

Audio-driven Neural Gesture Reenactment with Video Motion Graphs

Yang Zhou^{1,2} Jimei Yang² Dingzeyu Li² Jun Saito² Deepali Aneja² Evangelos Kalogerakis¹
¹University of Massachusetts Amherst ²Adobe Research

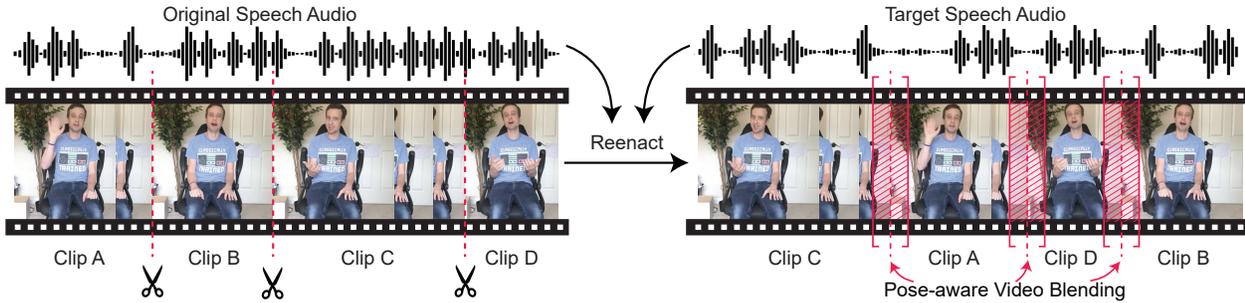


Figure 1. Given an input reference video of a speaker (left), our method reenacts it with gestures matching a target speech audio (right). The video is synthesized by re-assembling clips from the reference video and blending the inconsistent boundaries with a pose-aware neural network such that the synthesized video is coherent visually and consistent with both the rhythm and content of the target audio.

Abstract

Human speech is often accompanied by body gestures including arm and hand gestures. We present a method that reenacts a high-quality video with gestures matching a target speech audio. The key idea of our method is to split and re-assemble clips from a reference video through a novel video motion graph encoding valid transitions between clips. To seamlessly connect different clips in the reenactment, we propose a pose-aware video blending network which synthesizes video frames around the stitched frames between two clips. Moreover, we developed an audio-based gesture searching algorithm to find the optimal order of the reenacted frames. Our system generates reenactments that are consistent with both the audio rhythms and the speech content. We evaluate our synthesized video quality quantitatively, qualitatively, and with user studies, demonstrating that our method produces videos of much higher quality and consistency with the target audio compared to previous work and baselines. Our project page https://yzhou359.github.io/video_reenact includes code and data.

1. Introduction

Gesture is a key visual component for human speech communication [32]. It enhances the expressiveness of human performance and helps the audience to better comprehend the speech content [19]. Given the progress in talking head generation [17, 67, 85, 87], synthesizing plausible

gesture videos becomes increasingly important for applications such as digital voice assistants [48] and photo-realistic virtual avatars [24, 85]. In this paper, we propose an audio-driven gesture reenactment system that synthesizes speaker-specific human speech video from a target audio clip and a single reference speech video (Figure 1).

Unlike lip motions with specific phoneme-to-viseme mappings [20, 65, 88] or facial expressions mostly corresponding to low-frequency sentimental signals [69], gestures exhibit complex relationships with not only acoustics but also semantics of the audio [46]. Therefore, it is nontrivial to find a direct cross-modal mapping from audio waveform to gesture videos, even for the same speaker. To bridge the gap between audio and video, previous methods [24, 41] predict body pose (i.e., a jointed skeleton) as an intermediate low dimensional representation to drive the video synthesis. However, they dissect the problem into two independent modules (audio-to-pose, and pose-to-video) and produce results suffering from noticeable artifacts, e.g. distorted body parts and blurred appearance.

Our method introduces a video reenactment method that is able to synthesize high-resolution, high-quality speech gesture videos directly in the video domain by cutting, re-assembling, and blending clips from a single input reference video. The process is driven by a novel *video motion graph*, inspired by 3D motion graphs used in character animation [4, 35]. The graph nodes represent frames in the reference video, and edges encode possible transitions between them. We discover possible valid transitions between

frames, and also discover paths in the graph leading to the generation of a new video such that the re-enacted gestures are coherent and consistent with both the audio rhythms and speech content of the target audio.

Direct playback on the discovered paths for an output video can cause temporal inconsistency at the boundary of two disjoint raw frames. Existing frame blending methods cannot easily solve this problem, especially with fast moving and highly deformed human poses. Therefore, we also propose a novel human *pose-aware video blending* network to smoothly blend frames around the temporally inconsistent boundaries to produce naturally-looking video transitions. By doing so, we successfully transform the problem of audio-driven gesture reenactment into the search for valid paths that best match the given audio.

Our path discovery algorithm is motivated by psychological studies on co-speech gesture analysis. The studies show co-speech gestures can be categorized into rhythmic gestures and referential gestures [46]. While rhythmic gestures are well synchronized with audio onsets [9, 80], referential gestures mostly co-occur with certain phrases, e.g. a greeting gesture of hand-waving appears when a speaker says ‘hello’ or ‘hi’ [8, 15]. We analyze the speech of the reference video and detect the audio onset peaks [18] as well as a set of keywords from its transcript [76] as audio features added to the corresponding nodes on the video motion graph. Given the extracted audio onset peaks and keywords from a new audio clip, the optimal paths that best match audio features are used to drive our video synthesis.

Our contributions are summarized as follows:

- a new system that creates high-quality human speech videos with realistic gestures driven by audio only,
- a novel video motion graph that preserves the video realism and gesture subtleties,
- a pose-aware video blending neural network that synthesizes smooth transitions of two disjoint reference video clips along graph paths, and
- an audio-based search algorithm that drives the video synthesis to match the synthesized gesture frames with both the audio rhythms and the speech content.

2. Related Work

Our method is related to previous work on motion graphs, audio-driven 3D speech animation, and in particular human video synthesis, and video frame blending.

Motion Graph. The idea of motion graphs was first proposed in [4, 35] to create realistic and controllable animation based on a pre-captured motion. It is broadly used in generating 3D character animations [6, 26, 36, 38, 47, 52, 56, 59]. However, these approaches only work on 3D human skeleton representations and cannot be directly applied to video

animation in image space. While blending re-assembled motions requires interpolating 3D joint positions in character animation, in our case blending requires synthesizing whole image frames to create a coherent video.

[1, 57] propose motion graph in pixel space and solve this issue by de-ghosting [60] and gradient-domain compositing [68] based on pixel warping. However, these approaches focus on simple periodical scene scenarios, e.g. pendulum, waterfalls, etc. and cannot work on complex human motions. [23, 39, 77, 82] generate controllable human action videos by retrieving and warping nearest candidate frames. However, they require additional motion capture resources such as physical markers, multi-view or RGB-D cameras. [12, 13, 29] also introduce human video synthesis based on reconstruction of human meshes from pre-captured multi-view camera datasets. However, these methods are not suitable for monocular camera videos.

Audio-driven Speech Animation of 3D models. Several approaches for audio-driven speech animation of lips, heads, and body gestures have been proposed in the recent years [17, 24, 41, 67, 85, 87]. [2, 3, 37, 79] propose learning methods to solve the multimodal mapping from audio to 3D human gestures. They represent synthesized gestures with 3D skeletons, which can drive a 3D character model. Yet, these methods are not able to synthesize video of a target speaker unless they are also provided with a detailed, textured, and rigged 3D model for that speaker. When it is not available, their demonstrated results lack photorealism.

Human Video Synthesis. [24, 41] translate predicted skeletal gesture motions to photo-realistic speaker videos via recent neural image translation approaches [31, 34, 71, 72]. However, neural image translation is not artifact-free: disconnected moving object parts, as well as incoherent texture appearance are known issues in video generation [71]. Due to the large number of parameters in their network, these methods also require large datasets for training. Few-shot solutions [70, 81] do not have such dataset requirements, yet they suffer from various artifacts, in particular for human pose synthesis, such as blurred appearance and distorted body parts [70]. [42, 43, 61, 73, 74] fit human body model or/and texture parameters to a training video to improve the appearance of body shapes and texture at test time. Yet, inaccurate fitting easily results in artifacts and loss of subtleties, especially in the presence of loose clothing and detailed body parts, e.g. fingers. [62, 63, 84] warp learned features of each body to generate target pose frames based on estimated optical flow. They can handle large pose changes and texture hallucination, but often fail in blending frames naturally. Our method follows a largely different approach from all the above prior work: instead of per-frame neural translation, the video of a speaker is generated by re-assembling clips from a short, few minute long reference

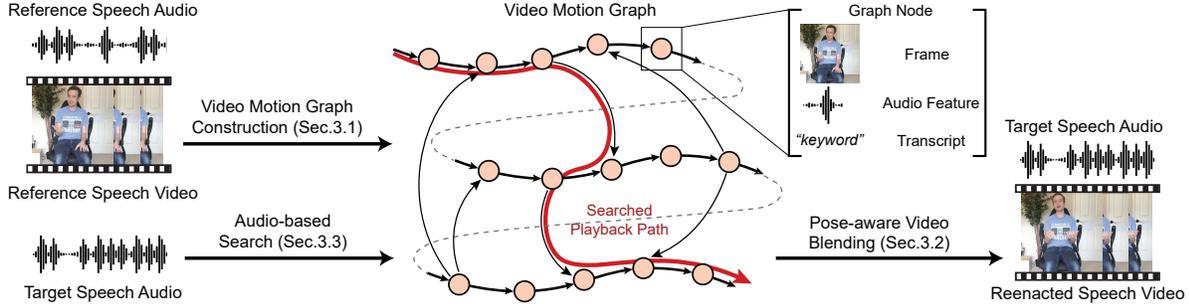


Figure 2. System overview. The reference video is first encoded into a directed graph where nodes represent video frames and audio features, and edges represent transitions. The transitions include original ones between consecutive reference frames, and synthetic ones between disjoint frames. Given a unseen target audio at test time, a beam search algorithm finds plausible playback paths such that gestures best match the target speech audio. Synthetic transitions along disjoint frames are neurally blended to achieve temporal consistency.

video. Since most of the frames originate from the reference video, the synthesized video preserves gesture realism as well as appearance subtleties. As a result, the problem is simplified to blending video frames. Our neural blending network focuses on solving this particular task, instead of generating all frames from scratch.

Video Frame Blending. The choice of the frame blending strategy significantly impacts the quality of the video generated from re-assembling clips. Naive weighted averaging of video frames easily result in ghost effects [49, 57]. More advanced frame interpolation methods [25, 33, 44, 50] based on optical flow estimation [5, 30, 66] have been proposed to synthesize intermediate frames between two consecutive frames, in particular for slow motion videos. However, such methods fail if two frames are very different from each other and the optical flow estimation is not accurate enough. They work for generic content, yet do not consider human motion as a prior for our task. Our method uses a human pose-aware neural network for frame blending that produces significantly better quality video compared to prior such work, as demonstrated in our experiments.

3. Method

Overview. The goal of our method is to synthesize a new video for a reference speaker given a target speech audio from the same or different speaker. Our video synthesis is guided by a novel *video motion graph* created from an input reference video of the speaker (Sec. 3.1). The video motion graph is a directed graph that encodes how the reference video may be split and re-assembled in different graph paths (see Fig. 2 for an illustration). The graph node representations are defined as the raw reference video frames and corresponding audio features. The edges are defined as the transitions between frames, including *natural transitions* in the input video and *synthetic transitions* connecting disjoint clips. Synthetic transitions are introduced to expand the graph connectivity and enable nonlinear video playback.

However, a direct nonlinear playback along synthetic transitions does not guarantee smooth video rendering due to the abrupt changes of disjoint frames in image space. Thus, we design a novel *pose-aware video blending* network to re-render and interpolate neighboring frames required by the synthetic transitions (Sec. 3.2). We develop an *audio-based searching* method to find optimal paths in the video motion graph that best match the target audio features both rhythmically and semantically (Sec. 3.3). To generate new videos, we retrieve the raw input video frames at natural transitions and synthesize neural blended frames at synthetic transitions.

3.1. Video Motion Graph

The key idea of our video motion graph is to create synthetic transitions based on the similarity of the speaker’s pose in the reference video frames. Our pose similarity metric relies on 3D space and image space cues. Given a reference video, our first step is to extract pose parameters θ of the SMPL model [45] for all frames with an off-the-shelf motion capture method [75]. We further smooth the pose parameters with [14] to promote temporally coherent results.

3D space pose similarity. Based on the pose parameters, we compute the 3D positions in world space for all joints via forward kinematics. For each pair of frames $\forall(m, n)$, we evaluate pose dissimilarity $d_{feat}(m, n)$ based on the Euclidean distance of their position and velocity of all joints.

Image space pose similarity. To obtain the pose similarity in image space, for each frame m , we project the fitted 3D SMPL human mesh onto image space using known camera parameters from [75], and mark the mesh surface area which is visible on image after projection as S_m . Then for each pair of frames (m, n) , the image space dissimilarity is estimated by the Intersection-over-Union (IoU) between their common visible surface areas: $d_{img}(m, n) = 1 - (S_m \cap S_n)/(S_m \cup S_n)$. The lower $d_{img}(m, n)$ is, the higher the IoU, thus larger overlap exists

in the surface area in two meshes, indicating higher pose similarity in terms of image rendering.

Based on these two distance measurements, we create graph synthetic transitions between any pair of reference video frames (nodes in our graph) if their distance $d_{feat}(m, n)$ and $d_{img}(m, n)$ are below predefined thresholds (both distance for natural transitions are defined as 0). Here we follow [78] to set the thresholds as the average distance between close frames ($m, m + l$) in the reference video. Larger frame offset l results in higher thresholds, thus more synthetic transitions, increasing the possible number of paths in the motion graph. This also results in larger computational cost for the path search algorithm of Sec. 3.3. Our experiments use $l = 4$ which practically achieves a balance between computational cost and number of available paths in the graph.

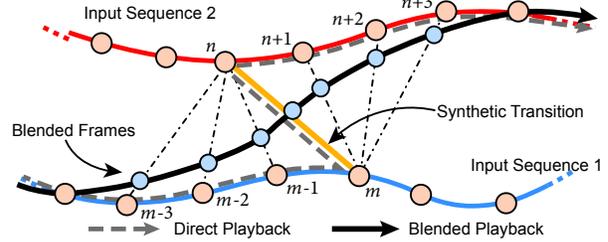
3.2. Pose-aware Video Blending

A mere playback of connected frames at synthetic transitions easily results in noticeable jittering artifacts (see direct playback in Fig. 3(a) grey dashed path and Fig. 3(b) third column). To solve this problem, we synthesize blended frames to replace original frames around a small temporal neighborhood of a synthetic transition so that the video can smoothly transit from the first sequence to the other (see Fig. 3(a) solid black path and Fig. 3(b) last column). For a synthetic transition connecting frames m, n , we define the neighborhood using the frame range $[m - k, m]$ and $[n, n + k]$ with a neighborhood size k .

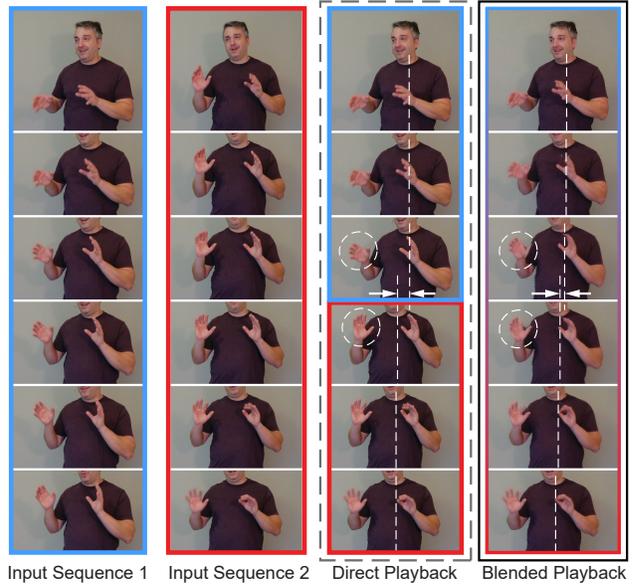
We designed a *pose-aware video blending network* to synthesize frames within the above neighborhood. Given two frames with indices i, j (where $i \in [m - k, m]$ and $j \in [n, n + k]$) and their corresponding raw RGB image representations I_i and $I_j \in \mathbb{R}^{H \times W \times 3}$ from the reference video, the network synthesizes each blended frame in the neighborhood with a target blended weight $\alpha \in [0, 1/K, 2/K, \dots, 1]$, where $K = 2k$.

As a first step, we use the blending weight to estimate the SMPL pose parameter θ_t for a blended frame t as: $\theta_t = (1 - \alpha)\theta_i + \alpha\theta_j$, where θ_i and θ_j are the SMPL pose parameters captured from two input frames respectively.

Our network processes the images I_i, I_j , the body foreground masks, and the pose parameter $\theta_i, \theta_j, \theta_t$. Processing takes place in two stages. The first stage warps foreground human body image features based on a 3D motion field computed from vertex displacements of the fitted SMPL meshes. The second stage further refines the warping by computing the residual optical flow between the warped image features produced by the first stage, and the optical flow from the rest of the image (i.e., background). Finally, an image translation network transforms the refined warped image features to the image I_t representing the target output frame t . The network architecture is shown in Fig. 4.



(a) Illustration of pose-aware blended playback.



(b) Our blended playback generates smoother transition.

Figure 3. Compared to direct playback along synthetic transitions which have severe horizontal shift for body poses and abrupt change in hand rotations (see dashed lines and circles in (b)), our blending strategy generates natural transitions between clips.

Mesh Flow Stage. The first stage has two parallel streams, each producing image deep feature maps encoding the warping for the input images I_i and I_j . To produce these features, we first compute an initial 3D motion field, which we refer to as initial “mesh flow”, from the SMPL body mesh displacements between the two frames. To this end, we first find the body mesh vertex positions $\mathbf{v}_i, \mathbf{v}_j, \mathbf{v}_t$ from the SMPL pose parameters $\theta_i, \theta_j, \theta_t$ respectively. Then we obtain the initial mesh flow $F_{t \rightarrow i}^{init}$ and $F_{t \rightarrow j}^{init}$ as the displacement of the corresponding mesh vertices $\mathbf{v}_t - \mathbf{v}_i$ and $\mathbf{v}_t - \mathbf{v}_j \in \mathbb{R}^{N \times 3}$ respectively. We note that we only consider here the displacements from visible vertices found via perspective projection onto image plane. These displacements are projected and rasterized as image-space motion field $\mathbb{R}^{N \times 3} \rightarrow \mathbb{R}^{H \times W \times 2}$. Since the vertex sampling does not match the image resolution, the resulting flow fields are rather sparse. Thus, we diffuse them with a Gaussian kernel with σ set to 8 in our experiments.

These initial motion fields are far from perfect. This is

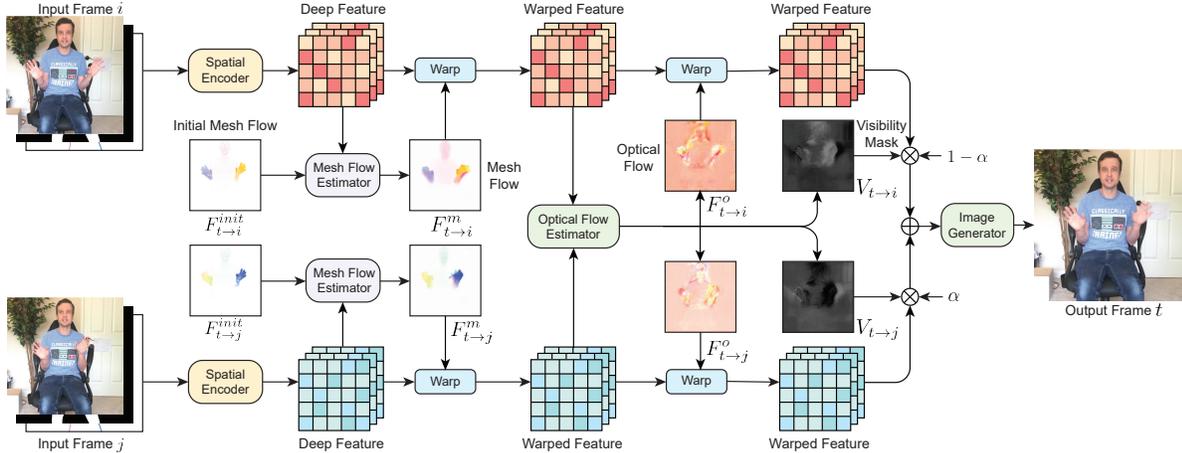


Figure 4. Pose-aware neural blending network architecture. Two source frames are encoded into deep feature maps and then warped based on the predicted flows from two stages: a 3D mesh-based flow stage for coarse feature map alignment, followed by an optical flow-based stage further refining the warping. Finally, the warped features are blended with predicted visibility masks to generate the target frame.

because the boundaries of the projected mesh often do not exactly align with the boundaries of the human body in the input frames. Thus, we refine these fields with a neural module. The module has two streams, each refining the corresponding motion field for frame i and j . The first stream processes as inputs the RGB image I_i , the foreground mask I_{mask} , and an image containing the rendered skeleton I_{skel} representing the SMPL pose parameters. It then encodes them into an image deep feature map \mathbf{x}_i :

$$\mathbf{x}_i = E_s(I_i, I_{mask}, I_{skel}; \mathbf{w}_s) \quad (1)$$

where \mathbf{w}_s are learnable weights. Similarly, the second stream produces an image deep feature map \mathbf{x}_j for frame j . The two streams share the same network based on 8 stacked CNN residual blocks [10]. More details are provided in the supplementary material.

We then estimate the refined motion fields through another network E_m ,

$$F_{t \rightarrow i}^m = E_m(\mathbf{x}_i, F_{t \rightarrow i}^{init}; \mathbf{w}_m), \quad (2)$$

$$F_{t \rightarrow j}^m = E_m(\mathbf{x}_j, F_{t \rightarrow j}^{init}; \mathbf{w}_m). \quad (3)$$

where \mathbf{w}_m are learnable weights. This network is designed based on UNet [53]. More details are provided in the supplementary material. We then backwards warp the above image feature maps with the above motion fields to obtain the warped deep features \mathbf{x}'_i and \mathbf{x}'_j .

Optical Flow Stage. Synthesizing the final target frame directly from the two warped feature maps \mathbf{x}'_i and \mathbf{x}'_j suffers from ghost effect (Fig. 5). This is because the motion field calculated in the previous stage is based on the SMPL model which ignores details such as textures on clothing.

Our second stage aims to further warp the deep feature maps \mathbf{x}'_i and \mathbf{x}'_j based on optical flow computed throughout the image including the background. At this stage, the

warped features already represent bodies that are roughly aligned. We found that an off-the-shelf frame interpolation network based on optical flow [33] can reproduce the missing pixel-level details and remedy the ghost effect. The network predicts optical flow $F_{t \rightarrow i}^o$ and $F_{t \rightarrow j}^o$ to further warp the features from \mathbf{x}'_i and \mathbf{x}'_j to \mathbf{x}''_i and \mathbf{x}''_j respectively. It also estimates soft visibility maps [33] $V_{t \rightarrow i}$ and $V_{t \rightarrow j}$ used for blending to obtain a deep feature map for frame t :

$$\mathbf{x}''_t = (1 - \alpha)V_{t \rightarrow i} \odot \mathbf{x}''_i + \alpha V_{t \rightarrow j} \odot \mathbf{x}''_j. \quad (4)$$

Finally, we take as input the above blended deep feature map to synthesize the target image I_t . This is performed with a generator network G following a UNet image translation network architecture [87]: $\hat{I}_t = G(\mathbf{x}''_t; \mathbf{w}_g)$, where \mathbf{w}_g are learnable weights. More details and output examples are provided in the supplementary material.

Training. To train our pose-aware video blending network, we sample triplets of frames in the reference video. Given a target frame e.g., frame t , we randomly sample two other nearby frames with indices $t - k_0$ and $t + k_1$, $k_0, k_1 \in [1, 8]$ to form triplets. The corresponding blending weight α is computed as $k_0 / (k_0 + k_1)$. We train the entire network end-to-end with losses defined to better estimate the flows and reconstruct the final image. More details are provided in the supplementary material.

3.3. Audio-based Search

Given a speech audio at test time, we develop a graph search algorithm to find plausible paths along which gestures match the speech audio both rhythmically and semantically. Previous studies have shown that speech gestures can be classified into two categories: 1) referential gestures that appear together with specific, meaningful keywords, and 2) rhythmic gestures which respond to the audio



Figure 5. A ghost effect example. **Left:** two input frames. **Top-right:** ghost effect from using mesh flow only. **Bottom-right:** sharp features with further warping by optical flow.

prosody features [46]. More specifically, the key stroke of a rhythmic gesture appear at the same time as (or within a very short of period of) an *audio onset* within a phonemic clause [21]. To find precise gestures on the right timings, or frame indices, we define a pair of audio features for input speech: *audio onset feature* and *keyword feature*. The audio frame indices match the video frame rate.

Audio onset and keyword feature. We define the audio onset feature as a binary value indicating the activation of an audio onset for each frame detected with a standard audio processing algorithm [7]. To extract keyword features, we first use the Microsoft Azure speech-to-text engine [76] to convert the input audio into transcripts with corresponding start and end time for each word. We create a dictionary of common words for referential gestures, which we call *keywords* (see supplementary for a list). If a keyword appears at a frame (or node), we set its keyword feature to that word. Otherwise, we simply set it to *empty* (no keyword).

Target speech audio segmentation. We split the target speech audio into segments starting and ending with the frames where the audio onset or keyword feature is activated. Let $\{a_s\}_{s=1}^S$ be the frame indices of such frames, where S is their total number. Segments are represented as $a_s \rightarrow a_{s+1}$, and their duration are $L_s = a_{s+1} - a_s$ (number of frames). We also add two extra endpoints $a_0 = 1$ and $a_{S+1} = N_t$ indicating the first and last frame of the target audio respectively to form the complete segment list, i.e. $a_s \rightarrow a_{s+1}$, $s = 0, 1, \dots, S$.

Beam search. We utilize the beam search [55] in the video motion graph to find K plausible paths matching the target speech audio segments. The beam search initializes K paths starting with K random nodes as the first frame a_0 for the target audio. Next, we apply breadth-first search to find path segments ending on nodes whose feature matches

the target audio segment feature at frame a_1 . We continue with the same search procedure as above to find full graph paths matching the rest of the target segments $a_s \rightarrow a_{s+1}$, $s = 1, \dots, S$ iteratively. All searched K paths can be used to generate various plausible results for the same target speech audio. Detailed search criteria and result variants can be found in the supplementary material and our project page.

Video synthesis. We generate a video along with the final path in the motion graph discovered by the beam search executions, and use the blending network to handle synthetic transitions (see Fig. 3 for an example). As explained above, for each synthesized video segment corresponding to target audio segment $a_s \rightarrow a_{s+1}$, we adjust its speed to match the target duration. Finally, we post-process our result by adopting [51] to synchronize the lips of the speaker to match the corresponding speech audio.

4. Results and Evaluation

Dataset. We evaluate our neural blending network and produced audio-driven reenactment results on two datasets.

Personal Story Dataset. Since our approach works for speaker-specific speech gesture reenactment, we collected seven speech videos. Each speaker is asked to tell a personal story in front of a static camera, either standing or sitting. Speakers are encouraged to use their gestures while telling the stories. The length of the video varies between 2-10 minutes depending on the story. We split each video into 90%/10% for training and testing purpose.

TED-talks dataset [63]. We also demonstrate the generalization of our neural blending network on the TED-talks dataset. It contains 1265 talk speech videos with 393 unique speakers. Each video contains the upper part of speaker body and the video length ranges from 2 to 60 seconds. We use the same train/test split proposed in [63]. We evaluated the generalization ability of our model on this dataset since the test speakers are unseen during training.

4.1. Video Blending Evaluation

We firstly numerically evaluate the proposed video blending network on both dataset. Given two frames $t - k$ and $t + k$ in the test split of each video, we synthesize blended frames with the blending weight $\alpha = 0.5$ and compare its quality with the ground-truth frame t . All the compared frames are multiplied with ground-truth human masks to compare the foreground human results only.

We compare our method with the state-of-the-art frame interpolation methods **FeatureFlow** [25] and **SuperSIMO** [33], as well as human pose-based image synthesis methods **vUnet** [22]. We also compare with methods based on the pix2pix [72] backbone: the **EBDance** [16] method for speaker-specific Personal story dataset and the **Fewshot-vid2vid** [70] for the speaker-varying TED-talks dataset. For

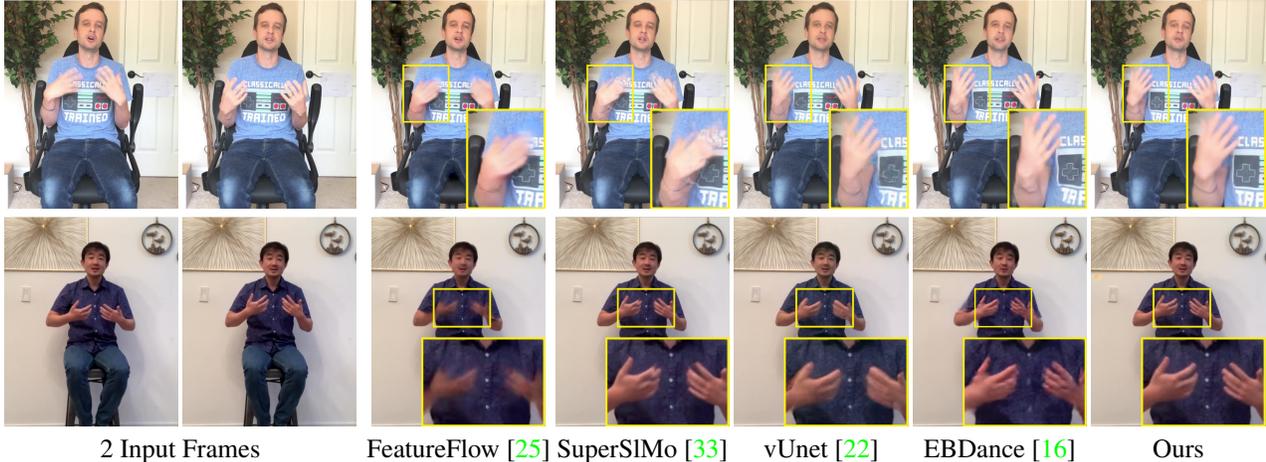


Figure 6. Comparison of blended frame synthesis using different methods. Note the natural look of details such as fingers in our method.

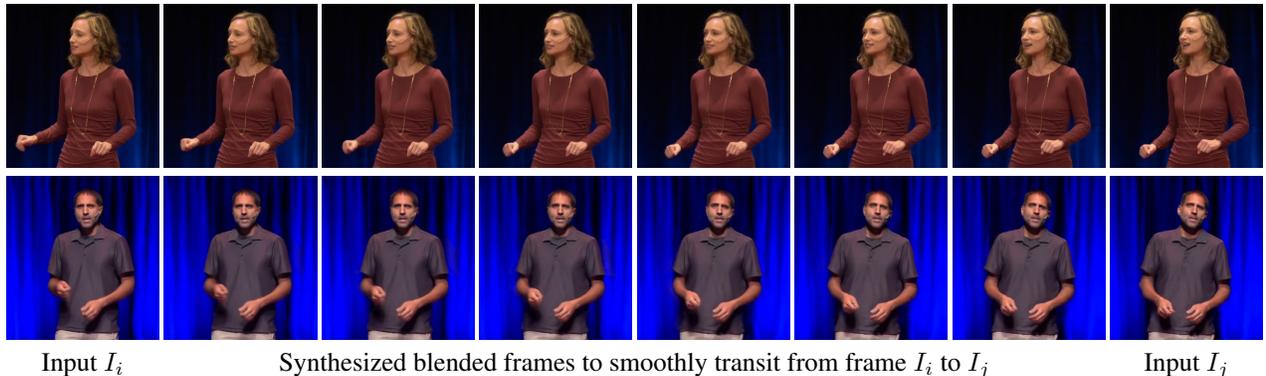


Figure 7. Blended frames for transition edges of our video motion graph on TED-talks dataset (see demo videos on our project page).

the pose-based image synthesis methods, we interpolate human skeleton by averaging joint positions. We retrain all the comparison methods on our dataset for a fair comparison. We also evaluate two network alternatives: **Ours w/ mesh** which only uses mesh-based warping flows and **Ours w/ optical** which only uses optical flows.

Image Quality. We evaluate the quality of synthesized images via four common metrics: Image Error (IE) - average absolute pixel difference between two images; Peak Signal-to-Noise Ratio (PSNR) and LPIPS [83].

Table 1 shows our model consistently outperforms all other methods for speaker-specific videos on the Personal story dataset. It also demonstrates the generalization of our model for unseen speakers on the TED-talks dataset. Fig. 6 shows examples of synthesized frames by different methods. In the top example, the inputs are two frames with larger gesture difference. The frame interpolation methods [25, 33] cannot estimate the flow field, and thus result in broken and blurred hand results. The pose-based image synthesis methods [16, 22] preserve hand structures but have artifacts around fingers and clothing. Ours achieves the best quality for both hands and clothing. The lower

example shows frames with smaller gesture differences. [16, 22, 33] preserve hands better but still suffer from broken and blurred texture. Ours generates clear and sharp results.

Video Quality. To evaluate the quality of the generated video, we adopt the metric, MOVIE [58] index, to evaluate the video distortion in spatio-temporal aspects. We also follow [71] to evaluate the visual quality of the video and temporal consistency with Fréchet Inception Distance (FID) scores [27]. We use the pre-trained video recognition CNN model to get features from synthesized video clips [11]. Table 1 relative columns show our method can achieve the best video quality in the temporal domain. It demonstrates that the synthesized blended frames seamlessly connect reenacted frames with much less temporal artifacts. In Fig. 7, we show detailed blended frames on the selected transition edges of our video motion graph from the TED-talks dataset. We provide additional synthesized clips to show-case blending results on our project page.

4.2. Audio-driven Reenactment Results

Given a reference video from speaker A and a target audio clip randomly from another speaker B , we can reenact

Method	Personal story dataset					TED-talks dataset				
	IE↓	PSNR↑	LPIPS↓	MOVIE↓	FID↓	IE↓	PSNR↑	LPIPS↓	MOVIE↓	FID↓
FeatureFlow [25]	1.18	33.5	0.015	0.22	19.1	5.2	19.7	0.267	1.29	33.6
SuperSliMo [33]	1.04	35.0	0.012	0.17	15.4	1.18	28.6	0.052	0.50	12.6
vUnet [22]	1.20	33.6	0.013	0.19	15.6	1.19	28.8	0.058	0.52	14.0
EBDance [16]	1.75	30.7	0.020	0.43	20.5	-	-	-	-	-
Fewshot-vid2vid [70]	-	-	-	-	-	10.7	15.1	0.159	1.06	21.5
Ours w/ mesh	0.87	35.2	0.009	0.14	15.1	1.36	27.9	0.072	0.64	11.5
Ours w/ optical	0.97	34.6	0.009	0.16	13.2	1.25	28.2	0.069	0.57	11.9
Ours	0.76	36.1	0.007	0.13	13.0	0.93	30.7	0.040	0.43	11.8

Table 1. Image and video quality assessment for Personal story dataset and TED-talks dataset.

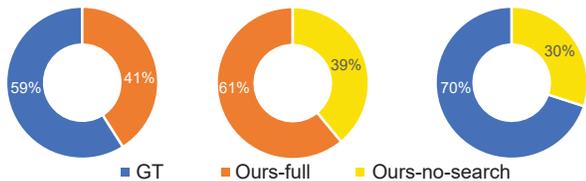


Figure 8. Pairwise comparison results from our user study. The comparison of *Ours-full* against *Ours-no-search* shows the effectiveness of the proposed audio-based search algorithm.

the reference video to generate a new speech video with *A*'s appearance and *B*'s voice based on our pipeline. The reenacted results on both the Personal story dataset and the TED-talks dataset are provided on the project page.

User Study. To further quantitatively evaluate the consistency of such reenacted videos to target speech audio, we perform a perceptual user study on the reenacted videos on the Personal story dataset. We generate 127 such videos of 25 seconds in length. Each of them contains expressive speech gestures for every speaker in the dataset. The study was conducted via the Amazon Mechanical Turk service. We compare the results from our full system (**Ours-full**) against ground-truth (**GT**), which are original reference video clips of speaker *B*, and results from a baseline system (**Ours-no-search**), which randomly finds paths along video motion graphs without audio-based search.

We design the user study questionnaire by providing a list of queries involving pairwise comparisons of results from two out of three methods mentioned above. The participants were asked to choose which gestures in those two results are more consistent with the speech audio. Detailed setup to prevent biases and invalid answers can be found in the supplementary material. Finally, 1130 valid choices from 113 valid participants are gathered. We plot the statistics in Fig. 8. The preference (61% vs. 39%) of *Ours-full* over *Ours-no-search* shows the effectiveness of the audio-based search algorithm. Although no audio guidance is used, 30% votes received by *Ours-no-search* against *GT* also suggest our video motion graph and frame blending approach is able to generate high-quality and realistic videos.

The relative higher votes (41%) given to *Ours-full* against *GT* demonstrates our full system generates better though not perfect gesture videos that are coherent with the audio.

5. Conclusion and Future Work

We propose a novel system based on video motion graphs to generate new videos that best preserve high image synthesis quality and speaker gesture motion subtleties. To seamlessly reenact disjoint frames from the input video, we introduce a neural pose-aware video blending method to smoothly blend inconsistent transition frames. We show the superior performance of the proposed system comparing to the state-of-the-art methods and baselines via both numerical experiments and perceptual user studies.

Limitations. We use a pre-defined common keyword dictionary for keyword features, which may fail on uncommon individual vocabulary. Using richer audio features learnt through data might help with accurate gesture matching. There is an inevitable trade-off between the quality and variety of synthesized animations: increasing the graph edge density can increase the transition variety, yet may retrieve frames harder to blend. The proposed video blending network can blend the foreground human poses and slight background changes, but it fails on dramatically changed backgrounds (see the supplementary for examples).

Future work. Neural blending shows its strengths on reenacting human videos in the pose-aware embedding space. We believe our hybrid framework of video motion graph and neural reenactment is a promising direction for high-quality controllable digital human animations.

Potential negative societal impacts. Our approach enables synthesis of talking people. This offers the ability of creating fake videos for malicious purposes. Detecting deep fakes videos [40, 54, 86] is an active area of research.

Acknowledgements. Our research was partially funded by NSF (EAGER-1942069) and Adobe.

References

- [1] Aseem Agarwala, Ke Colin Zheng, Chris Pal, Maneesh Agrawala, Michael Cohen, Brian Curless, David Salesin, and Richard Szeliski. Panoramic video textures. In *ACM Trans. on Graphics (TOG)*, 2005. 2
- [2] Chaitanya Ahuja, Dong Won Lee, Yukiko I Nakano, and Louis-Philippe Morency. Style transfer for co-speech gesture animation: A multi-speaker conditional-mixture approach. In *Proc. ECCV*, 2020. 2
- [3] Simon Alexanderson, Gustav Eje Henter, Taras Kucherenko, and Jonas Beskow. Style-controllable speech-driven gesture synthesis using normalising flows. In *Computer Graphics Forum*, 2020. 2
- [4] Okan Arıkan and David A Forsyth. Interactive motion generation from examples. *ACM Trans. on Graphics (TOG)*, 2002. 1, 2
- [5] Simon Baker, Daniel Scharstein, JP Lewis, Stefan Roth, Michael J Black, and Richard Szeliski. A database and evaluation methodology for optical flow. *IJCV*, 2011. 3
- [6] Philippe Beaudoin, Stelian Coros, Michiel van de Panne, and Pierre Poulin. Motion-motif graphs. In *Proc. ACM SCA*, 2008. 2
- [7] Juan Pablo Bello, Laurent Daudet, Samer Abdallah, Chris Duxbury, Mike Davies, and Mark B Sandler. A tutorial on onset detection in music signals. *IEEE Trans on Speech and Audio Processing*, 2005. 6
- [8] Kirsten Bergmann and Stefan Kopp. Gnetic—using bayesian decision networks for iconic gesture generation. In *International Workshop on Intelligent Virtual Agents*, 2009. 2
- [9] Elif Bozkurt, Yücel Yemez, and Engin Erzin. Multimodal analysis of speech and arm motion for prosody-driven synthesis of beat gestures. *Speech Communication*, 2016. 2
- [10] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale gan training for high fidelity natural image synthesis. In *Proc. ICLR*, 2018. 5
- [11] Joao Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset. In *Proc. CVPR*, 2017. 7
- [12] Dan Casas, Christian Richardt, John Collomosse, Christian Theobalt, and Adrian Hilton. 4d model flow: Precomputed appearance alignment for real-time 4d video interpolation. In *Computer Graphics Forum*, 2015. 2
- [13] Dan Casas, Marco Volino, John Collomosse, and Adrian Hilton. 4d video textures for interactive character appearance. In *Computer Graphics Forum*, 2014. 2
- [14] Géry Casiez, Nicolas Roussel, and Daniel Vogel. 1 € filter: A simple speed-based low-pass filter for noisy input in interactive systems. In *Proc. SIGCHI on Human Factors in Computing Systems*, 2012. 3
- [15] Justine Cassell, Matthew Stone, and Hao Yan. Coordination and context-dependence in the generation of embodied conversation. In *Proc. International Conference on Natural Language Generation*, 2000. 2
- [16] Caroline Chan, Shiry Ginosar, Tinghui Zhou, and Alexei A Efros. Everybody dance now. In *Proc. ICCV*, 2019. 6, 7, 8
- [17] Lele Chen, Guofeng Cui, Celong Liu, Zhong Li, Ziyi Kou, Yi Xu, and Chenliang Xu. Talking-head generation with rhythmic head motion. In *Proc. ECCV*, 2020. 1, 2
- [18] Abe Davis and Maneesh Agrawala. Visual rhythm and beat. *ACM Trans. on Graphics (TOG)*, 2018. 2
- [19] James E Driskell and Paul H Radtke. The effect of gesture on speech production and comprehension. *Human factors*, 2003. 1
- [20] Pif Edwards, Chris Landreth, Eugene Fiume, and Karan Singh. Jali: an animator-centric viseme model for expressive lip synchronization. *ACM Trans. on Graphics (TOG)*, 2016. 1
- [21] Mohamed Elhoseiny, Scott Cohen, Walter Chang, Brian Price, and Ahmed Elgammal. Sherlock: Scalable fact learning in images. In *Proc. AAAI*, 2017. 6
- [22] Patrick Esser, Ekaterina Sutter, and Björn Ommer. A variational u-net for conditional appearance and shape generation. In *Proc. CVPR*, 2018. 6, 7, 8
- [23] Matthew Flagg, Atsushi Nakazawa, Qiushuang Zhang, Sing Bing Kang, Young Kee Ryu, Irfan Essa, and James M Rehg. Human video textures. In *Proc. Symposium on Interactive 3D Graphics and Games*, 2009. 2
- [24] Shiry Ginosar, Amir Bar, Gefen Kohavi, Caroline Chan, Andrew Owens, and Jitendra Malik. Learning individual styles of conversational gesture. In *Proc. CVPR*, 2019. 1, 2
- [25] Shurui Gui, Chaoyue Wang, Qihua Chen, and Dacheng Tao. Featureflow: robust video interpolation via structure-to-texture generation. In *Proc. CVPR*, 2020. 3, 6, 7, 8
- [26] Rachel Heck and Michael Gleicher. Parametric motion graphs. In *Proc. Symposium on Interactive 3D Graphics and Games*, 2007. 2
- [27] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *Proc. NeurIPS*, 2017. 7
- [28] Chien-Ming Huang and Bilge Mutlu. Modeling and evaluating narrative gestures for humanlike robots. In *Robotics: Science and Systems*, 2013. 13
- [29] Peng Huang, Margara Tejera, John Collomosse, and Adrian Hilton. Hybrid skeletal-surface motion graphs for character animation from 4d performance capture. *ACM Trans. on Graphics (TOG)*, 2015. 2
- [30] Eddy Ilg, Nikolaus Mayer, Tonmoy Saikia, Margret Keuper, Alexey Dosovitskiy, and Thomas Brox. Flownet 2.0: Evolution of optical flow estimation with deep networks. In *Proc. CVPR*, 2017. 3
- [31] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. In *Proc. CVPR*, 2017. 2
- [32] Jana M Iverson and Susan Goldin-Meadow. Why people gesture when they speak. *Nature*, 1998. 1
- [33] Huaizu Jiang, Deqing Sun, Varun Jampani, Ming-Hsuan Yang, Erik Learned-Miller, and Jan Kautz. Super slo-mo: High quality estimation of multiple intermediate frames for video interpolation. In *Proc. CVPR*, 2018. 3, 5, 6, 7, 8, 12
- [34] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks.

- In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 2
- [35] Lucas Kovar, Michael Gleicher, and Frédéric Pighin. Motion graphs. In *ACM Trans. on Graphics (TOG)*. 2008. 1, 2
- [36] Björn Krüger, Jochen Tautges, Andreas Weber, and Arno Zinke. Fast local and global similarity searches in large motion capture databases. In *Proc. ACM SCA*, 2010. 2
- [37] Taras Kucherenko, Patrik Jonell, Sanne van Waveren, Gustav Eje Henter, Simon Alexandersson, Iolanda Leite, and Hedvig Kjellström. Gesticulator: A framework for semantically-aware speech-driven gesture generation. In *Proc. ICMI*, 2020. 2
- [38] Yongjoon Lee, Kevin Wampler, Gilbert Bernstein, Jovan Popović, and Zoran Popović. Motion fields for interactive character locomotion. In *ACM Trans. on Graphics (TOG)*. 2010. 2
- [39] Kun Li, Jingyu Yang, Leijie Liu, Ronan Boulic, Yu-Kun Lai, Yebin Liu, Yubin Li, and Eray Molla. Spa: Sparse photorealistic animation using a single rgb-d camera. *IEEE Trans. on CSVT*, 2016. 2
- [40] Yuezun Li, Xin Yang, Pu Sun, Honggang Qi, and Siwei Lyu. Celeb-df: A large-scale challenging dataset for deepfake forensics. In *Proc. CVPR*, 2020. 8
- [41] Miao Liao, Sibao Zhang, Peng Wang, Hao Zhu, Xinxin Zuo, and Ruigang Yang. Speech2video synthesis with 3d skeleton regularization and expressive body poses. In *Proc. ACCV*, 2020. 1, 2
- [42] Lingjie Liu, Weipeng Xu, Michael Zollhoefer, Hyeonwoo Kim, Florian Bernard, Marc Habermann, Wenping Wang, and Christian Theobalt. Neural rendering and reenactment of human actor videos. *ACM Trans on Graphics (TOG)*, 2019. 2
- [43] Wen Liu, Zhixin Piao, Jie Min, Wenhan Luo, Lin Ma, and Shenghua Gao. Liquid warping gan: A unified framework for human motion imitation, appearance transfer and novel view synthesis. In *Proc. ICCV*, 2019. 2
- [44] Ziwei Liu, Raymond A Yeh, Xiaoou Tang, Yiming Liu, and Aseem Agarwala. Video frame synthesis using deep voxel flow. In *Proc. ICCV*, 2017. 3
- [45] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J Black. Smpl: A skinned multi-person linear model. *ACM Trans. on Graphics (TOG)*, 2015. 3
- [46] David McNeill. *Hand and mind: What gestures reveal about thought*. University of Chicago press, 1992. 1, 2, 6, 13
- [47] Jianyuan Min and Jinxiang Chai. Motion graphs++ a compact generative model for semantic motion analysis and synthesis. *ACM Trans. on Graphics (TOG)*, 2012. 2
- [48] Lucie Naert, Caroline Larboulette, and Sylvie Gibet. A survey on the animation of signing avatars: From sign representation to utterance synthesis. *Computers & Graphics*, 2020. 1
- [49] Simon Niklaus and Feng Liu. Softmax splatting for video frame interpolation. In *Proc. CVPR*, 2020. 3
- [50] Simon Niklaus, Long Mai, and Feng Liu. Video frame interpolation via adaptive separable convolution. In *Proc. ICCV*, 2017. 3
- [51] KR Prajwal, Rudrabha Mukhopadhyay, Vinay P Namboodiri, and CV Jawahar. A lip sync expert is all you need for speech to lip generation in the wild. In *Proc. ACM International Conference on Multimedia*, 2020. 6
- [52] Paul SA Reitsma and Nancy S Pollard. Evaluating motion graphs for character animation. *ACM Trans. on Graphics (TOG)*, 2007. 2
- [53] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, 2015. 5
- [54] Andreas Rossler, Davide Cozzolino, Luisa Verdoliva, Christian Riess, Justus Thies, and Matthias Nießner. Faceforensics++: Learning to detect manipulated facial images. In *Proc. ICCV*, 2019. 8
- [55] Steven M Rubin and Raj Reddy. The locus model of search and its use in image interpretation. In *IJCAI*, 1977. 6, 13
- [56] Alla Safonova and Jessica K Hodgins. Construction and optimal search of interpolated motion graphs. *ACM Trans. on Graphics (TOG)*, 2007. 2
- [57] Arno Schödl, Richard Szeliski, David H Salesin, and Irfan Essa. Video textures. In *Proc. Conference on Computer Graphics and Interactive Techniques*, 2000. 2, 3
- [58] Kalpana Seshadrinathan and Alan Conrad Bovik. Motion tuned spatio-temporal quality assessment of natural videos. *IEEE Trans. Image Processing*, 2009. 7
- [59] Hyun Joon Shin and Hyun Seok Oh. Fat graphs: constructing an interactive character with continuous controls. In *Proc. ACM SCA*, 2006. 2
- [60] H-Y Shum and Richard Szeliski. Construction of panoramic image mosaics with global and local alignment. In *Panoramic vision*. 2001. 2
- [61] Aliaksandra Shysheya, Egor Zakharov, Kara-Ali Aliev, Renat Bashirov, Egor Burkov, Karim Isakov, Aleksei Ivakhnenko, Yury Malkov, Igor Pasechnik, Dmitry Ulyanov, et al. Textured neural avatars. In *Proc. CVPR*, 2019. 2
- [62] Aliaksandr Siarohin, Enver Sangineto, Stéphane Lathuilière, and Nicu Sebe. Deformable gans for pose-based human image generation. In *Proc. CVPR*, 2018. 2
- [63] Aliaksandr Siarohin, Oliver J Woodford, Jian Ren, Menglei Chai, and Sergey Tulyakov. Motion representations for articulated animation. In *Proc. CVPR*, 2021. 2, 6, 12, 15
- [64] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *ICLR*, 2015. 12
- [65] Sarah Taylor, Taehwan Kim, Yisong Yue, Moshe Mahler, James Krahe, Anastasio Garcia Rodriguez, Jessica Hodgins, and Iain Matthews. A deep learning approach for generalized speech animation. *ACM Trans. on Graphics (TOG)*, 2017. 1
- [66] Zachary Teed and Jia Deng. Raft: Recurrent all-pairs field transforms for optical flow. In *Proc. ECCV*, 2020. 3
- [67] Justus Thies, Mohamed Elgharib, Ayush Tewari, Christian Theobalt, and Matthias Nießner. Neural voice puppetry: Audio-driven facial reenactment. In *Proc. ECCV*, 2020. 1, 2
- [68] Hongcheng Wang, Ramesh Raskar, and Narendra Ahuja. Seamless video editing. In *Proc. ICPR*, 2004. 2

- [69] Kaisiyuan Wang, Qianyi Wu, Linsen Song, Zhuoqian Yang, Wayne Wu, Chen Qian, Ran He, Yu Qiao, and Chen Change Loy. Mead: A large-scale audio-visual dataset for emotional talking-face generation. In *Proc. ECCV*, 2020. 1
- [70] Ting-Chun Wang, Ming-Yu Liu, Andrew Tao, Guilin Liu, Bryan Catanzaro, and Jan Kautz. Few-shot video-to-video synthesis. In *Proc. NeurIPS*, 2019. 2, 6, 8
- [71] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Guilin Liu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. Video-to-video synthesis. In *Proc. NeurIPS*, 2018. 2, 7
- [72] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. High-resolution image synthesis and semantic manipulation with conditional gans. In *Proc. CVPR*, 2018. 2, 6
- [73] Chung-Yi Weng, Brian Curless, and Ira Kemelmacher-Shlizerman. Photo wake-up: 3d character animation from a single photo. In *Proc. CVPR*, 2019. 2
- [74] Chung-Yi Weng, Brian Curless, and Ira Kemelmacher-Shlizerman. Vid2actor: Free-viewpoint animatable person synthesis from video in the wild. *arXiv preprint arXiv:2012.12884*, 2020. 2
- [75] Donglai Xiang, Hanbyul Joo, and Yaser Sheikh. Monocular total capture: Posing face, body, and hands in the wild. In *Proc. CVPR*, 2019. 3
- [76] Wayne Xiong, Lingfeng Wu, Fil Allewa, Jasha Droppo, Xuedong Huang, and Andreas Stolcke. The microsoft 2017 conversational speech recognition system. In *Proc. ICASSP*, 2018. 2, 6
- [77] Feng Xu, Yebin Liu, Carsten Stoll, James Tompkin, Gaurav Bharaj, Qionghai Dai, Hans-Peter Seidel, Jan Kautz, and Christian Theobalt. Video-based characters: creating new human performances from a multi-view video database. In *ACM Trans. on Graphics (TOG)*. 2011. 2
- [78] Yanzhe Yang, Jimei Yang, and Jessica Hodgins. Statistics-based motion synthesis for social conversations. In *Computer Graphics Forum*, 2020. 4
- [79] Youngwoo Yoon, Bok Cha, Joo-Haeng Lee, Minsu Jang, Jaeyeon Lee, Jaehong Kim, and Geehyuk Lee. Speech gesture generation from the trimodal context of text, audio, and speaker identity. *ACM Trans. on Graphics (TOG)*, 2020. 2, 13
- [80] Fajriyan Yunus, Chloé Clavel, and Catherine Pelachaud. Sequence-to-sequence predictive models: from prosody to communicative gestures. In *Workshop sur les Affects, Compagnons artificiels et Interactions*, 2020. 2
- [81] Egor Zakharov, Aliaksandra Shysheya, Egor Burkov, and Victor Lempitsky. Few-shot adversarial learning of realistic neural talking head models. In *Proc. ICCV*, 2019. 2
- [82] Haotian Zhang, Cristobal Sciotto, Maneesh Agrawala, and Kayvon Fatahalian. Vid2player: Controllable video sprites that behave and appear like professional tennis players. *arXiv preprint arXiv:2008.04524*, 2020. 2
- [83] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proc. CVPR*, 2018. 7
- [84] Haitian Zheng, Lele Chen, Chenliang Xu, and Jiebo Luo. Unsupervised pose flow learning for pose guided synthesis. *arXiv*, 2019. 2
- [85] Hang Zhou, Yu Liu, Ziwei Liu, Ping Luo, and Xiaogang Wang. Talking face generation by adversarially disentangled audio-visual representation. In *Proc. AAAI*, 2019. 1, 2
- [86] Tianfei Zhou, Wenguan Wang, Zhiyuan Liang, and Jianbing Shen. Face forensics in the wild. In *Proc. CVPR*, 2021. 8
- [87] Yang Zhou, Xintong Han, Eli Shechtman, Jose Echevarria, Evangelos Kalogerakis, and Dingzeyu Li. Makeltalk: speaker-aware talking-head animation. *ACM Trans. on Graphics (TOG)*, 2020. 1, 2, 5
- [88] Yang Zhou, Zhan Xu, Chris Landreth, Evangelos Kalogerakis, Subhransu Maji, and Karan Singh. Visemenet: Audio-driven animator-centric speech animation. *ACM Trans. on Graphics (TOG)*, 2018. 1

Supplementary Materials

Implementation and demo videos. Our implementation and demo videos can be found at our project page https://yzhou359.github.io/video_reenact.

Training details. We train the entire network end-to-end with losses promoting better flow estimation and final frame reconstruction. Specifically, we first have an L1 reconstruction loss L_{rec} and a perceptual loss L_{per} between the synthesized image \hat{I}_t and I_t :

$$L_{rec} = \mathcal{L}_1(I_t, \hat{I}_t) \quad (5)$$

$$L_{per} = \mathcal{L}_1(\phi(I_t), \phi(\hat{I}_t)) \quad (6)$$

where $\phi(\cdot)$ concatenates feature map activations from a pre-trained VGG19 network [64].

We then adopt another L1 reconstruction loss \mathcal{L}_{rec}^b promoting better frame reconstruction directly from the warped deep features x''_i and x''_j after these pass through our generator network G . This helped predict warped deep features such that they lead to generating frames as close as possible to ground-truth in the first place. We also empirically observed faster convergence with this loss:

$$L_{rec}^b = \mathcal{L}_1(I_t, G(x''_i)) + \mathcal{L}_1(I_t, G(x''_j)) \quad (7)$$

Further, we have warping loss L_{warp}^m and L_{warp}^o by measuring the L1 reconstruction error between the target image and the source images I_i and I_j after being warped through the motion field $F_{t \rightarrow i}^m$ (Equations 2 and 3 in the main paper) and also the optical flow $F_{t \rightarrow i}^o$:

$$L_{warp}^m = \mathcal{L}_1(I_t, \mathcal{W}(I_i, F_{t \rightarrow i}^m)) + \mathcal{L}_1(I_t, \mathcal{W}(I_j, F_{t \rightarrow j}^m)) \quad (8)$$

$$L_{warp}^o = \mathcal{L}_1(I_t, \mathcal{W}(\mathcal{W}(I_i, F_{t \rightarrow i}^m), F_{t \rightarrow i}^o)) + \mathcal{L}_1(I_t, \mathcal{W}(\mathcal{W}(I_j, F_{t \rightarrow j}^m), F_{t \rightarrow j}^o)) \quad (9)$$

where $\mathcal{W}(I, F)$ applies backward warping flow F on image I .

Finally, we follow [33] and include a smoothness loss for both mesh flow and optical flow:

$$L_{sm} = \|\nabla F_{t \rightarrow i}^m\|_1 + \|\nabla F_{t \rightarrow j}^m\|_1 + \quad (10)$$

$$\|\nabla F_{t \rightarrow i}^o\|_1 + \|\nabla F_{t \rightarrow j}^o\|_1 \quad (11)$$

The overall loss \mathcal{L} is defined as the weighted sum of all losses described above, then averaged over all training frames.

Category	Keywords
greeting	hey, hi, hello
counting	one, two, three, first, second, third
direction	east, west, north, south, back, front, away, here, around
sentiment	crazy, incredible, surprising, screaming
action	walk, drive, ride, enter, open, attach, take, move
relative	more, less, much, few
others	called

Table 2. Dictionary of common keywords.

$$\mathcal{L} = \mathcal{L}_{rec} + \lambda_p \mathcal{L}_{per} + \lambda_b \mathcal{L}_{rec}^b + \lambda_m \mathcal{L}_{warp}^m + \lambda_o \mathcal{L}_{warp}^o + \lambda_s \mathcal{L}_{sm} \quad (12)$$

The weights have been set empirically based on [33] as $\lambda_p = 0.01$, $\lambda_b = 0.25$, $\lambda_m = 0.25$, $\lambda_o = 0.25$, $\lambda_s = 0.01$.

To train the entire model, we first train the mesh flow estimator network with L_{warp}^m as a ‘‘warming’’ stage. Then we load a pre-trained optical flow model from [33]. Finally, we train the entire network end-to-end with the loss mentioned above. The network weights are optimized with Adam optimizer using PyTorch. The learning rate is set to 10^{-4} and weight decay to 10^{-6} . The training process is performed on 4 Nvidia GeForce 1080Ti GPUs.

We show the detailed training procedure for our numerical evaluation. For the Personal story dataset, we train the model on each individual speaker video and report the evaluation numbers accordingly. The compared methods are trained on each speaker video for a fair comparison. For the TED-talks dataset, we train a single model on the entire training split of the dataset. We evaluate our model generalization on the testing split which contains unseen speakers. The comparison methods are also trained on multiple speakers on this dataset for a fair comparison.

The TED-Talks dataset proposed by [63] contains a list of Youtube video URL links, corresponding frame indices, and cropped areas with auto-detected upper bodies of speakers inside. They are not directly helpful to create the audio-driven reenacted video results. This is because 1) the original dataset only contains very short video clips, e.g., with a duration of a few seconds, which are not sufficient to create rich video motion graphs; 2) the frames in the original dataset are processed to 384×384 resolution by cropping and scaling the upper body of the speaker from a zoomed-out full body frame. As a result, the frames do not have high resolution and high quality. In this case, we use such dataset for numerical evaluation and easier reproduction purpose. To achieve high resolution and high quality audio-driven reenacted TED-talks videos shown on our project page, we use the original full Youtube videos. We manually select the frames with the zoomed-in camera

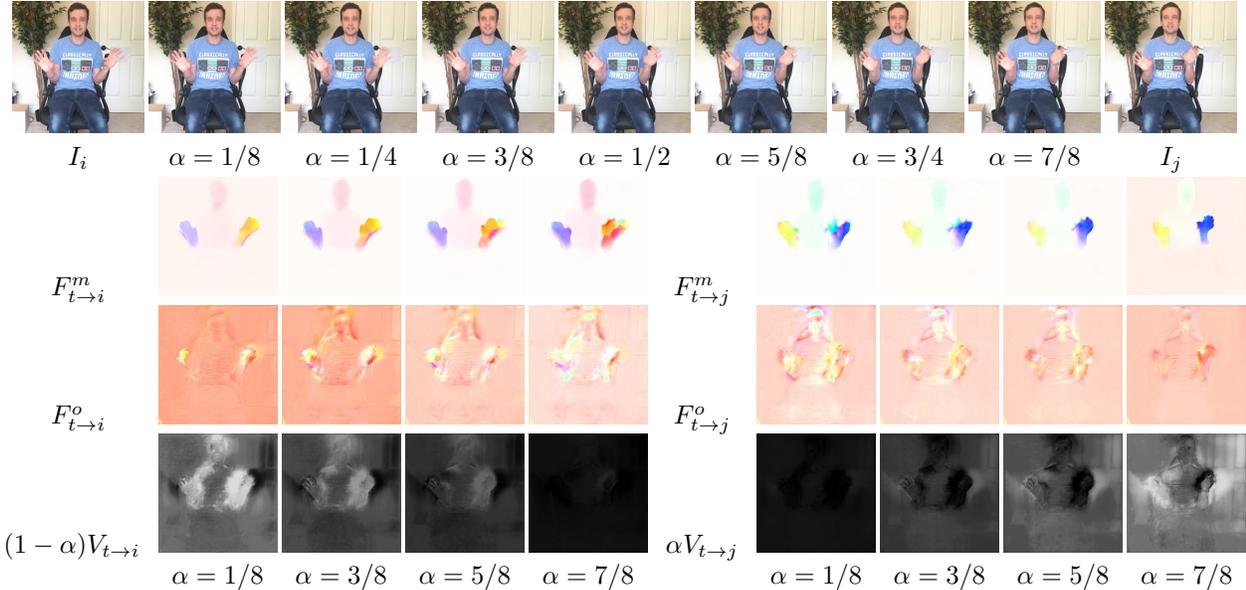


Figure 9. Pose-aware video blending results for target blending weights $\alpha \in (0, 1)$. **Top row**: synthesized in-between frames with blended human gestures for different blending weights. **Bottom rows**: intermediate mesh flows, optical flows and visibility maps results for corresponding blending weights.

where the upper body of the speaker appears at high resolution and high quality (see examples in our HTML files). The selected frames have sufficient length to create reasonable video motion graphs. Finally we fine-tune the model on these frames from each specific speaker and generate reenacted video results given test audios.

Pose-aware video blending network results. Fig. 9 shows output images from the video blending network for different blending weights, along with results from our intermediate stages.

Dictionary of common keywords. Referential gestures, especially iconic and metaphoric gestures, have strong correlations with the transcript [46, 79]. They usually appear together with certain keywords, such as action verbs, concrete objects, abstract concepts, and relative quantities to co-express the speech content [28]. We gather a few frequently used such keywords co-occurring with referential gestures in our speaker videos, as shown in Table. 2.

Network architecture details. The spatial encoder network E_s takes as input the RGB image I_i , the foreground mask I_{mask} , and an image containing the rendered skeleton I_{skel} representing the SMPL pose parameters. Fig. 10 shows an example of these input images.

We show our *Spatial Encoder* network structure for generating the mesh flow warping field in Table 3. In this table, the left column indicates the spatial resolution of the feature map output. The *ResBlock down* block is a 2-strided convo-

lutional layer with a 3×3 kernel followed by two residual blocks. The *ResBlock up* block is a nearest-neighbor up-sampling with a scale of 2, followed by a 3×3 convolutional layer and then two residual blocks. The term *Skip* means skip connection that concatenates the feature maps of an encoding layer and decoding layer with the same spatial resolution. For Personal story dataset, the input and generated images are in 512×512 resolution, while for TED-talks dataset, the image resolution is 384×384 .

The *Mesh Flow Estimator* and *Image Generator* network follows the structure of the *Spatial Encoder* network (see Table 3), but the input and output number of channels are different. For the *Mesh Flow Estimator* network, the number of input feature channel is 13 and output feature channel is 2. For the *Image Generator* network, the number of input feature channel is 19 and output feature channel is 3. Besides, the *Image Generator* network uses in the end a $\tanh(\cdot)$ activation to regularize the image values between $[0, 1]$.

Audio-driven Beam search details. We initialize a beam search [55] procedure in the video motion graph to find K plausible paths matching the target speech audio segments. We set K to 20. The beam search initializes K paths starting with K random nodes as the first frame a_0 for the target audio, then expands in a breadth-first-search manner to find paths ending at a *target graph node* whose audio feature matches the target audio feature at the endpoint of the first segment a_1 , associated with either an activated audio onset

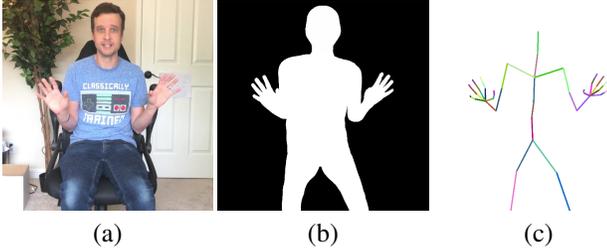


Figure 10. An example of inputs to our spatial encoder network (a) input image frame, (b) corresponding foreground human mask and (c) rendered skeleton image.

or the same non-empty keyword feature. Note that there can be multiple target graph nodes sharing the same audio feature with a_1 .

During the beam search, all the explored paths are sorted based on a *path transition cost*, plus a *path duration cost*. The path transition cost is defined as the sum of node distances between all consecutive nodes m, n along the path, i.e. $\sum_{m,n} (d_{feat}(m, n) + d_{img}(m, n))$. The cost of synthetic transitions are always higher than natural ones. Thus, the path cost prevents using implausible paths with too many synthetic transitions.

When a path reaches a *target* graph node, we check its duration. Due to the sparsity of the graph, there may not be any path matching exactly the target audio segment length L_i . Still, the path length should be similar to L_i , otherwise one would need to accelerate or decelerate the path too much to adjust it to the exact length, leading to unnaturally fast or slow gestures. We only accept paths with duration $L'_s \in [0.9L_s, 1.1L_s]$ since these can be slightly adjusted, e.g. re-sampled, to match the target segment duration. For the above range, we observed that the motion still looks natural. Nevertheless, we also add a path duration cost $|1 - L'_s/L_s|$ to favor paths during beam search with duration closer to the target duration.

When the speech audio is silent, the searched motion graph paths go through nodes without audio onset features, which are often the frames with rest poses.

After processing the first segment $a_0 \rightarrow a_1$, we start another beam search for the next segment $a_1 \rightarrow a_2$. Here, the path expansion starts with the last node of the K paths discovered from previous iteration. The expansion continues with the same search procedure as above. In order, the searches run iteratively for all the rest segments $a_s \rightarrow a_{s+1}$, $s \in [1, S]$ while always keeping the most plausible K paths. All searched K paths can be used to generate various plausible results for the same target speech audio (see demo videos on our project page). The best path is picked in our experiments.

User study details. We provide here more details about the user study.

We have a pool of 381 queries (127 videos from each method \times 3 comparison pairs). For each query, we show two videos in parallel randomly placed at left/right positions. The participants are asked which speaker’s gestures are more consistent with the speech audio and vote for one of the two choices: “left animation”, “right animation”. Fig. 11 shows the webpage layout used in our questionnaires. The layout shows two video results to the participants, a question on the bottom and two choices (“left”/“right”). To enable the selection of either choice, the users must watch both videos until the end. We also explicitly instruct them to focus on the speakers’ hand gestures and ignore the masked facial area.

Our questionnaires also include a similar page layout showing tutorial examples in the beginning. The tutorial shows a pair of videos with clear differences: one video is from ground-truth in which the speaker’s gestures are naturally consistent with the audio; the other video is a failure case, which shows gestures that are inconsistent with audio at some places. For these tutorial cases only, we let the participants pick an answer first and then let them know whether their answer is correct or wrong and explain why.

We also adapt a user validation check to filter out unreliable MTurkers. Specifically, after the tutorial, our questionnaires showed 10 queries in a random order. 3 of the queries were repeated twice (i.e., we had 7 unique queries per questionnaire). We randomly flipped the two videos each time to detect unreliable participants giving inconsistent answers. We filter out unreliable MTurk participants who give different answers to two (or more) of the repeated queries in the questionnaire or took less than 5 minutes to complete it. Each participant was allowed to answer one questionnaire maximum to ensure participant diversity. We collected answers from 113 reliable participants for our user study. We paid \$1 per questionnaire. All comparison outcomes are statistically significant using a z-test ($p < .05$).

Importance of the reference video. The key idea of using reference video is that it provides personalized gestures. Directly animating a single portrait is hard since it is not clear what are the ‘correct’ gestures. There are many applications of our setup. For example, in video production, there is a need to add or remove sentences from existing clips. In online education, different video lessons can be created based on a reference video.

Runtime speed Generating a video from a 15 second input audio and a 2 minute reference video takes about 43 seconds in total. Here is the breakdown: (a) 8 seconds are used for audio-driven search to find graph paths in the video motion graph, (b) 35 seconds are used for synthesizing all transitions. Specifically, for a 15 second input audio, there are maximum of 4 synthetic transitions in our examples, with 8 blended frames created per transition. For blending,

Question 1 out of 10

Press play to start each of the two videos. Watch them **UNTIL THE END**, then you will be able to answer below



Please look carefully at the **speakers' hand gesture** in each video above and **listen to the audio** while the video is playing. Feel free to rewind it. Which of the two speakers' gestures are **MORE consistent with the speech audio**?

- LEFT
- RIGHT

NEXT

Figure 11. User study questionnaire page.

512×512	Input RGB image, foreground mask image, and rendered skeleton image
256×256	ResBlock down $(16 + 2) \rightarrow 32$
128×128	ResBlock down $32 \rightarrow 64$
64×64	ResBlock down $64 \rightarrow 128$
32×32	ResBlock down $128 \rightarrow 256$
16×16	ResBlock down $256 \rightarrow 512$
8×8	ResBlock down $512 \rightarrow 512$
8×8	ResBlock up $512 \rightarrow 512$
16×16	Skip + ResBlock up $(512 + 512) \rightarrow 512$
32×32	Skip + ResBlock up $(512 + 512) \rightarrow 256$
64×64	Skip + ResBlock up $(256 + 256) \rightarrow 128$
128×128	Skip + ResBlock up $(128 + 128) \rightarrow 64$
256×256	Skip + ResBlock up $(64 + 64) \rightarrow 32$
512×512	Skip + ResBlock up $(32 + 32) \rightarrow 16$

Table 3. Spatial Encoder network structure.

obtaining the initial mesh flow from human fitting takes 1 second, then synthesizing each blended frame takes 0.1 seconds measured on a single Tesla V100 GPU.

Personal data / human subjects. The Personal story dataset contains 7 videos with 6 different speakers (5 male, 1 female). The number of frames ranges from 4465 to 19176 (148 to 639 seconds). We collected it under the permission from each speaker to include frames, clips and full video in the paper submission. We also used the TED-talks dataset from the previous work [63]. The perceptual user study is collected with the approval of IRB.