

NeuWigs: A Neural Dynamic Model for Volumetric Hair Capture and Animation

Ziyan Wang^{1,2} Giljoo Nam² Tuur Stuyck² Stephen Lombardi^{3†} Chen Cao²
 Jason Saragih² Michael Zollhöfer² Jessica Hodgins^{1†} Christoph Lassner^{4†}
¹Carnegie Mellon University ²Reality Labs Research ³Google ⁴Epic Games

Abstract

The capture and animation of human hair are two of the major challenges in the creation of realistic avatars for the virtual reality. Both problems are highly challenging, because hair has complex geometry and appearance, as well as exhibits challenging motion. In this paper, we present a two-stage approach that models hair independently from the head to address these challenges in a data-driven manner. The first stage, state compression, learns a low-dimensional latent space of 3D hair states containing motion and appearance, via a novel autoencoder-as-a-tracker strategy. To better disentangle the hair and head in appearance learning, we employ multi-view hair segmentation masks in combination with a differentiable volumetric renderer. The second stage learns a novel hair dynamics model that performs temporal hair transfer based on the discovered latent codes. To enforce higher stability while driving our dynamics model, we employ the 3D point-cloud autoencoder from the compression stage for de-noising of the hair state. Our model outperforms the state of the art in novel view synthesis and is capable of creating novel hair animations without having to rely on hair observations as a driving signal. Project page is here <https://ziyanw1.github.io/neuwigs/>.

1. Introduction

The ability to model the details of human hair with high fidelity is key to achieving realism in human avatar creation because hair establishes part of our personal identity. The realism of human hair involves not only geometry, appearance and interaction with light, but also motion. The sheer number of hair strands leads to a very complex geometry, while the interactions between light and hair strands lead to non-trivial view-dependent appearance changes. Capturing the dynamics of hair is difficult due to the complexity of the motion space as well as severe self-occlusions. Generating realistic dynamics given only the head motion is similarly hard because the motion of hair is not solely controlled by

[†]Work done while at Meta.



Figure 1. **Animation from Single View Captures.** Our model can generate realistic hair animation from single view video based on head motion and gravity direction. Original captures of subjects wearing a wig cap are shown in red boxes.

the head position but also influenced by gravity and inertial forces. From a control perspective, hair does not respond linearly to the head position and a zero-order system is inadequate for modeling hair dynamics.

These challenges lead to two major problems for creating realistic hair avatars: 3D hair capture and dynamics modeling. While modern capture systems reconstruct the hair geometry and appearance from a sparse and discrete set of real world observations with high fidelity, they do not

directly solve the problem of novel motion generation. To achieve that, we need to go beyond reconstructing to create a *controllable* dynamic hair model using captured data.

In conventional animation techniques, hair geometry is created by an artist manually preparing 3D hair grooms. Motion of the 3D hair groom is created by a physics simulator where an artist selects the parameters for the simulation. This process requires expert knowledge. In contrast, data-driven methods aim to achieve hair capture and animation in an automatic way while preserving metric photo-realism. Most of the current data-driven hair capture and animation approaches learn to regress a dense 3D hair representation that is renderable directly from per-frame driving signals, without modeling dynamics.

However, there are several factors that limit the practical use of these data-driven methods for hair animation. First of all, these methods mostly rely on sophisticated driving signals, like multi-view images [21, 33], a tracked mesh of the hair [23], or tracked guide hair strands [45], which are hard to acquire. Furthermore, from an animation perspective, these models are limited to rendering hair based on hair observations and cannot be used to generate novel motion of hair. Sometimes it is not possible to record the hair driving signals at all. We might want to animate hair for a person wearing accessories or equipment that (partially) obstructs the view of their hair, for example VR glasses; or animate a novel hair style for a subject or animate hair for a bald person.

To address these limitations of existing data-driven hair capture and animation approaches, we present a neural dynamic model that is able to animate hair with high fidelity conditioned on head motion and relative gravity direction. By building such a dynamic model, we are able to generate hair motions by evolving an initial hair state into a future one, without relying on per-frame hair observation as a driving signal. We utilize a two-stage approach for creating this dynamic model: in the first stage, we perform state compression by learning a hair autoencoder from multi-view video captures with an evolving tracking algorithm. Our method is capable of capturing a temporally consistent, fully renderable volumetric representation of hair from videos with both head and hair. Hair states with different time-stamps are parameterized into a semantic embedding space via the autoencoder. In the second stage, we sample temporally adjacent pairs from the semantic embedding space and learn a dynamic model that can perform the hair state transition between each state in the embedding space given the previous head motion and gravity direction. With such a dynamic model, we can perform hair state evolution and hair animation in a recurrent manner which is not driven by existing hair observations. As shown in Fig 1, our method is capable of generating realistic hair animation with different hair styles on single view captures of a mov-

ing head with a bald cap. In summary, the contributions of this work are

- We present NeuWigs, a novel end-to-end data-driven pipeline with a volumetric autoencoder as the backbone for real human hair capture and animation, learnt from multi-view RGB images.
- We learn the hair geometry, tracking and appearance end-to-end with a novel autoencoder-as-a-tracker strategy for hair state compression, where the hair is modeled separately from the head using multi-view hair segmentation.
- We train an animatable hair dynamic model that is robust to drift using a hair state denoiser realized by the 3D autoencoder from the compression stage.

2. Related Work

We discuss related work in hair capture, data-driven hair animation, and volumetric avatars.

Static Hair Capture. Reconstructing static hair is challenging due to its complex geometry. Paris *et al.* [31] and Wei *et al.* [47] reconstruct 3D hair from multi-view images and 2D orientation maps. PhotoBooth [32] further improves this technique by applying calibrated projectors and cameras with patterned lights to reconstruct finer geometry. Luo *et al.* [24] and Hu *et al.* [12] leverage hair specific structure (strand geometry) and physically simulated strands for more robust reconstruction. The state-of-the-art work on static hair capture is Nam *et al.* [29]. They present a line-based patch method for reconstructing 3D line clouds. Sun *et al.* [42] extends this approach to reconstruct both geometry and appearance with OLAT images. Rosu *et al.* [38] presents a learning framework for hair reconstruction that utilizes a hair strand prior trained on synthetic data.

Dynamic Hair Capture. Hair dynamics is hard to capture as a result of the complex motion patterns and self-occlusions. Zhang *et al.* [53] refine hair dynamics by applying physical simulation techniques to per-frame hair reconstruction. Hu *et al.* [11] invert the problem of hair tracking by directly solving for hair dynamic parameters through iterative grid search on thousands of simulation results under different settings. Xu *et al.* [51] perform hair strand tracking by extracting hair strand tracklets from spatio-temporal slices of a video volume. Liang *et al.* [20] and Yang *et al.* [52] design a learning framework fueled by synthetic hair data to regress 3D hair from video. Winberg *et al.* [48] perform dense facial hair and underlying skin tracking under quasi-rigid motion with a multi-camera system. While achieving good results on the capture of the hair geometry from dynamic sequences, those methods either do not model hair appearance with photo-realism or do not solve the problem of drivable animation.

Data-driven Hair Animation. Using physics-based simulation for hair animation is a common practice in both,

academia and the film/games industry [1, 46]. However, generating hair animations with physics-based simulation can be computationally costly. To remedy this problem, reduced data-driven methods [5, 6, 9] simulate only a small portion of guide hair strands and interpolate the rest using skinning weights learned from full simulations. With the latest advances in deep learning, the efficiency of both dynamic generation [25] and rendering [4, 30] of hair has been improved using neural networks. Lyu *et al.* [25] uses deep neural networks for adaptive binding between normal hair and guide hair. Olszewski *et al.* [30] treats hair rendering as an image translation problem and generates realistic rendering of hair conditioned on 2D hair masks and strokes. Similarly, Chai *et al.* [4] achieves faster rendering with photorealistic results by substituting the rendering part in the animation pipeline with screen-space neural rendering techniques. Temporal consistency is enforced in this pipeline by conditioning on hair flow. However, those methods still build on top of conventional hair simulation pipelines and use synthetic hair wigs, which require manual efforts by an artist to set up and are non-trivial to metrically evaluate. Wu *et al.* [49] propose to use a secondary motion graph (SDG) for hair animation without relying on a conventional hair simulation pipeline at runtime. However, this method is limited by artist’s design of hair wigs and control of hair simulation parameters and is not able to capture or animate hair with realistic motion.

Volumetric Avatars. With the recent advent of differentiable volumetric raymarching [21, 28], many works attempt to directly build volumetric avatars from images or videos. To the best of our knowledge, Neural Volumes [22] is the earliest work that creates a volumetric head avatar from multiview images using differentiable volumetric raymarching. One of the many strengths of this work is that it directly optimizes a volume grid from multiview images while still producing high quality renders for semi-transparent objects like hair. One followup work [44] combines volumetric and coordinate-based representations into a hybrid form for better rendering quality and drivability. However, the level of detail either methods can capture is limited by the resolution of the volume grid.

Another early work on differentiable volumetric rendering is NeRF [28], which parameterizes radiance fields implicitly with MLPs instead of using a volumetric grid. Due to the success of NeRF [28] for modeling 3D scenes from multiple images, there are many works that build avatars with NeRF [7, 8, 10, 14, 16, 17, 27, 33, 34, 36, 55]. PVA [36] and KeypointNeRF [27] utilize pixel aligned information to extend NeRF’s drivability and generalization over sequence data. Nerfies [33], HyperNeRF [34] and TAVA [17] optimize a deformation field together with a NeRF in a canonical space given videos. NeRFace [7], IM Avatar [55] and HeadNeRF [10] substitute the deformation field with face

models like 3DMM [2] or FLAME [18] for better controllability. However, those methods mostly assume hair to be rigidly attached to the head without motion and most of the approaches suffer from prohibitively long rendering time.

In contrast to NeRF-based avatars, a mixture of volumetric primitives (MVP) [23] builds a volumetric representation that can generate high-quality *real time* renderings that look realistic even for challenging materials like hair and clothing. Several follow up works extend it to model body dynamics [37], moderate hair dynamics [45] and even for in-the-wild captures [3]. However, animating such volumetric representations with dynamics is still an unsolved problem.

3. Method

Our method for hair performance capture and animation consists of two stages: state compression and dynamic modeling (see also Fig. 2). The goal of the first stage is to perform dynamic hair capture from multi-view video of a head in motion. To be more specific, in this stage, we aim to distill a 3D renderable representation of hair from multi-view images at each frame into an embedding space. To achieve that, we train a volumetric autoencoder in a self-supervised manner to model the hair geometry, tracking and appearance. The output of this model is a set of tracked hair point clouds p_t , their corresponding local radiance fields in the form of volumetric primitives V_t and a compact embedding space that is spanned by the 1D hair state encoding z_t . In the second stage, we perform modeling of hair dynamics based on the hair capture from the first stage. The goal of this stage is to create a controllable, self-evolving representation of hair without relying on online observations of hair. We achieve this by learning a neural network to regress the next possible hair state, which is conditioned on the previous hair state as well as the previous head motion and head-relative gravity direction. Equipped with the hair encoding space acquired from the first stage, we can train that model in a supervised manner by simply sampling data pairs of temporally adjacent hair states from the encoding space. Using both stages, we can perform dynamic hair animation at test time given an initialization using a recurrent strategy, without relying on direct observations of hair as a per-frame driving signal.

3.1. State Compression

We assume that multi-view image captures I_{cam_i} with their corresponding calibrated cameras are given. We denote the extrinsics of each camera i as \mathbf{R}_i and \mathbf{t}_i . We then run l-MVS [29] and non-rigid tracking [50] to obtain the per-frame hair reconstructions $p_t \in \mathbb{R}^{N_{p_t} \times 3}$ and head tracked vertices $x_t \in \mathbb{R}^{N_{x_t} \times 3}$ at time frame t , where N_{p_t} and N_{x_t} denote the size of each. The p_t and x_t together serves as a coarse representation for the hair and head. Dif-

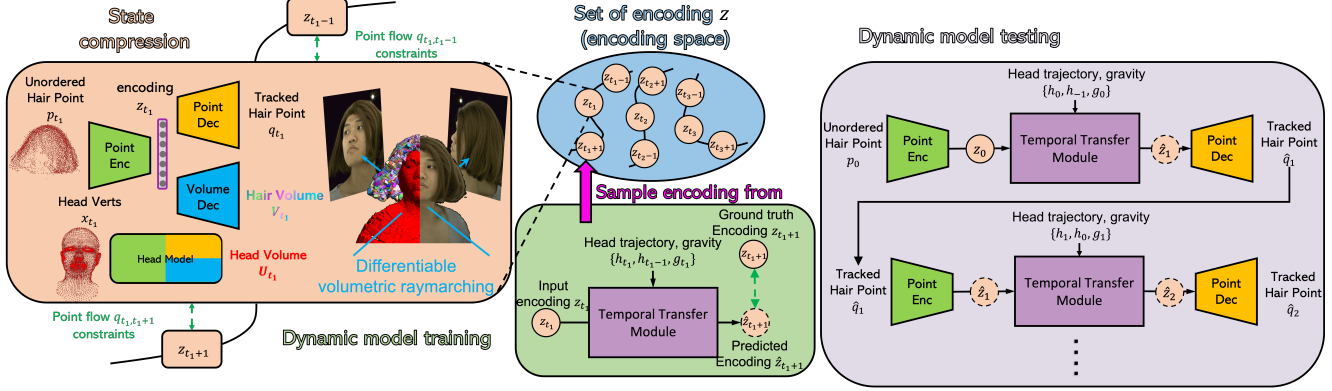


Figure 2. **Method Overview.** Our method is comprised of two stages: state compression and dynamic modeling. In the first stage, we train an autoencoder for hair and head appearance from multiview RGB images using differentiable volumetric raymarching; at the same time we create an encoding space of hair states. In the dynamic modeling stage, we sample temporally adjacent hair encodings to train a temporal transfer module (T²M) that performs the transfer between the two, based on head motion and head-relative gravity direction.

ferent from the head vertices, here the p_t represent an unordered set of hair point clouds. Due to the difference between hair and head dynamic patterns, we model them separately by training two different volumetric autoencoders. For the head model, we use an autoencoder to regress the volumetric texture in an unwrapped UV layout conditioned on the tracked head mesh, similar to [23, 45]. For the hair model, we optimize the hair volumetric texture and its tracking simultaneously. To better enforce the disentanglement between hair and head, we attach head volumes only to the head mesh and hair volumes only to the hair point clouds. Moreover, we use a segmentation loss to constrain each of them to only model the texture of their assigned category of hair or head.

Autoencoder as a Tracker. Learning to track hair in a supervised manner with manual annotation is infeasible. We automatically discover the hair keypoints as well as their tracking information by optimizing a variational autoencoder (VAE) [15] in a semi-supervised manner. By doing autoencoding on hair point clouds, we find that the VAE representing the hair point cloud can automatically align hair shapes along the temporal axis and is capable of tracking through both long and discontinuous hair video segments such as captures of different hair motions.

The input to the point encoder \mathcal{E} is the point coordinates of the unordered hair point cloud p_t . Given its innate randomness in terms of point coverage and order, we use PointNet [35] to extract the corresponding encoding $z_t \in \mathbb{R}^{256}$. Besides being agnostic to the order of p_t , it also can process varying numbers of points and aggregate global information from the input point cloud. The point decoder \mathcal{D} is a simple MLP that regresses the coordinate and point tangential direction of the tracked point cloud $q_t \in \mathbb{R}^{N_{prim} \times 3}$, $\mathbf{dir}(q_t) \in \mathbb{R}^{N_{prim} \times 3}$ and $s_t \in \mathbb{R}^{N_{prim}}$ from z_t , where N_{prim} is the number of tracked hair points. We denote

$\mathbf{dir}(x)$ as the tangential direction of x . s_t is a per-point scale factor which will be used later. We optimize the following loss to train the point autoencoder:

$$\mathcal{L}_{geo} = \mathcal{L}_{cham} + \omega_{temp} \mathcal{L}_{temp} + \omega_{KL} \mathcal{L}_{KL}.$$

The first term is the Chamfer distance loss which aims to align the shape of tracked point cloud q_t to p_t :

$$\begin{aligned} \mathcal{L}_{cham} = & \|q_t - N_{q_t, p_t}\|_2 - \cos(\mathbf{dir}(q_t), \mathbf{dir}(N_{q_t, p_t})) \\ & + \|p_t - N_{p_t, q_t}\|_2 - \cos(\mathbf{dir}(p_t), \mathbf{dir}(N_{p_t, q_t})), \end{aligned}$$

where $\cos(\cdot, \cdot)$ is the cosine similarity and $N_{x,y} \in \mathbb{R}^{N_x \times 3}$ are the coordinates of the nearest neighbor of each point of x in y . To further enforce temporal smoothness, we use point flow $\vec{fl}(p_t)$ and $\overleftarrow{fl}(p_t)$ denoting forward and backward flow from p_t to p_{t+1} and p_{t-1} as additional supervision and formulate \mathcal{L}_{temp} as follows:

$$\begin{aligned} \mathcal{L}_{temp} = & \|\vec{fl}(\hat{p}_t) - \vec{fl}(N_{q_t, p_t})\|_2 + \|\overleftarrow{fl}(p_t) - \overleftarrow{fl}(N_{p_t, q_t})\|_2 \\ & + \|\vec{fl}(\hat{p}_t) - \vec{fl}(N_{q_t, p_t})\|_2 + \|\overleftarrow{fl}(p_t) - \overleftarrow{fl}(N_{p_t, q_t})\|_2, \end{aligned}$$

where as q_t is the tracked point, we can simply have $\vec{fl}(q_t) = q_t - q_{t-1}$ and $\overleftarrow{fl}(q_t) = q_t - q_{t+1}$. Please see the supplemental materials for how we estimate $\vec{fl}(N_{q_t, p_t})$. The last term \mathcal{L}_{KL} is the KL-divergence loss [15] on the encoding z_t to enforce similarity with a normal distribution $\mathcal{N}(0, 1)$.

Hair Volumetric Decoder. In parallel to the point decoder \mathcal{D} , we optimize a hair volumetric decoder that regresses a volumetric radiance field around each of the hair points. The hair volumetric primitives $V_t \in \mathbb{R}^{N_{prim} \times 4 \times m^3}$ store RGB and alpha in resolution of m^3 . We use a decoder similar to HVH [45] to regress the volume payload. The pose of each volumetric primitive is directly determined by

the output of the point decoder: q_t and $\text{dir}(q_t)$. We denote $\mathbf{R}_t^{p,n} \in SO(3)$ and $\mathbf{d}_t^{p,n} \in \mathbb{R}^3$ as the n th volume-to-world rotation and translation of hair volume and per-hair-volume scale $s_t^{p,n}$ as the n th element in scale \mathbf{s}_t . Similarly, $\mathbf{d}_t^n = q_t^n$ is the n th element of q_t . Given the head center \mathbf{x}_t^c extracted from the head vertices \mathbf{x}_t and hair head direction as $\bar{\mathbf{h}}_t^n = q_t^n - \mathbf{x}_t^c$, we formulate the rotation $\mathbf{R}_t^{p,n}$ as $\mathbf{R}_t^{p,n} = [\mathbf{l}(q_t^n), \mathbf{l}(q_t^n \times \bar{\mathbf{h}}_t^n), \mathbf{l}(q_t^n \times (q_t^n \times \bar{\mathbf{h}}_t^n))]^T$, where $\mathbf{l}(x) = x/\|x\|_2$ is the normalization function. The output of the head model is similar to the hair part except that it is modeling head related (non-hair) regions. We denote the head volume payload as $\mathbf{U}_t \in \mathbb{R}^{N_{prim} \times 3 \times m^3}$ and head related rotation $\mathbf{R}_t^{x,n} \in SO(3)$, translation $\mathbf{d}_t^{x,n} \in \mathbb{R}^3$ and scale $s_t^{x,n} \in \mathbb{R}$.

Differentiable Volumetric Raymarching. Given all volume rotations $\mathbf{R}_t^{all} = [\mathbf{R}_t^{x,n}, \mathbf{R}_t^{p,n}]$, translations $\mathbf{d}_t^{all} = [\mathbf{d}_t^{x,n}, \mathbf{d}_t^{p,n}]$, scales $\mathbf{s}_t^{all} = [s_t^{x,n}, s_t^{p,n}]$ and local radiance fields $\mathbf{V}_t^{all} = [\mathbf{U}_t, \mathbf{V}_t]$, we can render them into image \mathcal{I}_{cam_i} and compare it with \mathcal{I}_{cam_i} to optimize all volumes. Using an optimized BVH implementation similar to MVP [23], we can efficiently determine how each ray intersects with each volume. We define a ray as $\mathbf{r}(\mathbf{p}, l) = \mathbf{o}(\mathbf{p}) + l\mathbf{v}(\mathbf{p})$ shooting from pixel \mathbf{p} in direction of $\mathbf{v}(\mathbf{p})$ with a depth l in range of (l_{min}, l_{max}) . The differentiable formation of an image given the volumes can then be formulated as below:

$$\mathcal{I}_p = \int_{l_{min}}^{l_{max}} \mathbf{V}_{t,rgb}^{all}(\mathbf{r}_p(l)) \frac{dT(l)}{dl} dl,$$

$$T(l) = \min(\int_{l_{min}}^l \mathbf{V}_{t,\alpha}^{all}(\mathbf{r}_p(l)) dl, 1),$$

where $\mathbf{V}_{t,rgb}^{all}$ is the RGB part of \mathbf{V}_t^{all} and $\mathbf{V}_{t,\alpha}^{all}$ is the alpha part of \mathbf{V}_t^{all} . To get the full rendering, we composite the rendered image as $\tilde{\mathcal{I}}_p = \mathcal{I}_p + (1 - \mathcal{A}_p)I_{p,bg}$ where $\mathcal{A}_p = T(l_{max})$ and $I_{p,bg}$ is the background image. We optimize the following loss to train the volume decoder:

$$\mathcal{L}_{pho} = \|\tilde{\mathcal{I}}_p - I_{p,gt}\|_1 + \omega_{VGG} \mathcal{L}_{VGG}(\tilde{\mathcal{I}}_p, I_{p,gt}),$$

where \mathcal{L}_{VGG} is the perceptual loss in [13] and $I_{p,gt}$ is the ground truth pixel value of \mathbf{p} . We find that the usage of a perceptual loss yields more salient rendering results.

However, as we optimize both \mathbf{U}_t and \mathbf{V}_t from the images, texture bleeding between the hair volume \mathbf{V}_t and the head volume \mathbf{U}_t becomes a problem. The texture bleeding issue is especially undesirable when we want to treat the hair and head separately, for example, when we want to animate just the hair. To prevent this, we additionally render a hair mask map to regularize both $\mathbf{V}_{t,\alpha}$ and $\mathbf{U}_{t,\alpha}$. We denote the ground truth hair mask as $M_{p,gt}$ and the rendered hair

mask as \mathcal{M}_p :

$$\mathcal{M}_p = \int_{l_{min}}^{l_{max}} \mathbf{V}_{t,1}^{all}(\mathbf{r}_p(l)) \frac{dT(l)}{dl} dl,$$

$$T(l) = \min(\int_{l_{min}}^l \mathbf{V}_{t,\alpha}^{all}(\mathbf{r}_p(l)) dl, 1),$$

where $\mathbf{V}_{t,1}^{all}$ is all one volume if it belongs to \mathbf{V}_t otherwise zero. We formulate the segmentation loss as $\mathcal{L}_{mask} = \|\mathcal{M}_p - M_{p,gt}\|_1$. The final objective for training the whole autoencoder is $\mathcal{L} = \mathcal{L}_{geo} + \mathcal{L}_{pho} + \omega_{mask} \mathcal{L}_{mask}$.

3.2. Dynamic Model

In the second stage, we aim to build a dynamic model that can evolve hair states over time without relying on per-frame hair observation as a driving signal. To achieve that, we leverage the embedding space of hair states built at the state compression stage and train a model that performs the hair state transfer in a supervised manner. To this end, we build a temporal transfer module (T²M) of hair dynamic priors, that evolves the hair state and can produce hair animation based on the indirect driving signals of head motion and head-relative gravity direction in a self-evolving manner. The design of T²M is similar to a hair simulator except that it is fully data-driven. One of the inputs to T²M is the encoding \mathbf{z}_{t-1} of the previous time step. At the same time, T²M is also conditioned on the head per-vertex displacement $\mathbf{h}_t = \mathbf{x}_t - \mathbf{x}_{t-1}$ and $\mathbf{h}_{t-1} = \mathbf{x}_{t-1} - \mathbf{x}_{t-2}$ from the previous two time steps as well as head relative gravity direction $\mathbf{g}_t \in \mathbb{R}^3$ at the current time step. T²M will then predict the next possible state $\hat{\mathbf{z}}_t$ from T²M based on those inputs. Similar to the design of the VAE, the output of T²M is a distribution instead of a single vector. To be more specific, the mean and standard deviation of $\hat{\mathbf{z}}_t$ are $\mu(\hat{\mathbf{z}}_t), \delta(\hat{\mathbf{z}}_t) = \text{T}^2\text{M}(\mathbf{z}_{t-1} | \mathbf{h}_t, \mathbf{h}_{t-1}, \mathbf{g}_t)$. During training, we take $\hat{\mathbf{z}}_t = \mu(\hat{\mathbf{z}}_t) + \mathbf{n} \odot \delta(\hat{\mathbf{z}}_t)$ and during testing $\hat{\mathbf{z}}_t = \mu(\hat{\mathbf{z}}_t)$. The per point normal distribution vector, \mathbf{n} , is the same shape as $\delta(\hat{\mathbf{z}}_t)$ and \odot is the element-wise multiplication.

Training Objectives. We denote the point encoder as $\mathcal{E}(\cdot)$ and the point decoder as $\mathcal{D}(\cdot)$. Across the training of T²M, we freeze the parameters of both, \mathcal{E} and \mathcal{D} . We denote the unordered point cloud at frame t as p_t , its corresponding encoding as $\mu(\mathbf{z}_t), \delta(\mathbf{z}_t) = \mathcal{E}(p_t)$ and the tracked point cloud as $q_t = \mathcal{D}(\mathbf{z}_t)$. The following loss enforces the prediction of T²M to be similar to its ground truth:

$$\mathcal{L}_{mse} = \|\mu(\hat{\mathbf{z}}_{t+1}) - \mu(\mathbf{z}_{t+1})\|_2 + \|\delta(\hat{\mathbf{z}}_{t+1}) - \delta(\mathbf{z}_{t+1})\|_2$$

$$\mathcal{L}_{cos} = -\cos(\mu\hat{\mathbf{z}}_{t+1}, \mu\mathbf{z}_{t+1}) - \cos(\delta\hat{\mathbf{z}}_{t+1}, \delta\mathbf{z}_{t+1})$$

$$\mathcal{L}_{ptsmse} = \|\mathcal{D}(\hat{\mathbf{z}}_{t+1}) - \mathcal{D}(\mathbf{z}_{t+1})\|_2,$$

where we not only minimize the ℓ_2 distance between $\hat{\mathbf{z}}_{t+1}$ and \mathbf{z}_{t+1} , but also enforce the cosine similarity and the corresponding tracked point cloud to be equivalent. To adapt

T²M to the tracked point cloud q_t , we compute the above loss again but using q_t as input, where we generate the corresponding encoding as z'_t from $\mathcal{E}(q_t)$ and its prediction \hat{z}'_{t+1} from T²M($z'_t|h_t, h_{t-1}, g_t$):

$$\begin{aligned}\mathcal{L}_{mse,cyc} &= \|\mu(\hat{z}'_{t+1}) - \mu(z_{t+1})\|_2 + \|\delta(\hat{z}'_{t+1}) - \delta(z_{t+1})\|_2 \\ \mathcal{L}_{cos,cyc} &= -\cos(\mu\hat{z}'_{t+1}, \mu z_{t+1}) - \cos(\delta\hat{z}'_{t+1}, \delta z_{t+1}) \\ \mathcal{L}_{ptsmse,cyc} &= \|\mathcal{D}(\hat{z}'_{t+1}) - \mathcal{D}(z_{t+1})\|_2.\end{aligned}$$

Similar to how we train our autoencoder, we also enforce two KL divergence losses on both the predicted \hat{z}_{t+1} and \hat{z}'_{t+1} with a normal distribution \mathcal{N} . The final objective for training the T²M is a weighted sum of the above eight terms. **Animation.** Given an initialized hair state, our dynamic model T²M can evolve the hair state into future states conditioned on head motion and head-relative gravity direction. One straightforward implementation of the would be to directly propagate the hair state encoding z_t . However, in practice, we find this leads to severe drift in the semantic space. As a simple feed forward neural network, T²M can not guarantee that its output is noise free. The noise in the output becomes even more problematic when we use T²M in a recurrent manner, where the output noise will aggregate and lead to drift. To remedy this, instead of propagating the encoding z_t directly, we reproject the predicted encoding z_t by the point autoencoder \mathcal{E} and \mathcal{D} every time for denoising. To be more specific, we acquire the de-noised predicted hair encoding $\hat{z}_{t+1} = \mathcal{E}(\mathcal{D}(\hat{z}_{t+1}))$ from the raw prediction \hat{z}_{t+1} of T²M. The use of the point autoencoder \mathcal{E} and \mathcal{D} can help us remove the noise in z_{t+1} , as the point cloud encoder \mathcal{E} can regress the mean μ_{t+1} and standard deviation δ_{t+1} of z_{t+1} separately by using q_{t+1} as an intermediate variable. Thus, we can extract the noise free part of z_{t+1} by taking the mean μ_{t+1} regressed from \mathcal{E} . Please see our experiments for further details of this approach.

4. Experiments

In order to test our proposed model, we conduct experiments on both the hair motion data set presented in HVH [45] and our own dataset with longer sequences following a similar capture protocol as HVH [45]. We collect a total of four different hair wig styles with scripted head motions like nodding, swinging and tilting. We also collect an animation test set with the same scripted head motions performed by different actors wearing a wig cap, which we will refer to as “bald head motion sequence”. The animation test set contains both single view captures from a smart phone and multiview captures. The total length of each hair wig capture is around 1-1.5 minutes with a frame rate of 30Hz. 100 cameras are used during the capture where 93 of them are used to obtain training views and the rest are providing held-out test views. We split each sequence into two folds with similar amounts of frames and train our model

	seq01				seq02				seq03			
	MSE↓	PSNR↑	SSIM↑	LPIPS↓	MSE↓	PSNR↑	SSIM↑	LPIPS↓	MSE↓	PSNR↑	SSIM↑	LPIPS↓
PENeRF	51.25	31.16	0.9269	0.3717	103.41	28.15	0.8659	0.5067	76.59	29.50	0.9000	0.2949
NSFF	50.13	31.21	0.9346	0.3672	90.06	28.75	0.8885	0.4728	83.18	29.10	0.8936	0.3292
NRNeRF	56.78	30.78	0.9231	0.3554	132.16	27.13	0.8549	0.5241	79.83	29.33	0.8987	0.3067
MVP	47.54	31.60	0.9476	0.2587	77.23	29.62	0.9088	0.3051	73.78	29.66	0.9224	0.2455
HVH	41.89	32.17	0.9543	0.2019	59.84	30.69	0.9275	0.2353	71.58	29.81	0.9314	0.2021
Ours	40.34	32.28	0.9558	0.1299	56.47	30.94	0.9329	0.1254	73.65	29.69	0.9247	0.1496

Table 1. **Novel view synthesis.** We compute MSE↓, PSNR↑, SSIM↑ and LPIPS↓ comparing rendered and ground truth images on hold-out views. **First** and **second** best results are highlighted.

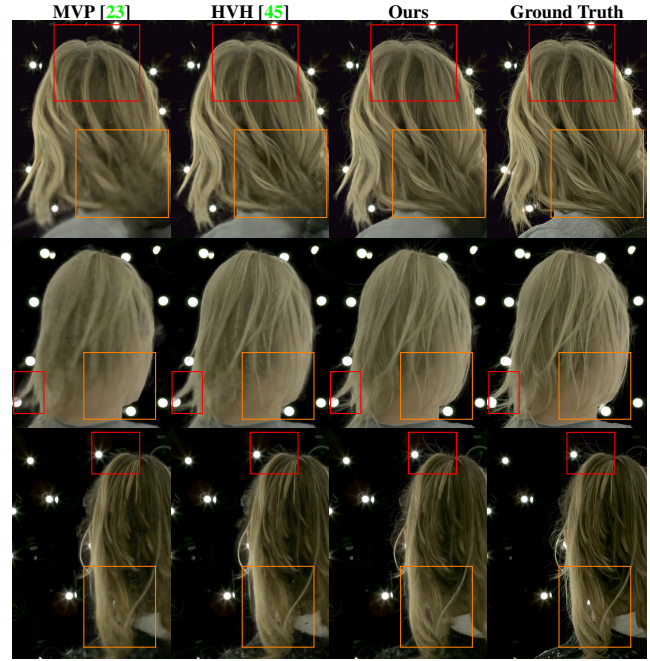


Figure 3. **Novel View Synthesis.** Compared with previous methods, our method captures hair with more details, including fly-away hair strands and creates an overall more accurate hair reconstruction with perceptually better rendering results.

exclusively on the training portion of each sequence.

4.1. Evaluation of the State Compression Model

We first test our state compression model to evaluate its ability to reconstruct the appearance of hair and head.

Novel View Synthesis. We compare with volumetric methods like NeRF based methods [19, 43] and volumetric primitives based methods [23, 45] on the data set from HVH [45]. In Tab 1, we show the reconstruction related metrics MSE, SSIM, PSNR and LPIPS [54] between predicted images and the ground truth image on hold out views. Our method yields a good balance between perceptual loss and reconstruction loss while keeping both of them relatively low. Furthermore, our method achieves a much higher perceptual similarity with ground truth images. In Fig. 3 we show that our method can capture high frequency details and even preserve some fly-away hair strands.

	Seq01	Seq02	Seq03
HVH [45]	0.6685	0.4121	0.3766
Ours	0.8289	0.9243	0.8571

Table 2. **IoU**(\uparrow) between rendered hair silhouette and ground truth hair segmentation.

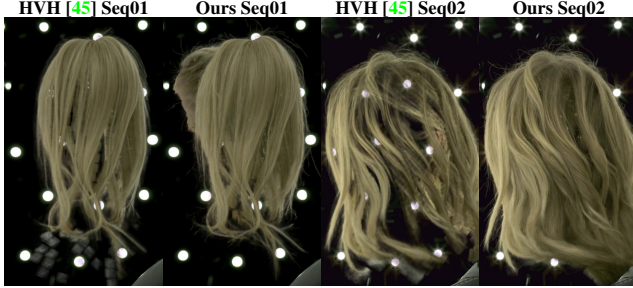


Figure 4. **Hair/Head Disentanglement.** By explicitly enforcing the semantic segmentation of head and hair through additional supervision, we learn a more opaque hair texture while the result suffers less from texture bleeding.

Ablation on Hair/Head Disentanglement. We test how well our model handles hair/head disentanglement compared to previous work. As hair and head are exhibiting different dynamic patterns, disentanglement is usually required, especially to achieve independent controllability for both. We compare with HVH, which implicitly separates the hair and head using the dynamic discrepancy between them and optical flow. In our model, we further facilitate the disentanglement by using semantic segmentation. In Tab 2, we show the IoU between the rendered silhouette of hair volumes and ground truth hair segmentation of the different methods. We visualize the difference in Fig. 4. Our method generates a more opaque hair texture with less texture bleeding between hair volumes and non-hair volumes. Moreover, our model creates the hair shape in an entirely data-driven fashion, which yields higher fidelity results than the artist prepared hair in HVH [45].

Ablation on \mathcal{L}_{VGG} . We examine the synergy between \mathcal{L}_{VGG} and the ℓ_1 loss for improving the rendering quality. As shown in Tab. 3, we find that the perceptual loss has positive effects on the reconstruction performance while the improvements are negated when the weight is too large. In Fig. 5, we compare the rendered images using different \mathcal{L}_{VGG} weights. The results are more blurry when not using \mathcal{L}_{VGG} and fewer details are reconstructed, such as fly-away strands.

Ablation on Point Flow Supervision. Although with \mathcal{L}_{cham} we can already optimize a reasonably tracked point cloud p_t , we find that point flow can help remove the jittering in appearance. We show the temporal smoothness enforced by the point flow supervision in Fig. 6. Our model

	vgg=0.0	vgg=0.1	vgg=0.3	vgg=1.0	vgg=3.0	vgg=10.0
MSE	42.84	42.40	39.94	40.34	40.98	42.82
PSNR	32.04	32.09	32.34	32.28	32.25	32.07
SSIM	0.9544	0.9564	0.9576	0.9558	0.9541	0.9518
LPIPS	0.2021	0.1765	0.1511	0.1299	0.1238	0.1257

Table 3. **Ablation on \mathcal{L}_{VGG} .** We find using an additional complementary perceptual loss leads to better appearance reconstruction.

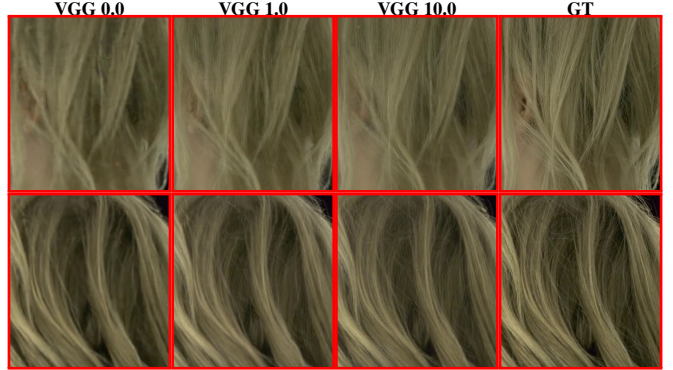


Figure 5. **Ablation on \mathcal{L}_{VGG} .** Adding a perceptual loss leads to sharper reconstruction results.

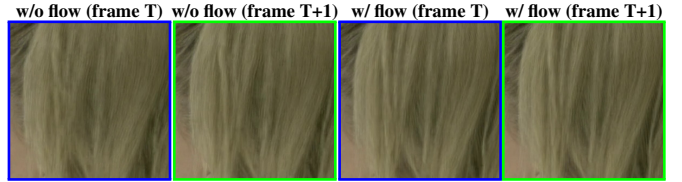


Figure 6. **Ablation on Point Flow.** We find that adding point flow to regularize the offsets between temporally adjacent tracked points prevents jittering.

learns a more consistent hair texture with less jittering when trained with point flow. Please see the videos in the supplemental material for a visualization over time.

4.2. Evaluation of the Dynamic Model

Lastly, we perform tests of our animation model. Compared to a per-frame model that takes hair observations as input, the input to our dynamic model and any of its variations is a subset of the head motion trajectory and hair point cloud at the initialization frame. As a quantitative evaluation, we compare our model with per-frame driven models using either hair observations or head observations as driving signals. For qualitative evaluation, we render new hair animations for the bald head motion sequences.

Quantitative Test of the Dynamic Model. We evaluate our dynamic model on the test sequences of scripted hair motion capture. The goal is to test whether our dynamic model generates reasonable novel content rather than only testing how well it reconstructs the test sequence. To test the performance of our dynamic model, we treat the model driven

	MSE(↓)	PSNR(↑)	SSIM(↑)	LPIPS(↓)	ChamDis(↓)
pf w/ hair img	37.36	32.94	0.9559	0.1333	10.47
dyn w/o cos	44.96	32.26	0.9458	0.1327	25.49
dyn w/o cyc	45.22	32.23	0.9453	0.1335	26.79
dyn w/o grav	40.12	32.64	0.9504	0.1268	13.76
dyn	38.49	32.80	0.9532	0.1211	11.12

Table 4. **Ablation of Different Dynamic Models.** We compare different models in terms of rendering quality and tracking accuracy.

by per-frame (**pf**) hair observations as an oracle paradigm to compare with, while our dynamic model does not use any per-frame hair observation as a driving signal. In Tab. 4, we compare the rendering quality of our dynamic model with **pf** models and ablate several designs. The best performing dynamic model (**dyn**) has similar performance with the **pf** model even without the per-frame hair observation as driving signal. We find that both adding a cosine similarity loss as an additional objective to MSE and adding a cycle consistency loss helps improve the stability of the dynamic model. Meanwhile, we find adding gravity as an auxiliary input stabilizes our model on slow motions. This improvement might be related to the fact that during slow motions of the head, the hair motion is primarily driven by gravity.

Ablation for the Point Autoencoder $\mathcal{E} + \mathcal{D}$. Here, we analyze how well the point autoencoder acts as a stabilizer for the dynamic model. We compare two different models in Fig. 7: **encprop**, that propagates the encoding directly, and **ptsprop** that propagates the regressed hair point cloud and generates the corresponding encoding from the point encoder. The results are more salient with **ptsprop**. The improvement of the **ptsprop** over the **encprop** is partially because the mapping from point cloud to encoding is an injective mapping and the point encoder serves as a noise canceller in the encoding space. To further study this behavior, we perform a cycle test on the point encoder, where we add noise n to a certain encoding z and get a noisy version of the encoding $\hat{z} = z + n$ and its corresponding noisy point cloud \hat{p} . Then, we predict a cycled encoding $\bar{z} = \mathcal{E}(\mathcal{D}(\hat{z}))$. We compare z , \hat{z} and \bar{z} in Fig. 7. z and \bar{z} are consistently close while \hat{z} jitters. This result suggests that the remapping of \hat{z} using \mathcal{E} and \mathcal{D} counteracts the noise n .

Animation on Bald Head Sequences. We show animation results driven by head motions on in-the-wild phone video captures in Fig 1. We find our model generates reasonable motions of hair under natural head motions like swinging or nodding. Please refer to the supplemental materials and videos for visualizations over time.

5. Discussion

We present a two-stage data-driven pipeline for volumetric hair capture and animation. The core of our method is a 3D volumetric autoencoder that we find useful for both, au-

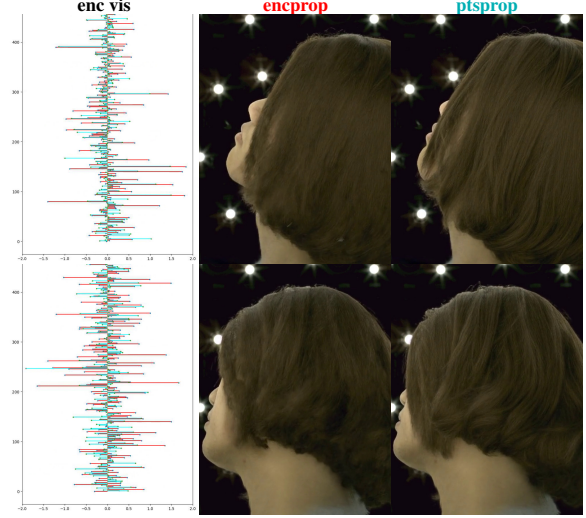


Figure 7. **encprop v.s. ptsprop.** ptsprop generates sharper results with less drifting than encprop.

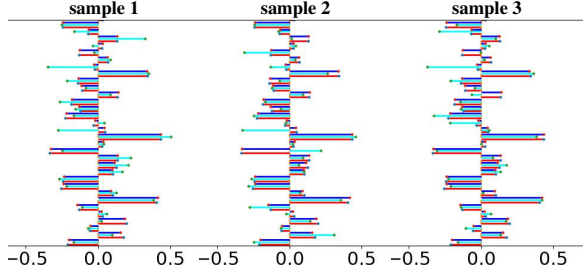


Figure 8. **Point Encoder \mathcal{E} as a Stabilizer.** We sample several \hat{z} and corresponding $\bar{z} = \mathcal{E}(\mathcal{D}(\hat{z}))$ with a fixed z and visualize part of them above. As we can see, \bar{z} stays similar to z while \hat{z} jitters.

tomatic hair state acquisition and stable hair dynamic generation. The first stage of our pipeline simultaneously performs hair tracking, geometry and appearance reconstruction in a self-supervised manner via an autoencoder-as-a-tracker strategy. The second stage leverages the hair states acquired from the first stage and creates a recurrent hair dynamics model that is robust to moderate drift with the autoencoder as a denoiser. We empirically show that our method performs stable tracking of hair on long and segmented video captures while preserving high fidelity for hair appearance. Our model also supports generating new animations in both, lab- and in-the-wild conditions and does not rely on hair observations.

Limitations. Like many other data-driven methods, our method requires a large amount of diverse sstraining data and might fail with data that is far from the training distribution. Our model is currently not relightable and can not animate new hairstyles. One future direction is to separately model appearance and lighting, which can be learnt from varied lighting captures. Another interesting direction is to

learn a morphable hair model for new hairstyle adaptation.

6. Appendix

6.1. Network Architecture

Here we provide details about how we design our neural networks and further information about training.

Encoder. As training a point cloud encoder solely is extremely unstable, we first train an image encoder and use it as a teacher model to train the point cloud encoder. In practice, we train two encoders for our hair branch together. Here we first illustrate the structure of both encoders. We will go back to how we train them and use them later. One of the encoders is an image encoder which is a convolutional neural network (CNN) that takes multiple view images as input. We denote the image encoder as \mathcal{E}_{img} . The other one is \mathcal{E} which is a PointNet encoder that takes either an unordered hair point cloud p_t or a tracked hair point cloud q_t as input. Positional encoding [28] is applied to the raw point cloud coordinate before it is used as the input to the network. We find this is very effective to help the network capturing high frequency details. In practice, we use frequencies of x^2 where x ranges from 1 to 7. We show the detailed architecture of \mathcal{E}_{img} in Tab 5. The architecture of the point cloud encoder \mathcal{E} is shown in Tab 6. Both of the two encoders \mathcal{E} and \mathcal{E}_{img} can produce a latent vector in size of 256, which are supposed to describe the same content. Their output will be passed to \mathcal{E}_μ and \mathcal{E}_σ which are two linear layers that produce μ and σ of z_t respectively.

	Encoder \mathcal{E}_{img}
1	Conv2d(3, 64)
2	Conv2d(64, 64)
3	Conv2d(64, 128)
4	Conv2d(128, 128)
5	Conv2d(128, 256)
6	Conv2d(256, 256)
7	Conv2d(256, 256)
8	Flatten()
9	Linear($256 \times n_{inimg} \times 15$, 256)

Table 5. **Encoder \mathcal{E}_{img} architecture.** Each Conv2d layer in the encoder has a kernel size of 3, stride of 1 and padding of 1. Weight normalization [41] and untied bias are applied. After each layer, except for the last two parallel fully-connected layers, a Leaky ReLU [26] activation with a negative slope of 0.2 is applied. Then a downsample layer with a stride of 2 is applied after every conv2d layer. The first linear layer takes the concatenation of all towers from different image views as input. n_{inimg} stands for much many views we take.

Point Decoder. We use a 3-layer MLP as the point decoder \mathcal{D} , which takes a 1d latent code z_t as input and outputs the

	Encoder \mathcal{E}
1	Conv2d(3, 128)
2	Conv2d(128, 256)
3	Conv2d(256, 256)
4	Conv2d(256, 256)
5	Conv2d(256, 512)
6	Conv2d(512, 512)
7	Conv2d(512, 512)
8	Conv2d(512, 1024)
8	MAM pooling()
9	Linear(1024×3 , 512)
10	Linear(512, 256)
11	Linear(256, 256)

Table 6. **Encoder \mathcal{E} architecture.** We use a \mathcal{E} structure similar to PointNet [35]. All Conv2d uses a kernel of 1 and stride of 1, which serves as a shared MLP. We only use Conv2d for simpler implementation. After each Conv2d layer, a Leaky ReLU [26] activation with a negative slope of 0.2 is applied. Then we use a MAM pool layer to aggregate features from all points. MAM stands for min, average and max pooling, which concatenates the results of min, average and max pooling into one. Then, two linear layers are applied to the output of MAM pooling and generate a 256 latent vector.

coordinate of the corresponding tracked point cloud q_t . We show the architecture of \mathcal{D} in Tab 7.

	Decoder \mathcal{D}
1	Linear(256, 256)
2	Linear(256, 256)
2	Linear(256, 4096×3)

Table 7. **Decoder \mathcal{D} architecture.** We use an MLP with three Linear layers as the decoder \mathcal{D} . After each layer except the last layer, a Leaky ReLU [26] activation with a negative slope of 0.2 is applied.

Volume Decoder. The volumetric model is a stack of 2D deconv layers. We align the x-axis and y-axis of each volume and put them onto a 2D imaginary UV-space. Then we convolve on them to regress the z-axis content for each of the x,y position. We show the architecture of the volume decoder in Tab 8. In our setting, we have two separate volume decoder for both RGB volume and alpha volume.

Dynamic Model. We use three different inputs to the dynamic model T^2M , namely the hair encoding z_{t-1} at the previous frame $t - 1$, the head velocity h_{t-1} and h_{t-2} from the previous two frames $t - 1$ and $t - 2$, and the head relative gravity direction g_t at the current frame t . We first encode $\{h_{t-1}, h_{t-2}\}$ and g_t into two 1d vectors with 128 dimensions respectively. Then, we concatenate them together with encoding z_{t-1} as the input to another MLP to regress

Volume Decoder	
global encoding z_t	per-point hair feature
repeat	
concat	
1	Linear(320, 512)
2	deconv2d(512, 256)
3	conv2d(256, 256)
4	deconv2d(256, 256)
5	conv2d(256, 256)
6	deconv2d(256, 128)
7	conv2d(128, 128)
8	deconv2d(128, 16×ch)

Table 8. **Architecture of the Volume Decoder.** We first repeat the global encoding z_t into the shape of the per-point hair feature. The per-point hair feature is a tensor that is shared across all time frames. We then concatenate those two into one. Each layer except for the last one is followed by a Leaky ReLU layer with a negative slope of 0.2. Each `deconv2d` layer has a filter size of 4, stride size of 2 and padding size of 1. Each `conv2d` layer has a filter size of 3, stride size of 1 and padding size of 1. `ch` stands for the channel size of the output. It is set to 3 if it is an rgb decoder and 1 for an alpha decoder.

the next possible hair state encoding z_t . As in Tab. 9, we show the flow of T²M. For the head velocity branch, we first extract the per-vertex velocity $h_{t-1} = x_t - x_{t-1}$ where x_t is the coordiante of the tracked head mesh at frame t . To be noted, here the h_{t-1} contains only the information of the rigid head motion but not any other non-rigid motion like expression change. This representation of head motion is redundant theoretically, but we find it helps our network to converge better comparing to just using the pure 6-DoF head rotation and translation. We then reshape it and use it as the input to a two layer MLP to extract a 1d encoding of size 128. For the gravity branch, we first encode the gravity direction g_t with cosine encoding [28]. The output of the dynamic model is the mean μ_{t+1} and standard deviation σ_{t+1} of the predicted hair state z_{t+1} .

6.2. Training details

Dataset and Capture Systems. Following the setting in HVH [45], we also captured several video sequences with scripted hair motion performed under different hair styles for animation tests. During the capture, we ask the participants to put on different kind of hair wigs and perform a variety of head motions like nodding, swinging and tilting for multiple times under both slow and fast speed. To collect a demonstration set for animation, we also ask the participants to put on a hair net (bare head) and perform the same set of motions as when they are wearing a hair wig.

Hair Point Flow Estimation. There are three steps for computing the hair point flow, namely per-point feature de-

Temporal Transfer Module (T ² M)			
1	head velocity $\{h_{t-1}, h_{t-2}\}$	head relative gravity g_t	hair state z_{t-1}
2	Linear(7306×3, 256)	cosine encoding	
3	Linear(256, 128)		
4	Linear(539, 256)		
5	Linear(256, 256)		
6	Linear(256, 256)		
7	Linear(256, 256)	Linear(256, 256)	

Table 9. **Temporal Transfer Module (T²M).** We first encode the head velocity $\{h_{t-1}, h_{t-2}\}$ and head relative gravity g_t into 1d vectors, with a 2-layer MLP and cosine encoding respectively. Then we concatenate hair state z_{t-1} with those vectors to serve as the input to another MLP. The last two layers will be regressing the mean μ_{t+1} and standard deviation σ_{t+1} of the predicted hair state z_{t+1} . All Linear expect for the last two are followed by a Leaky ReLU activation with a negative slope of 0.2.

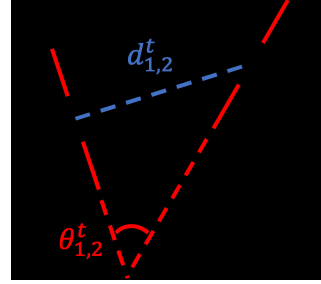


Figure 9

scriptor extraction, feature matching and flow filtering. In the first step, we compute a per-point feature descriptor based on the distribution of each point’s local neighboring. In the second step, we match the points from two adjacent time steps based on the similarity between their feature descriptor. In the last step, we filter out outlier flows that are abnormal.

To compute the point feature descriptor, we construct Line Feature Histograms (LFH) inspired by Point Feature Histograms (PFH) [40]. The LFH is a histogram of a 4-tuple that describes the spatial relationship between a certain point p_1^t and its neighboring point p_2^t . As shown in Fig. 9, we visualize two points $p_1^t \in \mathbb{R}^3$ and $p_2^t \in \mathbb{R}^3$ from the same time step t . Given p_1^t and p_2^t , we define the following four properties that describe their spatial relationship. The first one is the relative position of p_2^t with respect to p_1^t , which is $d_{1,2}^t = p_2^t - p_1^t$. Then we could compute the relative distance as $\|d_{1,2}^t\|_2 \in \mathbb{R}$. The second term is the angle $\theta_{1,2}^t$ between $dir(p_1^t)$ and $dir(p_2^t)$, where $dir(x)$ is the line direction of x from [29]. The last two terms are the angles $\alpha_{1,2}^t$ and $\beta_{1,2}^t$ between $(dir(p_1^t), d_{1,2}^t)$ and $(dir(p_2^t), d_{1,2}^t)$ respectively. For all intersections, we take the acute angle, which means $\theta_{1,2}^t, \alpha_{1,2}^t$ and $\beta_{1,2}^t$ are in $[0, \pi/2]$. Thus, the 4-tuple we used to create $LFH(p_1^t)$ is $(\|d_{1,2}^t\|_2, \theta_{1,2}^t, \alpha_{1,2}^t, \beta_{1,2}^t)$ and we normalize the histogram by its l2 norm. The designed LFH has three good proper-

ties. As we use the normalized feature, it is density invariant. Since $\theta_{1,2}^t$, $\alpha_{1,2}^t$ and $\beta_{1,2}^t$ are always acute angles, the feature is also rotation and flip invariant, where if we flip or rotate $\text{dir}(\mathbf{p}_1^t)$ the histogram are still the same. This design helps us getting a more robust feature descriptor for matching. We set the resolution for each entries of the 4-tuple to be 4 and it results in a descriptor in size of 256.

In the second step, we compute the correspondence between points from adjacent time frames t and $t + t_\delta$ where $t_\delta \in \{-1, 1\}$. We use the method from Rusu *et al.* [39] to compute the correspondence between two point clouds from t and $t + t_\delta$. To further validate the flow we get, we use several heuristic to filter out some obvious outliers. We first discard all the flows that have a large magnitude. As the flow is computed between two adjacent frames, it is supposed to be not abrupt. The second heuristic we use to filter the outliers is cycle consistency, where we compute the flow both forward and backward to see if we can map back to the origin. If the mapped back point departs too far away from the origin, we discard their flow.

Training of Encoder. As mentioned before, we train two encoders \mathcal{E} and \mathcal{E}_{img} together. In practice, we find that directly training \mathcal{E} is not very stable and might lead to not able to converge. Thus, we learn the two encoders in a teach-student manor, where we use \mathcal{E}_{img} as a teach model to train \mathcal{E} . We denote $\mathbf{x}_{img,t}$ as the output of \mathcal{E}_{img} and $\mathbf{x}_{pt,t}$ as the output of \mathcal{E} . Then, we formulate the following MSE loss to enforce the \mathcal{E} to output similarly to \mathcal{E}_{img} :

$$\mathcal{L}_{ts} = \|\mathbf{x}_{img,t} - \mathbf{x}_{pt,t}\|_2,$$

where we restraint the gradient from \mathcal{L}_{ts} from back-propagating to \mathcal{E}_{img} while training.

6.3. Ablation on Different Encoders

We show quantitative evaluations on rendering quality of different encoders on both SEEN and UNSEEN sequences in Tab. 10. As we can see, our \mathcal{E} performs similarly to the \mathcal{E}_{img} on the novel views of the SEEN sequence. This is as expected due to the nation of teach-student model and we train our model on the SEEN sequence with the training views. On the UNSEEN sequence, we find our \mathcal{E} performs better than \mathcal{E}_{img} . The reason could be that there are smaller domain gap between the point clouds from SEEN sequence and UNSEEN sequence while the multi-view images might be varying a lot due to the head motion. And the CNN is not good for handling such changes due to the head motion while point encoder can process point clouds with better structure awareness.

	MSE↓	PSNR↑	SSIM↑	LPIPS↓
\mathcal{E} on SEEN	29.48	34.05	0.9657	0.1109
\mathcal{E}_{img} on SEEN	29.44	34.05	0.9657	0.1109
\mathcal{E} on UNSEEN	34.97	33.21	0.9587	0.1209
\mathcal{E}_{img} on UNSEEN	37.36	32.94	0.9559	0.1333

Table 10. **Metrics on Novel Views.** We show quantitative results of different encoders under both SEEN and UNSEEN sequence of the same hair styles.

6.4. Ablation on Different Designs of the Dynamic Model

We show the comparisons of different dynamic models and per-frame driven models in Tab. 11. We find that our model offers a significant improvement over the per-frame driven model that takes head pose or motion as input. This is, because innately the hair motion is not only determined by the head pose or the previous history of head pose but also the initial status of the hair. In Fig. 10, we visualize how each model drifts by plotting the Chamfer distance between the regressed point cloud and the ground truth point cloud. We find that adding head relative gravity direction can improve the model performance on slow motions.

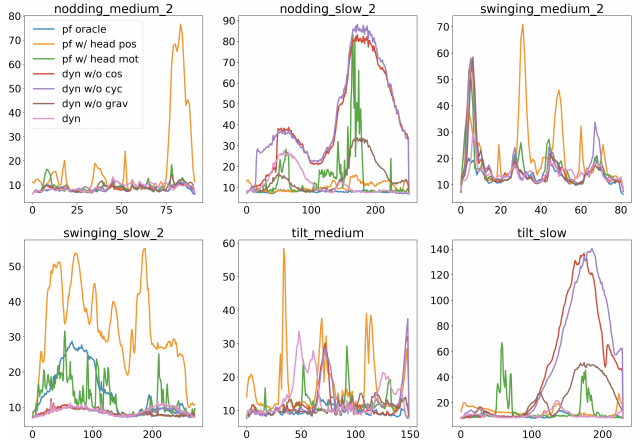


Figure 10. **ChamDist v.s. time.** We plot Chamfer distance v.s. time of different dynamic models to show drifting.

	MSE(↓)	PSNR(↑)	SSIM(↑)	LPIPS(↓)	ChamDis(↓)
pf w/ hair img	37.36	32.94	0.9559	0.1333	10.47
pf w/ hair pts	34.97	33.21	0.9587	0.1209	10.46
pf w/ head pos	47.43	32.01	0.9458	0.1522	18.94
pf w/ head mot	40.25	32.64	0.9508	0.1333	13.31
dyn w/o cos	44.96	32.26	0.9458	0.1327	25.49
dyn w/o cyc	45.22	32.23	0.9453	0.1335	26.79
dyn w/o grav	40.12	32.64	0.9504	0.1268	13.76
dyn	38.49	32.80	0.9532	0.1211	11.12

Table 11. **Ablation of Different Dynamic Models.**

6.5. Effect of the Initialization

We test how robust our model is to the initialization of the hair point cloud. In Fig. 11, we show animation results from models of two different hair styles (**hs**) with different initialization hair point clouds. The results look sharp when the model is matched with the correct hair style, but blurry when we use mismatched hair point clouds for initialization. However, we find our model self-rectifies and returns to a stable state after a certain number of iterations. This could be partially due to the model prior stored in the point encoder.



Figure 11. **Effect of the Initialization.** We initialize two different models (hs1 mod. and hs2 mod.) with two different hair point clouds (hs1 and hs2) in two time steps. The green box indicates matched initialization while orange indicates mismatched initialization. Although the mismatched initialization starts shows blurry results at first, the model automatically corrects itself when there is no head motion.

6.6. Further Ablation on the Point Encoder

To further study the point encoder \mathcal{E} 's ability on denoising the encoding, we tested the encoder with inputs that containing different level of noise. Similar to the study we did in the main paper, we first extract a fixed encoding z and add noise n to it as $\hat{z}=z+n$ and do noise removal as $\bar{z}=\mathcal{E}(\mathcal{D}(\hat{z}))$. We extend this by adding different level of noise by multiplying n with different scalars. Please refer to the video navigation page for more details.

6.7. Further Ablation on Novel View Synthesis

We compare our method with MVP [23] on the longer sequences we captured with scripted head motion. reconstruction related metrics are shown in Tab.12. We found that NeRF-based methods can not fit to longer sequences properly. This problem might be due to the large range of motion exhibited in the videos as well as the length of the video. HVH is not applicable because it does not support

	MSE(↓)	PSNR(↑)	SSIM(↑)	LPIPS(↓)
MVP	66.21	30.36	0.9291	0.2830
Ours	29.44	34.05	0.9657	0.1109

Table 12

hair tracking across segmented sequences of different hair motion. Compared to MVP, we achieve better reconstruction accuracy and improved perceptual similarity between rendered image and ground truth with the hair specific modeling in our design.

6.8. Animation on Bald Head Sequences

We show animation results driven by head motions both on lab conditioned multi-view video captures and in-the-wild phone video captures. For results on in-the-wild phone video captures, please refer to the supplemental videos. For phone captures, we ask the participants to face the frontal camera of the phone and perform different head motions. Then, we apply the face tracking algorithm in [3] to obtain face tracking data that serves as the input to our method. The initial hair state of the phone animation is sampled from the lab captured dataset. We find our model generates reasonable motions of hair under head motions like swinging, nodding and etc.

We also test our model on lab conditioned multi-view video captures. As shown in Figs. 12, 13, our model generates reasonable hair motions with respect to the head motion while preserving multi-view consistency.

References

- [1] Florence Bertails, Sunil Hadap, Marie-Paule Cani, Ming Lin, Tae-Yong Kim, Steve Marschner, Kelly Ward, and Zoran Kačić-Alesić. Realistic hair simulation: animation and rendering. In *ACM SIGGRAPH 2008 classes*, pages 1–154. 2008. 3
- [2] Volker Blanz and Thomas Vetter. A morphable model for the synthesis of 3d faces. In *Proceedings of the 26th annual conference on Computer graphics and interactive techniques*, pages 187–194, 1999. 3
- [3] Chen Cao, Tomas Simon, Jin Kyu Kim, Gabe Schwartz, Michael Zollhoefer, Shun-Suke Saito, Stephen Lombardi, Shih-En Wei, Danielle Belko, Shou-I Yu, Yaser Sheikh, and Jason Saragih. Authentic volumetric avatars from a phone scan. *ACM Trans. Graph.*, 41(4), jul 2022. 3, 12
- [4] Menglei Chai, Jian Ren, and Sergey Tulyakov. Neural hair rendering. In *European Conference on Computer Vision*, pages 371–388. Springer, 2020. 3
- [5] Menglei Chai, Changxi Zheng, and Kun Zhou. A reduced model for interactive hairs. *ACM Transactions on Graphics (TOG)*, 33(4):1–11, 2014. 3
- [6] Menglei Chai, Changxi Zheng, and Kun Zhou. Adaptive skinning for interactive hair-solid simulation. *IEEE transac-*



Figure 12. Animation on Bald Sequence.



Figure 13. Animation on Bald Sequence.

- tions on visualization and computer graphics, 23(7):1725–1738, 2016. 3
- [7] Guy Gafni, Justus Thies, Michael Zollhöfer, and Matthias Nießner. Dynamic neural radiance fields for monocular 4d facial avatar reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8649–8658, June 2021. 3
- [8] Philip-William Grassal, Malte Prinzler, Titus Leistner, Carsten Rother, Matthias Nießner, and Justus Thies. Neural head avatars from monocular rgb videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 18653–18664, June 2022. 3
- [9] Peng Guan, Leonid Sigal, Valeria Reznitskaya, and Jessica K Hodgins. Multi-linear data-driven dynamic hair model with efficient hair-body collision handling. In *Proceedings of the 11th ACM SIGGRAPH/Eurographics conference on Computer Animation*, pages 295–304, 2012. 3
- [10] Yang Hong, Bo Peng, Haiyao Xiao, Ligang Liu, and Juyong Zhang. Headnerf: A real-time nerf-based parametric head model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 20374–20384, June 2022. 3
- [11] Liwen Hu, Derek Bradley, Hao Li, and Thabo Beeler. Simulation-ready hair capture. In *Computer Graphics Forum*, volume 36, pages 281–294. Wiley Online Library, 2017. 2
- [12] Liwen Hu, Chongyang Ma, Linjie Luo, and Hao Li. Robust hair capture using simulated examples. *ACM Transactions on Graphics (TOG)*, 33(4):1–10, 2014. 2
- [13] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *European conference on computer vision*, pages 694–711. Springer, 2016. 5
- [14] Kacper Kania, Kwang Moo Yi, Marek Kowalski, Tomasz Trzciński, and Andrea Tagliasacchi. CoNeRF: Controllable Neural Radiance Fields. In *Proceedings of the IEEE Con-*



Figure 14. Animation on Bald Sequence.



Figure 15. Animation on Bald Sequence.

ference on Computer Vision and Pattern Recognition, 2022. 3

- [15] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013. 4
- [16] Anastasiia Kornilova, Marsel Faizullin, Konstantin Pakulev, Andrey Sadkov, Denis Kukushkin, Azat Akhmetyanov, Timur Akhtyamov, Hekmat Taherinejad, and Gonzalo Ferrer. Smartportraits: Depth powered handheld smartphone dataset of human portraits for state estimation, reconstruction and synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 21318–21329, June 2022. 3
- [17] Ruilong Li, Julian Tanke, Minh Vo, Michael Zollhofer, Jurgen Gall, Angjoo Kanazawa, and Christoph Lassner. Tava:

Template-free animatable volumetric actors. 2022. 3

- [18] Tianye Li, Timo Bolkart, Michael J Black, Hao Li, and Javier Romero. Learning a model of facial shape and expression from 4d scans. *ACM Transactions on Graphics (TOG)*, 36(6):194–1, 2017. 3
- [19] Zhengqi Li, Simon Niklaus, Noah Snavely, and Oliver Wang. Neural scene flow fields for space-time view synthesis of dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6498–6508, 2021. 6
- [20] Shu Liang, Xiufeng Huang, Xianyu Meng, Kunyao Chen, Linda G Shapiro, and Ira Kemelmacher-Shlizerman. Video to fully automatic 3d hair model. *ACM Transactions on Graphics (TOG)*, 37(6):1–14, 2018. 2



Figure 16. Animation on Bald Sequence.



Figure 17. Animation on Bald Sequence.

- [21] Stephen Lombardi, Tomas Simon, Jason Saragih, Gabriel Schwartz, Andreas Lehrmann, and Yaser Sheikh. Neural volumes: Learning dynamic renderable volumes from images. *ACM Trans. Graph.*, 38(4), 2019. [2](#), [3](#)
- [22] Stephen Lombardi, Tomas Simon, Jason Saragih, Gabriel Schwartz, Andreas Lehrmann, and Yaser Sheikh. Neural volumes: Learning dynamic renderable volumes from images. *ACM Transactions on Graphics (TOG)*, 38(4), July 2019. [3](#)
- [23] Stephen Lombardi, Tomas Simon, Gabriel Schwartz, Michael Zollhoefer, Yaser Sheikh, and Jason Saragih. Mixture of volumetric primitives for efficient neural rendering. *ACM Transactions on Graphics (TOG)*, 40(4), July 2021. [2](#), [3](#), [4](#), [5](#), [6](#), [12](#)
- [24] Linjie Luo, Hao Li, Sylvain Paris, Thibaut Weise, Mark Pauly, and Szymon Rusinkiewicz. Multi-view hair capture using orientation fields. In *2012 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1490–1497. IEEE, 2012. [2](#)
- [25] Qing Lyu, Menglei Chai, Xiang Chen, and Kun Zhou. Real-time hair simulation with neural interpolation. *IEEE Transactions on Visualization and Computer Graphics*, 2020. [3](#)
- [26] Andrew L Maas, Awni Y Hannun, and Andrew Y Ng. Rectifier nonlinearities improve neural network acoustic models. In *Proc. icml*, volume 30, 2013. [9](#)
- [27] Marko Mihajlovic, Aayush Bansal, Michael Zollhoefer, Siyu Tang, and Shunsuke Saito. Keypointnerf: Generalizing image-based volumetric avatars using relative spatial encoding of keypoints, 2022. [3](#)

- [28] 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*. Springer, 2020. 3, 9, 10
- [29] Giljoo Nam, Chenglei Wu, Min H Kim, and Yaser Sheikh. Strand-accurate multi-view hair capture. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 155–164, 2019. 2, 3, 10
- [30] Kyle Olszewski, Duygu Ceylan, Jun Xing, Jose Echevarria, Zhili Chen, Weikai Chen, and Hao Li. Intuitive, interactive beard and hair synthesis with generative models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7446–7456, 2020. 3
- [31] Sylvain Paris, Hector M Briceno, and François X Sillion. Capture of hair geometry from multiple images. *ACM Transactions on Graphics (TOG)*, 23(3):712–719, 2004. 2
- [32] Sylvain Paris, Will Chang, Oleg I Kozhushnyan, Wojciech Jarosz, Wojciech Matusik, Matthias Zwicker, and Frédo Durand. Hair photobooth: geometric and photometric acquisition of real hairstyles. *ACM Transactions on Graphics (TOG)*, 27(3):30, 2008. 2
- [33] Keunhong Park, Utkarsh Sinha, Jonathan T Barron, Sofien Bouaziz, Dan B Goldman, Steven M Seitz, and Ricardo Martin-Brualla. Nerfies: Deformable neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5865–5874, 2021. 2, 3
- [34] Keunhong Park, Utkarsh Sinha, Peter Hedman, Jonathan T. Barron, Sofien Bouaziz, Dan B Goldman, Ricardo Martin-Brualla, and Steven M. Seitz. Hypernerf: A higher-dimensional representation for topologically varying neural radiance fields. *arXiv preprint arXiv:2106.13228*, 2021. 3
- [35] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2017. 4, 9
- [36] Amit Raj, Michael Zollhofer, Tomas Simon, Jason Saragih, Shunsuke Saito, James Hays, and Stephen Lombardi. Pixel-aligned volumetric avatars. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11733–11742, June 2021. 3
- [37] Edoardo Remelli, Timur Bagautdinov, Shunsuke Saito, Chenglei Wu, Tomas Simon, Shih-En Wei, Kaiwen Guo, Zhe Cao, Fabian Prada, Jason Saragih, et al. Drivable volumetric avatars using texel-aligned features. In *ACM SIGGRAPH 2022 Conference Proceedings*, pages 1–9, 2022. 3
- [38] Radu Alexandru Rosu, Shunsuke Saito, Ziyang Wang, Chenglei Wu, Sven Behnke, and Giljoo Nam. Neural strands: Learning hair geometry and appearance from multi-view images. *ECCV*, 2022. 2
- [39] Radu Bogdan Rusu, Nico Blodow, Zoltan Csaba Marton, and Michael Beetz. Aligning point cloud views using persistent feature histograms. In *2008 IEEE/RSJ international conference on intelligent robots and systems*, pages 3384–3391. IEEE, 2008. 11
- [40] Radu Bogdan Rusu, Zoltan Csaba Marton, Nico Blodow, and Michael Beetz. Persistent point feature histograms for 3d point clouds. In *Proc 10th Int Conf Intel Autonomous Syst (IAS-10)*, Baden-Baden, Germany, pages 119–128, 2008. 10
- [41] Tim Salimans and Durk P Kingma. Weight normalization: A simple reparameterization to accelerate training of deep neural networks. *Advances in neural information processing systems*, 29, 2016. 9
- [42] Tiancheng Sun, Giljoo Nam, Carlos Aliaga, Christophe Hery, and Ravi Ramamoorthi. Human Hair Inverse Rendering using Multi-View Photometric data. In Adrien Bousseau and Morgan McGuire, editors, *Eurographics Symposium on Rendering - DL-only Track*. The Eurographics Association, 2021. 2
- [43] Edgar Tretschk, Ayush Tewari, Vladislav Golyanik, Michael Zollhöfer, Christoph Lassner, and Christian Theobalt. Non-rigid neural radiance fields: Reconstruction and novel view synthesis of a dynamic scene from monocular video. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12959–12970, 2021. 6
- [44] Ziyang Wang, Timur Bagautdinov, Stephen Lombardi, Tomas Simon, Jason Saragih, Jessica Hodgins, and Michael Zollhofer. Learning compositional radiance fields of dynamic human heads. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5704–5713, 2021. 3
- [45] Ziyang Wang, Giljoo Nam, Tuur Stuyck, Stephen Lombardi, Michael Zollhöfer, Jessica Hodgins, and Christoph Lassner. Hvh: Learning a hybrid neural volumetric representation for dynamic hair performance capture. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6143–6154, June 2022. 2, 3, 4, 6, 7, 10
- [46] Kelly Ward, Florence Bertails, Tae-Yong Kim, Stephen R Marschner, Marie-Paule Cani, and Ming C Lin. A survey on hair modeling: Styling, simulation, and rendering. *IEEE transactions on visualization and computer graphics*, 13(2):213–234, 2007. 3
- [47] Yichen Wei, Eyal Ofek, Long Quan, and Heung-Yeung Shum. Modeling hair from multiple views. In *ACM SIGGRAPH 2005 Papers*, pages 816–820. 2005. 2
- [48] Sebastian Winberg, Gaspard Zoss, Prashanth Chandran, Paulo Gotardo, and Derek Bradley. Facial hair tracking for high fidelity performance capture. *ACM Transactions on Graphics (TOG)*, 41(4):1–12, 2022. 2
- [49] Chenlei Wu and Takashi Kanai. Data-driven detailed hair animation for game characters. *Computer Animation and Virtual Worlds*, 27(3-4):221–230, 2016. 3
- [50] Chenglei Wu, Takaaki Shiratori, and Yaser Sheikh. Deep incremental learning for efficient high-fidelity face tracking. *ACM Transactions on Graphics (TOG)*, 37(6):1–12, 2018. 3
- [51] Zexiang Xu, Hsiang-Tao Wu, Lvdi Wang, Changxi Zheng, Xin Tong, and Yue Qi. Dynamic hair capture using spacetime optimization. *ACM Transactions on Graphics (TOG)*, 33(6), nov 2014. 2
- [52] Lingchen Yang, Zefeng Shi, Youyi Zheng, and Kun Zhou. Dynamic hair modeling from monocular videos using deep neural networks. *ACM Transactions on Graphics (TOG)*, 38(6):1–12, 2019. 2

- [53] Qing Zhang, Jing Tong, Huamin Wang, Zhigeng Pan, and Ruigang Yang. Simulation guided hair dynamics modeling from video. In *Computer Graphics Forum*, volume 31, pages 2003–2010. Wiley Online Library, 2012. [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 *CVPR*, 2018. [6](#)
- [55] Yufeng Zheng, Victoria Fernández Abrevaya, Marcel C. Bühler, Xu Chen, Michael J. Black, and Otmar Hilliges. I m avatar: Implicit morphable head avatars from videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 13545–13555, June 2022. [3](#)