects specific to metals and makes the facial texture more uniform. Note the sharp transition between the edited red eyes and the skin in the third column. All results preserve the original identity and are 3D- and temporally consistent, which can be observed in our supplemental video.

investigate the impact of not incorporating any time stamps in our fine-tuning. To achieve this, we used the same optimization discussed in Sec. 3.2, however, we incorporated just the different camera viewpoints at time 0. Fig. 14-15 shows that this strategy leads to clear temporal inconsistencies in the output. Please see the supplemental video. This is a critical problem, especially during video edits. Incorporating different time stamps, however, clearly produces better results with temporally coherent output.

The second part of our ablation study investigates the importance of our modified Score Distillation Sampling of Sec. 3.3. More specifically, we investigate the importance of using Eq. 10 during Score Distillation Sampling. To this end, we performed two main experiments. In the first experiment, we set v and K of Eq. 10 to 0. This examines the importance of the pre-trained model during Score Distillation Sampling. Here, Fig. 16 shows that by removing the pre-trained model, results undergo a significant drop in performance. In the second experiment, we attempted to use the original SDS loss of DreamFusion [Poole et al. 2022]. This approach, however, did not converge to any useful output. These experiments shows the importance of our modified Score Distillation Sampling. Last but

not lest, Fig. 17 shows that using the SDS annealing of Lin *et al.* [Lin et al. 2022] leads to better details in the final results.

## 5 DISCUSSION AND FUTURE WORK

We presented the first method for text-driven edits of dynamic head avatars. Our method utilizes a state-of-the-art NeRF based representation for dynamic heads and produces edits that are temporally coherent. Our results can be viewed from an arbitrary camera viewpoint in a 3D-consistent manner. At the heart of our method is a novel optimization strategy that incorporates multiple camera viewpoints, and multiple frames taken at different time stamps, in a pre-trained latent diffusion model. In addition, we proposed a new view-and-time-aware Score Distillation Sampling approach that combines knowledge from the pre-trained model, as well as our fine-tuned model. Our method enables a wide variety of text-driven edits and can produce both photorealistic and non-photorealistic edits. We compared against related methods and results show that our approach produces better edits that are more temporally stable and more inline with the text prompts. This is confirmed visually, as well as numerically via a user study.

Fig. 8. Our method produces compelling text-driven visual edits for different text prompts. Our results maintains the original identity and are 3D- and temporally consistent as shown in the supplemental video.

Our work pushes the boundaries of text-driven visual edits. Nevertheless, several interesting avenues are still open for future work. While our method can produce a wide variety of edits, it requires multi-view data captured in a uniform illumination. Thus a very interesting research direction will be to handle just a monocular video as input, shot with in-the-wild conditions. This could be followed up, for instance, with a method that takes just a single image as input and rigs the output according to a target motion. While some of our results show geometrical edits, future work could look into producing edits that change the head geometry more drastically. Currently our method is computationally expensive, requiring around 60 minutes train using a single A100 GPU. Future work could look into reducing this computational cost. We hope that our work in text-driven visual editing encourages further research in this interesting problem.

## REFERENCES

Rameen Abdal, Yipeng Qin, and Peter Wonka. 2019. Image2StyleGAN: How to Embed Images Into the StyleGAN Latent Space?. In *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*. 4431–4440. https://doi.org/10.1109/ICCV.2019.00453

LLC Agisoft. 2020. Metashape python reference. *Release* 1, 0 (2020), 1–199.

Shivangi Aneja, Justus Thies, Angela Dai, and Matthias Nießner. 2022. ClipFace: Text-guided Editing of Textured 3D Morphable Models. In *ArXiv preprint arXiv:2212.01406.*

ShahRukh Athar, Zexiang Xu, Kalyan Sunkavalli, Eli Shechtman, and Zhixin Shu. 2022. RigNeRF: Fully Controllable Neural 3D Portraits. In *Computer Vision and Pattern Recognition (CVPR).*

Aayush Bansal, Shugao Ma, Deva Ramanan, and Yaser Sheikh. 2018. Recycle-GAN: Unsupervised Video Retargeting. In *ECCV.*

Volker Blanz and Thomas Vetter. 1999. A morphable model for the synthesis of 3D faces. In *Proceedings of the 26th annual conference on Computer graphics and interactive techniques*. 187–194.

Tim Brooks, Aleksander Holynski, and Alexei A. Efros. 2023. InstructPix2Pix: Learning to Follow Image Editing Instructions. In *CVPR.*

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language Models are Few-Shot Learners. In *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (Eds.), Vol. 33. Curran Associates, Inc., 1877–1901. https://proceedings.neurips.cc/paper_files/paper/2020/file/1457c0d6bfcb4967418bfb8ac142f64a-Paper.pdf

CompVis. 2022. *Stable Diffusion.* https://github.com/CompVis/stable-diffusion

Prafulla Dhariwal and Alexander Nichol. 2021. Diffusion Models Beat GANs on Image Synthesis. In *Advances in Neural Information Processing Systems*, M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan (Eds.), Vol. 34. Curran Associates, Inc., 8780–8794. https://proceedings.neurips.cc/paper_files/paper/2021/file/49a42d3d1ec9fa4bd8d77d02681df5cfa-Paper.pdf

Bernhard Egger, William AP Smith, Ayush Tewari, Stefanie Wuhrer, Michael Zollhoefer, Thabo Beeler, Florian Bernard, Timo Bolkart, Adam Kortylewski, Sami Romdhani, et al. 2020. 3d morphable face models—past, present, and future. *ACM Transactions*

Fig. 9. Our method produces compelling text-driven visual edits for different text prompts. Note how our method can change the appearance (all columns) and the geometry (last two columns). Our results are 3D- and temporally consistent.

on Graphics (TOG) 39, 5 (2020), 1–38.

Jakub Fišer, Ondřej Jamriška, David Simons, Eli Shechtman, Jingwan Lu, Paul Asente, Michal Lukáč, and Daniel Sýkora. 2017. Example-Based Synthesis of Stylized Facial Animations. *ACM Transactions on Graphics* 36, 4, Article 155 (2017).

Xuan Gao, Chenglai Zhong, Jun Xiang, Yang Hong, Yudong Guo, and Juyong Zhang. 2022. Reconstructing Personalized Semantic Facial NeRF Models From Monocular Video. *ACM Transactions on Graphics (Proceedings of SIGGRAPH Asia)* 41, 6 (2022). https://doi.org/10.1145/3550454.3555501

Baris Gecer, Stylianos Ploumpis, Irene Kotsia, and Stefanos P Zafeiriou. 2021. Fast-GANFIT: Generative Adversarial Network for High Fidelity 3D Face Reconstruction. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2021).

Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2014. Generative Adversarial Nets. In *Advances in Neural Information Processing Systems*, Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence, and K.Q. Weinberger (Eds.), Vol. 27. Curran Associates, Inc. https://proceedings.neurips.cc/paper_files/paper/2014/file/5ca3e9b122f61f8f06494c97b1afccf3-Paper.pdf

Philip-William Grassal, Malte Prinzler, Titus Leistner, Carsten Rother, Matthias Nießner, and Justus Thies. 2021. Neural Head Avatars from Monocular RGB Videos. *arXiv preprint arXiv:2112.01554* (2021).

Ayaan Haque, Matthew Tancik, Alexei Efros, Aleksander Holynski, and Angjoo Kanazawa. 2023. Instruct-NeRF2NeRF: Editing 3D Scenes with Instructions. (2023).

Jonathan Ho, Ajay Jain, and Pieter Abbeel. 2020. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems* 33 (2020), 6840–6851.

Jonathan Ho and Tim Salimans. 2022. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598* (2022).

Ajay Jain, Ben Mildenhall, Jonathan T. Barron, Pieter Abbeel, and Ben Poole. 2022. Zero-Shot Text-Guided Object Generation with Dream Fields. (2022).

Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. 2017. Progressive Growing of GANs for Improved Quality, Stability, and Variation. *CoRR* abs/1710.10196 (2017). arXiv:1710.10196 http://arxiv.org/abs/1710.10196

Tero Karras, Samuli Laine, and Timo Aila. 2018. A Style-Based Generator Architecture for Generative Adversarial Networks. *CoRR* abs/1812.04948 (2018). arXiv:1812.04948 http://arxiv.org/abs/1812.04948

Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).

Tobias Kirschstein, Shenhan Qian, Simon Giebenhain, Tim Walter, and Matthias Nießner. 2023. NeRSemble: Multi-view Radiance Field Reconstruction of Human Heads. https://doi.org/10.48550/arXiv.2305.03027 arXiv:2305.03027 [cs.CV]

Jason Lee, Kyunghyun Cho, and Douwe Kiela. 2019. Countering language drift via visual grounding. *arXiv preprint arXiv:1909.04499* (2019).

Hao Li, Laura Trutoiu, Kyle Olszewski, Lingyu Wei, Tristan Trutna, Pei-Lun Hsieh, Aaron Nicholls, and Chongyang Ma. 2015. Facial Performance Sensing Head-Mounted Display. *ACM Transactions on Graphics (Proceedings SIGGRAPH 2015)* 34, 4 (July 2015).

Chen-Hsuan Lin, Jun Gao, Luming Tang, Towaki Takikawa, Xiaohui Zeng, Xun Huang, Karsten Kreis, Sanja Fidler, Ming-Yu Liu, and Tsung-Yi Lin. 2022. Magic3D: High-Resolution Text-to-3D Content Creation. *arXiv preprint arXiv:2211.10440* (2022).

Shanchuan Lin, Andrey Ryabtsev, Soumyadip Sengupta, Brian L Curless, Steven M Seitz, and Ira Kemelmacher-Shlizerman. 2021. Real-time high-resolution background matting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 8762–8771.

Stephen Lombardi, Jason Saragih, Tomas Simon, and Yaser Sheikh. 2018. Deep Appearance Models for Face Rendering. *ACM Trans. Graph.* 37, 4, Article 68 (July 2018), 13 pages.

Stephen Lombardi, Tomas Simon, Gabriel Schwartz, Michael Zollhoefer, Yaser Sheikh, and Jason Saragih. 2021. Mixture of Volumetric Primitives for Efficient Neural

Fig. 10. Our approach outperforms Dream Fields++ [Jain et al. 2022] and InstructeNeRF2NeRF [Haque et al. 2023] spatially and temporally. It maintains the original identity and generates edits that are inline with the target prompts. Please see the supplemental video and the user study (Tab. 1) for video coherency.
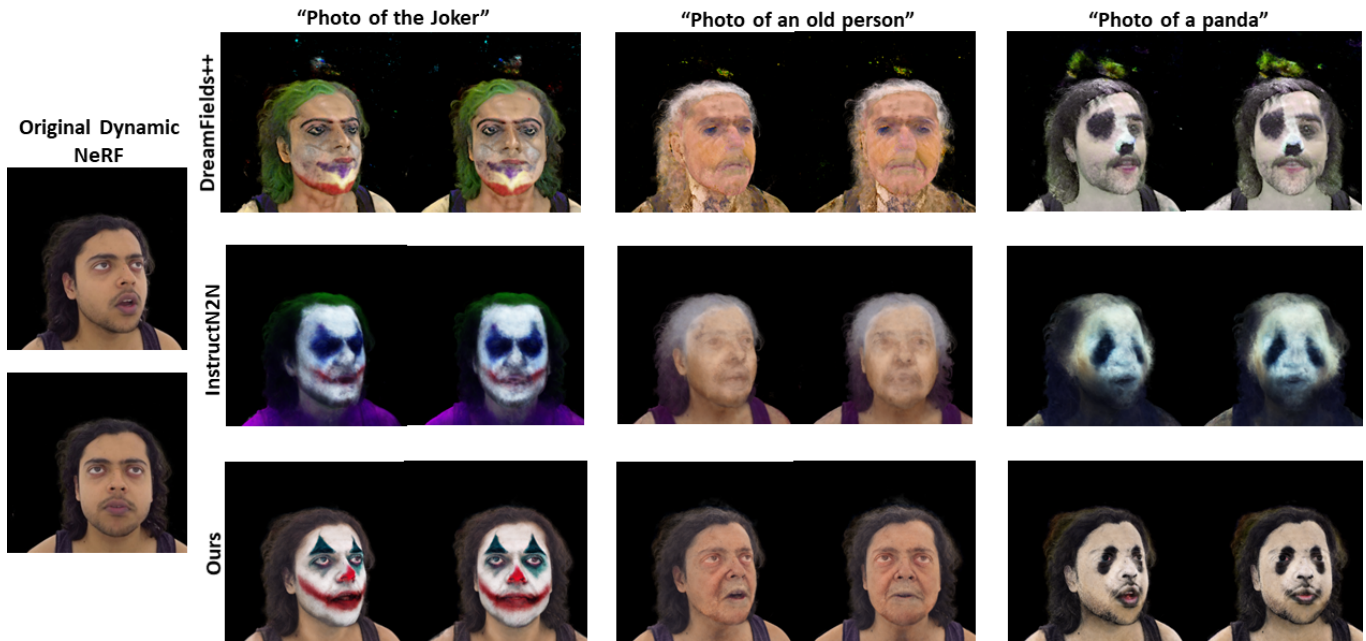


Fig. 11. Our approach outperforms Dream Fields++ [Jain et al. 2022] and InstructeNeRF2NeRF [Haque et al. 2023] spatially and temporally. It maintains the original identity and generates edits that are inline with the target prompts. Please see the supplemental video and the user study (Tab. 1) for video coherency.

Rendering. *ACM Trans. Graph.* 40, 4, Article 59 (jul 2021), 13 pages. https://doi.org/10.1145/3450626.3459863

Yuchen Lu, Soumye Singhal, Florian Strub, Aaron Courville, and Olivier Pietquin. 2020. Countering language drift with seeded iterated learning. In *International Conference on Machine Learning*. PMLR, 6437–6447.

B R Mallikarjun, Ayush Tewari, Abdallah Dib, Tim Weyrich, Bernd Bickel, Hans-Peter Seidel, Hanspeter Pfister, Wojciech Matusik, Louis Chevallier, Mohamed Elgharib, et al. 2021. PhotoApp: Photorealistic Appearance Editing of Head Portraits. *ACM Transactions on Graphics* 40, 4 (2021), 1–16.
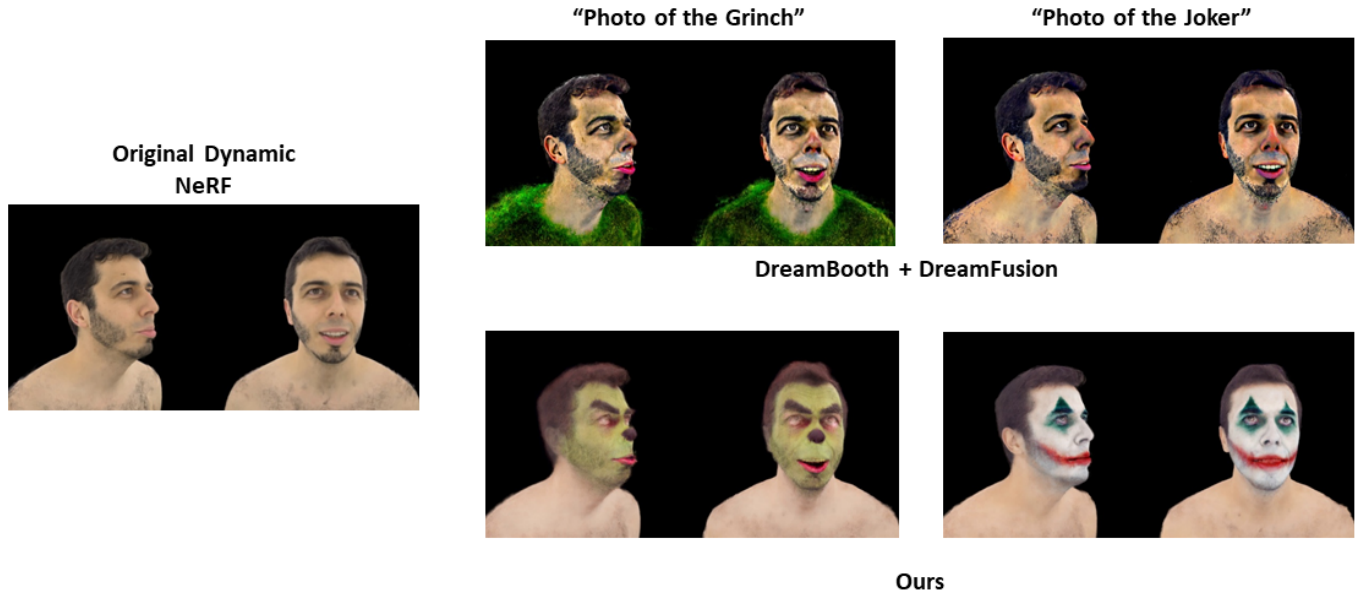
Fig. 12. An implementation that combines DreamBooth [Ruiz et al. 2022] with DreamFusion [Poole et al. 2022] has limited edibility (see top row). This is contrast to our method which produces edits that are clearly more inline with the target prompts.



Fig. 13. Comparison between our viewpoint-and-time-specific fine-tuning (Sec. 3.2.1) and DreamBooth fine-tuning. Ours has fewer artifacts and better preserves the identity, showing the importance of viewpoint-and-time-specific fine-tuning.

Luke Melas-Kyriazi, Christian Rupprecht, Iro Laina, and Andrea Vedaldi. 2023. RealFusion: 360° Reconstruction of Any Object from a Single Image. In *Arxiv*.

Marko Mihajlovic, Aayush Bansal, Michael Zollhoefer, Siyu Tang, and Shunsuke Saito. 2022. KeypointNeRF: Generalizing Image-based Volumetric Avatars using Relative Spatial Encoding of Keypoints. In *European conference on computer vision*.

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

Mehdi Mirza and Simon Osindero. 2014. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784* (2014).

Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. 2022. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Transactions on Graphics (ToG)* 41, 4 (2022), 1–15.

Omkar M. Parkhi, Andrea Vedaldi, and Andrew Zisserman. 2015. Deep Face Recognition. In *Proceedings of the British Machine Vision Conference (BMVC)*, Xianghua Xie, Mark W. Jones, and Gary K. L. Tam (Eds.). BMVA Press, Article 41, 12 pages. https://doi.org/10.5244/C.29.41

Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. 2022. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv preprint arXiv:2209.14988* (2022).

Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021. Learning Transferable Visual Models From Natural Language Supervision. In *Proceedings of the 38th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 139)*, Marina Meila and Tong Zhang (Eds.). PMLR, 8748–8763. https://proceedings.mlr.press/v139/radford21a.html

Alec Radford and Karthik Narasimhan. 2018. Improving Language Understanding by Generative Pre-Training.

Amit Raj, Michael Zollhoefer, Tomas Simon, Jason Saragih, Shunsuke Saito, James Hays, and Stephen Lombardi. 2021. PVA: Pixel-aligned Volumetric Avatars. arXiv:2101.02697 [cs.CV]

Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. 2022. Hierarchical Text-Conditional Image Generation with CLIP Latents. arXiv:2204.06125 [cs.CV]

Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. 2021. Zero-Shot Text-to-Image Generation. arXiv:2102.12092 [cs.CV]

Anurag Ranjan, Kwang Moo Yi, Jen-Hao Rick Chang, and Oncel Tuzel. 2023. FaceLit: Neural 3D Relightable Faces. In *CVPR*. https://arxiv.org/abs/2303.15437

Pramod Rao, Mallikarjun B R, Gereon Fox, Tim Weyrich, Bernd Bickel, Hans-Peter Seidel, Hanspeter Pfister, Wojciech Matusik, Ayush Tewari, Christian Theobalt, and Mohamed Elgharib. 2022. VoRF: Volumetric Relightable Faces. (2022).

Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. High-Resolution Image Synthesis with Latent Diffusion Models. In *Computer*
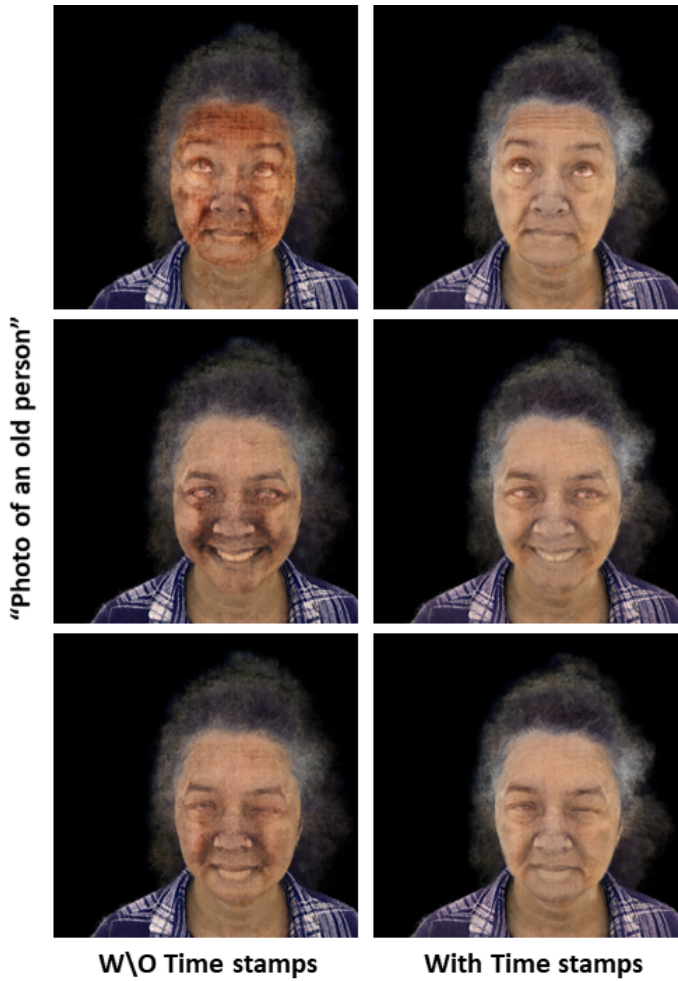
Fig. 14. Fine-tuning LDM (Sec. 3.2) without any temporal information leads to clear temporal inconsistencies. Notice the sharp temporal transition in the overall color.
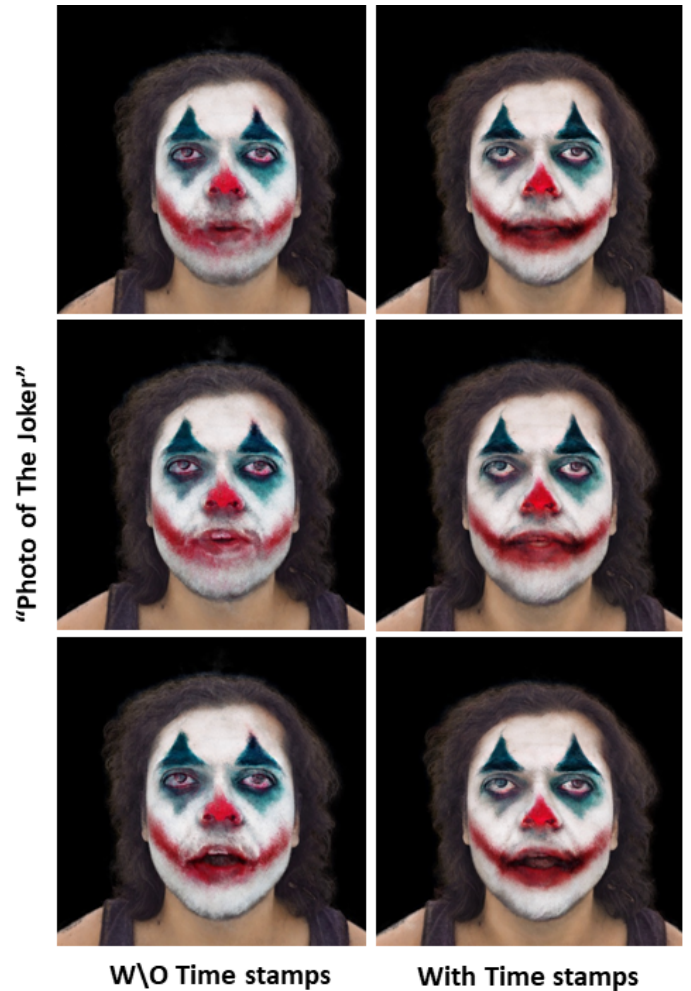


Fig. 15. Fine-tuning LDM (Sec. 3.2) without any temporal information leads to clear temporal inconsistencies (e.g. see mouth region).

Vision and Pattern Recognition (CVPR).

Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. 2022. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. *arXiv preprint arXiv:2208.12242* (2022).

Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. 2022. Laion-5b: An open large-scale dataset for training next generation image-text models. *arXiv preprint arXiv:2210.08402* (2022).

Ahmed Selim, Mohamed Elgharib, and Linda Doyle. 2016. Painting Style Transfer for Head Portraits using Convolutional Neural Networks. (2016), 129:1–129:18.

Aliaksandr Siarohin, Stéphane Lathuilière, Sergey Tulyakov, Elisa Ricci, and Nicu Sebe. 2019. First Order Motion Model for Image Animation. In *Conference on Neural Information Processing Systems (NeurIPS)*.

Tiancheng Sun, Jonathan T. Barron, Yun-Ta Tsai, Zexiang Xu, Xueming Yu, Graham Fyffe, Christoph Rhemann, Jay Busch, Paul Debevec, and Ravi Ramamoorthi. 2019. Single Image Portrait Relighting. *ACM Trans. Graph.* 38, 4, Article 79 (jul 2019), 12 pages. https://doi.org/10.1145/3306346.3323008

Supasorn Suwajanakorn, Steven M. Seitz, and Ira Kemelmacher-Shlizerman. 2017. Synthesizing Obama: Learning Lip Sync from Audio. *ACM Trans. Graph.* 36, 4, Article 95 (jul 2017), 13 pages. https://doi.org/10.1145/3072959.3073640

Feitong Tan, Sean Fanello, Abhimitra Meka, Sergio Orts-Escolano, Danhang Tang, Rohit Pandey, Jonathan Taylor, Ping Tan, and Yinda Zhang. 2022. VoLux-GAN: A Generative Model for 3D Face Synthesis with HDRI Relighting. arXiv:2201.04873 [cs.CV]

Kartik Teotia, Xingang Pan, Hyeongwoo Kim, Pablo Garrido, Mohamed Elgharib, Christian Theobalt, et al. 2023. HQ3DAvatar: High Quality Controllable 3D Head Avatar. *arXiv preprint arXiv:2303.14471* (2023).

Ayush Tewari, Mohamed Elgharib, Mallikarjun BR, Florian Bernard, Hans-Peter Seidel, Patrick Pérez, Michael Zöllhofer, and Christian Theobalt. 2020. PIE: Portrait Image Embedding for Semantic Control. *ACM Transactions on Graphics (Proceedings SIGGRAPH Asia)* 39, 6. https://doi.org/10.1145/3414685.3417803

Justus Thies, Mohamed Elgharib, Ayush Tewari, Christian Theobalt, and Matthias Nießner. 2020. Neural Voice Puppetry: Audio-driven Facial Reenactment. *ECCV 2020* (2020).

Aaron Van Den Oord, Oriol Vinyals, et al. 2017. Neural discrete representation learning. *Advances in neural information processing systems* 30 (2017).

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. 2017. Attention is All you Need. In *Advances in Neural Information Processing Systems*, I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (Eds.), Vol. 30. Curran Associates, Inc. https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf

Can Wang, Menglei Chai, Mingming He, Dongdong Chen, and Jing Liao. 2021b. CLIP-NeRF: Text-and-Image Driven Manipulation of Neural Radiance Fields. *arXiv preprint arXiv:2112.05139* (2021).

Can Wang, Ruixiang Jiang, Menglei Chai, Mingming He, Dongdong Chen, and Jing Liao. 2022. NeRF-Art: Text-Driven Neural Radiance Fields Stylization. *arXiv preprint arXiv:2212.08070* (2022).
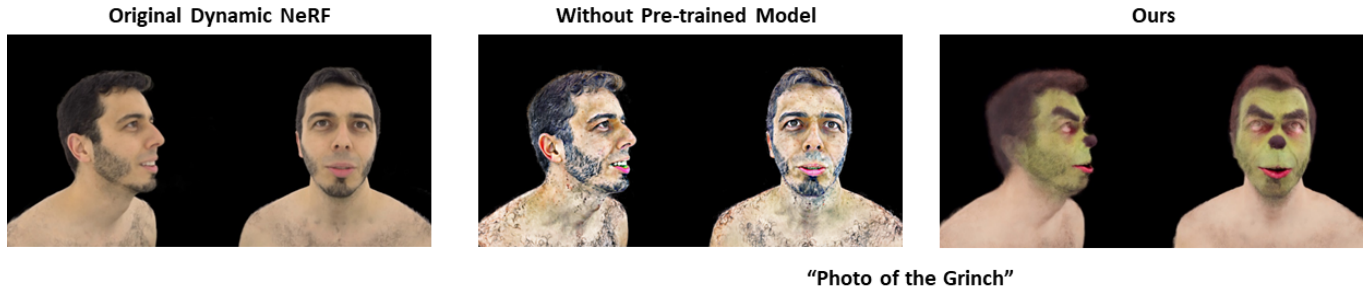
Fig. 16. Removing the pre-trained model during Score Distillation Sampling lead to a significant drop in performance.



Fig. 17. Using annealing during Score Distillation Sampling leads to better capturing of details.

Ziyan Wang, Timur Bagautdinov, Stephen Lombardi, Tomas Simon, Jason Saragih, Jessica Hodgins, and Michael Zollhofer. 2021a. Learning Compositional Radiance Fields of Dynamic Human Heads. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 5704–5713.

Shih-En Wei, Jason Saragih, Tomas Simon, Adam W. Harley, Stephen Lombardi, Michal Perdoch, Alexander Hypes, Dawei Wang, Hernan Badino, and Yaser Sheikh. 2019. VR Facial Animation via Multiview Image Translation. *ACM Trans. Graph.* 38, 4, Article 67 (jul 2019), 16 pages. https://doi.org/10.1145/3306346.3323030

Shuai Yang, Liming Jiang, Ziwei Liu, and Chen Change Loy. 2022. VToonify: Controllable High-Resolution Portrait Video Style Transfer. *ACM Transactions on Graphics (TOG)* 41, 6, Article 203 (2022), 15 pages. https://doi.org/10.1145/3550454.3555437

Egor Zakharov, Aliaksandra Shysheya, Egor Burkov, and Victor Lempitsky. 2019. Few-Shot Adversarial Learning of Realistic Neural Talking Head Models.

arXiv:1905.08233 [cs.CV]

Zhixing Zhang, Ligong Han, Arnab Ghosh, Dimitris Metaxas, and Jian Ren. 2022. SINE: SINgle Image Editing with Text-to-Image Diffusion Models. *arXiv preprint arXiv:2212.04489* (2022).

Yufeng Zheng, Victoria Fernández Abrevaya, Marcel C. Bühler, Xu Chen, Michael J. Black, and Otmar Hilliges. 2022. I M Avatar: Implicit Morphable Head Avatars from Videos. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR 2022)*. IEEE, Piscataway, NJ, 13535–13545. https://doi.org/10.1109/CVPR52688.2022.01318

Hao Zhou, Sunil Hadap, Kalyan Sunkavalli, and David Jacobs. 2019. Deep Single-Image Portrait Relighting. In *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*. 7193–7201. https://doi.org/10.1109/ICCV.2019.00729