

SPACE : Speech-driven Portrait Animation with Controllable Expression

Siddharth Gururani, Arun Mallya, Ting-Chun Wang, Rafael Valle, Ming-Yu Liu

NVIDIA

{sgururani, amallya, tingchunw, rafaelvalle, mingyul}@nvidia.com



Pose	Generated	Transferred	Generated	Generated
Emotion	Neutral	Neutral	Happy (top) / Sad (bottom)	Surprise (top) / Fear (bottom)

Fig. 1. **SPACE** enables speech-driven animation of a portrait, with control over the output pose, emotions, and intensities of expressions. It can handle a wide range of input poses and produce realistic and high-resolution outputs. *Click any result to play video.*

Abstract

Animating portraits using speech has received growing attention in recent years, with various creative and practical use cases. An ideal generated video should have good lip sync with the audio, natural facial expressions and head motions, and high frame quality. In this work, we present **SPACE**, which uses speech and a single image to generate high-resolution, and expressive videos with realistic head pose, without requiring a driving video. It uses a multi-stage approach, combining the controllability of facial landmarks with the high-quality synthesis power of a pretrained face generator. **SPACE** also allows for the control of emotions and their intensities. Our method outperforms prior methods in objective metrics for image quality and facial motions and is strongly preferred by users in pair-wise comparisons. The project website is available at <https://deepimagination.cc/SPACE/>

1. Introduction

Speech-driven portrait animation, which concerns animating a still image of a face using an arbitrary input speech signal, has a wide range of applications. For example, it can be used for driving characters in computer games, dubbing in movies, and animating avatars for virtual assistance, virtual reality, and telecommunications. It has to use the provided speech to predict all the nuances in human facial expressions while also guaranteeing that the animation looks natural, matches what is being said in the speech sample, and preserves the per-frame and video quality.

These requirements make the task of speech-driven portrait animation challenging. To make things harder, there exist a large number of languages and facial structures that can be provided as input, each with its own unique characteristics. Furthermore, the mapping from a given input speech to corresponding facial motions is inherently one-to-many

Table 1. **Comparison of functionality with prior works.** For fairness, we only compare with one-shot person-agnostic animation models.

	SPACE (ours)	Wav2Lip [23]	MakeItTalk [43]	PC-AVS [42]	Audio2Head [37]	MEAD [36]	EAMM [16]	GC-AVT [19]
Controllability								
• Emotion	✓	✗	✗	✗	✗	✓	✓	✓
• Facial landmarks	✓	✗	✓	✗	✗	✗	✗	✗
• Pose transfer	✓	✓	✗	✓	✗	✗	✓	✓
• Pose generation	✓	✗	✓	✗	✓	✗	✗	✗
Output resolution	512	96	256	224	256	384	256	256

due to variations in head poses, emotions, and expressions. Despite these challenges, the ability to control each of these aspects is vital. For example, a video game character animation application might prefer the generation of exaggerated expressions and head poses, while a newscaster animation application would prefer a neutral expression. In addition, natural motions and high-resolution outputs are desirable to provide the best end-user experience. Unfortunately, no prior framework supports the ability to control emotions, facial landmarks, and poses in a single framework while producing a high-resolution output video, as summarized in Table 1.

In this work, we present SPACE—a method for **Speech-driven Portrait Animation with Controllable Expression**. SPACE decomposes the task into several subtasks that allow for better interpretability and fine-grained controllability. Given an input speech and a facial image, we first predict facial landmark motions in a normalized space. Operating in the facial landmark space gives us the ability to modify facial features and add actions such as blinking when desired. Next, we apply the desired head pose to the facial landmarks and transform them into a latent keypoint space used by our pretrained face image generator [39]. This unsupervisedly-learned latent keypoint space has been shown to produce better synthesis quality than using conventional facial landmarks [25, 26, 39]. Finally, we feed these per-frame latent keypoints to our generator and produce an output video at 512×512 resolution. SPACE also introduces emotion conditioning, enabling control over the emotion types and intensities in the generated video.

Even though previous approaches have also followed similar strategies wherein the audio is mapped to an intermediate representation such as facial landmarks [43] or latent keypoints [16, 37], no previous work has utilized both facial landmarks and latent keypoints simultaneously as intermediate face representations. By using *both* explicit and latent keypoints, we are able to leverage the interpretability and direct controllability of explicit landmarks while also taking advantage of better motion transfer and image quality obtained with latent keypoints and a pretrained generator. The main contributions of our work are as follows:

- We achieve state-of-the-art quality for speech-driven portrait image animation. SPACE provides better quality in terms of FID and landmark distances compared to previous methods while also generating higher-resolution output videos.
- Our method can produce realistic head poses while also being able to transfer poses from another video. It also offers increased controllability by utilizing facial landmarks as an intermediate stage, allowing manipulations such as blinking, eye gaze control, *etc.*
- For the same set of inputs, our method allows for the manipulation of the emotion labels and their corresponding intensities in the output video.

2. Related work

Animating faces using speech signals was first explored by M. Brand in the 1990s, who presented a hidden Markov model-based approach for *Voice Puppetry* [1]. More recently, researchers have proposed multiple methods for the task of speech-driven facial animation using deep neural networks (DNNs). Below, we discuss prior works in audio-based talking-head synthesis, video-driven facial animation, and emotion manipulation in images. A high-level comparison of current methods is provided in Table 1.

Speech-driven facial animation. Approaches in this area learn to map the input speech to facial representations. Some methods animate 3D models of faces such as meshes or standard FACS-based face rigs [7, 18, 30], whereas others directly animate raw images of faces [6, 37, 43].

Among methods that animate 3D models of faces, Taylor *et al.* [30] first obtain phoneme sequences from audio and learn a mapping from phonemes to the mouth pose of a face model, which can be easily edited and re-targeted to other faces by animators. Zhou *et al.* [44] propose VisemeNet, which predicts JALI-based visemes in a multi-stage approach. Most recently, the transformer architecture has also been leveraged for the task by Fan *et al.* [7]. Karras *et al.* [18] animate 3D vertices of a face given a speech signal by utilizing an end-to-end deep network consisting of

a formant analysis network, an articulation network, and a learned emotion embedding. Unfortunately, these methods require special training data, such as 3D face models, which may not be available for many applications.

Another set of approaches learns to directly animate 2D face images. Zhou *et al.* [41] disentangle identity information and speech content and use adversarial training for learning a disentangled joint audio-visual representation. Finally, speech and identity embeddings are used to condition a temporal GAN-based generator [34]. Prajwal *et al.* [23] propose Wav2Lip, which focuses on re-dubbing videos by generating realistic lip motions by using a strong lip-sync discriminator to penalize incorrect lip shapes. However, their outputs lack realism when used with a single input image as the remainder of the face remains stationary. Zhou *et al.* [43] present MakItTalk, which leverages a voice conversion model to extract disentangled speech content and speaker identity features and learns to animate facial landmarks of a given portrait or cartoon image in a speaker-aware fashion. This line of research, while being able to generate reasonable lip motions, lacks the ability to control head poses since no 3D is involved. To circumvent this problem, Zhou *et al.* [42] propose Talking-Face PC-AVS, a method that controls the head pose of a talking face by disentangling identity, pose, and speech content, with the drawback that a driving video is required for pose control.

Our method combines the benefits of both approaches by using a two-stage prediction framework. We first predict pose-normalized facial landmarks for each time step given the input audio. This gives us control over the individual landmarks, as required for fine-grained editing. Further, any rotation and translation can be applied to these landmarks with either predicted or input poses. This allows us to have full control over the 3D head poses, as well as individual landmarks, while using only a single 2D image instead of sophisticated 3D face models.

More recently, the focus has expanded to control the emotion of the output subject. Wang *et al.* [36] collect the MEAD dataset and present a method to condition talking head generation on emotion labels. Ji *et al.* [17] extend the work to disentangle audio and emotion and apply the speech content and emotion to an input video. These methods are identity-specific and do not work in the one-shot setting. Ji *et al.* [16] present a one-shot generation method that disentangles emotions by utilizing a driving emotion video, which is used to guide expressions of the source image being animated. Our method is able to extend emotional control to the one-shot setting without the need for additional emotion video input.

Video-driven facial animation. Facial animation using videos has seen profound success recently. Early works primarily focused on single subjects, where each network can only handle the specific subject seen during training [8, 29, 31–33]. These models usually generate high-

quality results, but have limited use cases since they do not easily extend to new subjects. More recent works can perform subject-agnostic facial animation based on motions from another video [2, 3, 9, 13, 15, 21, 24, 28, 35, 38, 41]. For example, the first-order motion model [25] learns 2D latent keypoints and uses them to predict affine motions to animate the image. Face-vid2vid [39] learns 3D latent keypoints to warp the image, with control over the output poses. However, these methods rely on another video to borrow driving motions from, while we only need audio input to animate an image.

In our framework, we map the input audio to the latent space of face-vid2vid [39], the state-of-the-art video-driven facial animation network, which is fixed and focuses on synthesizing the final images. Using a pretrained face synthesis network gives us several benefits. First, the training time is largely reduced as we only need to focus on generating latent keypoints. Second, it allows us to synthesize a high-resolution output at 512×512 pixel resolution. Finally, by modulating the inputs with emotion labels, we learn the ability to control emotion and emotion intensity.

3. Method

SPACE takes an input speech clip and a face image, as well as (optionally) an emotion label, and produces an output video. It decomposes this task into three stages:

1. **Speech2Landmarks** (Sec. 3.2): Given an input image, it extracts normalized 3DDFA [12, 45] and MTCNN [40] facial landmarks, and predicts their per-frame motions based on the input speech and emotion label.
2. **Landmarks2Latents** (Sec. 3.4): This step translates the per-frame posed facial landmarks into latent keypoints used by face-vid2vid [39], a pretrained image-based facial animation model.
3. **Video synthesis** (Sec. 3.6): Given the input image and the per-frame latent keypoints predicted in the previous step, the face-vid2vid generator outputs an animated video.

This decomposition has multiple advantages. First, it allows for fine-grained control on the output facial expressions. For example, facial landmarks can be modified to introduce eye blinking or manipulated to apply the desired head pose, either using provided pose inputs or predicted poses (Sec. 3.3). Further, latent keypoints can be modulated with emotion labels (Sec. 3.5) to change the expression intensity, as well as control the gaze direction. By leveraging a pretrained face generator, we are able to reduce the training cost, as well as obtain high-quality output videos. The overall framework is illustrated in Fig. 2, and components are described below.

3.1. Preliminaries

Pretrained talking-head animation model. Instead of learning an image generator from scratch, we rely on pretrained face-vid2vid [39]: a state-of-the-art framework for

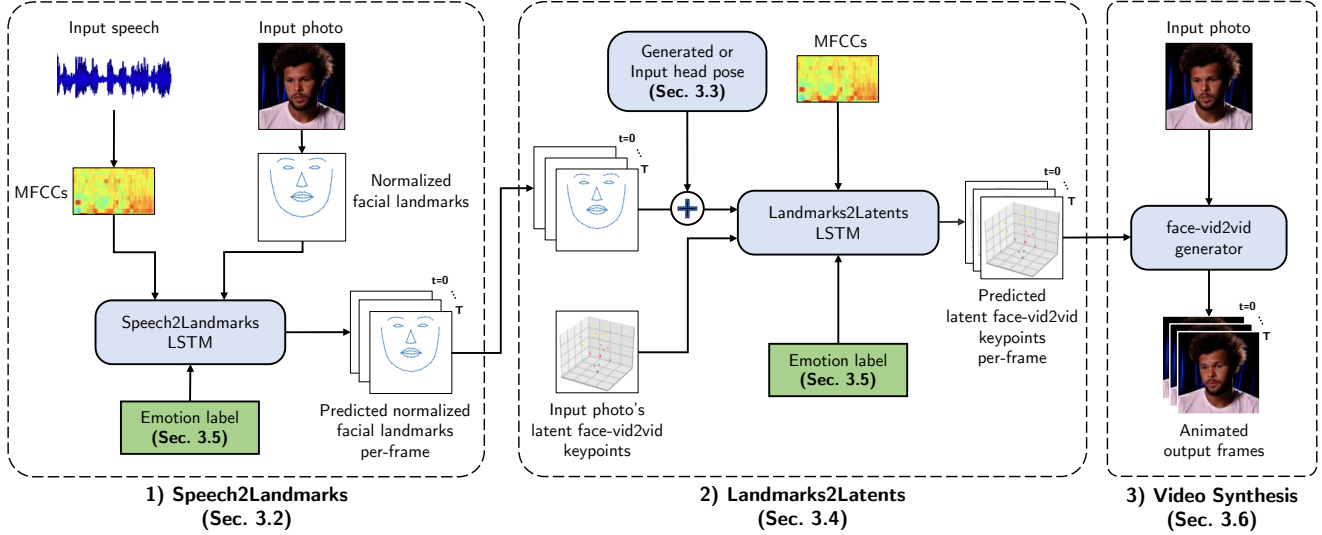


Fig. 2. **Overview of SPACE.** Our framework consists of three stages. The first stage, Speech2Landmarks, predicts per-frame normalized facial landmarks given the initial facial landmarks and input speech. The second stage, Landmarks2Latents, predicts the latent face-vid2vid keypoints for each frame, using the input audio and previously predicted facial landmarks. The last stage warps the input image using the latent keypoints to produce animated output frames at 512×512 px. Note that the intermediate predictions allow for increased controllability – facial landmarks can be modified to add eye blinking and change head pose, while latent keypoints can be used to modulate emotions.

animating a single source image using motions from a chosen driving video. To transfer motions from a driving video to a source image, it relies on unsupervisedly-learned latent keypoints extracted from the input image and video frames.

For each input frame, the face-vid2vid encoder predicts 20 latent keypoints. Given the latent keypoints of the input source image (kp_s) and the current driving video frame (kp_t), the decoder predicts a flow-based warping field. By applying the warping field to the source image features, it produces an output image that contains the identity of the source image in the pose of the driving frame. Further details are available in the supplementary and the face-vid2vid paper [39].

As described later in Sec. 3.4, we predict latent keypoints per frame given the input (source) image and speech. Using source image features and keypoints, along with the predicted per-frame keypoints, the face-vid2vid decoder produces an output image for each time step. This allows us to produce outputs at 512×512 resolution, higher than all prior works on speech-based photo animation.

Dataset preprocessing. Given a talking-head video, we first extract per-frame facial landmarks. We use the 3DDFA [11, 12, 45] landmark detector to extract 68 3D facial landmarks and head poses from each frame. We then frontalize the 3D facial landmarks by using the estimated head pose. Specifically, we rotate the face such that the nose tip is facing straight to the camera, aligned with the camera axis. The frontalized 3D facial landmarks are then orthographically projected to 2D, and each frame is normalized such that the distance between the two ear landmarks is fixed. To

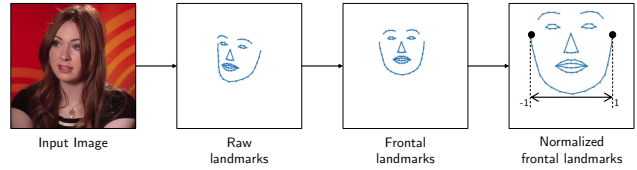


Fig. 3. **Input data preprocessing.** In order to obtain normalized facial landmarks per-frame for training, we run 3DDFA on each image to obtain 3D raw landmarks and the head pose (rotation and translation). We obtain frontal landmarks by undoing the rotation. Finally, we normalize the landmarks such that the distance between the landmarks of the left ear and the right ear is 2.

obtain accurate eye landmarks, we use MTCNN [40]. The 52 eye landmarks are required for controlling blinks and gazes. An overview of our facial landmark normalization is shown in Fig. 3. Additionally, we perform various steps of data filtering to remove noisy data from our training set.

In addition to landmarks, we also extract latent keypoints per frame using the pretrained face-vid2vid encoder. This gives per-frame (facial landmarks, latent keypoints) pairs.

From the audio, we extract 40 Mel-Frequency Cepstral Coefficients (MFCCs) using a 1024 sample FFT window size at 30 fps in order to align the audio features with the video frames. We also use audio data augmentation methods such as pitch-shifting, equalization, loudness variation, *etc.*

For emotion labels, we use those provided with the training dataset. If unavailable, we utilize a pretrained emotion classifier [27] to predict the emotion for each frame.

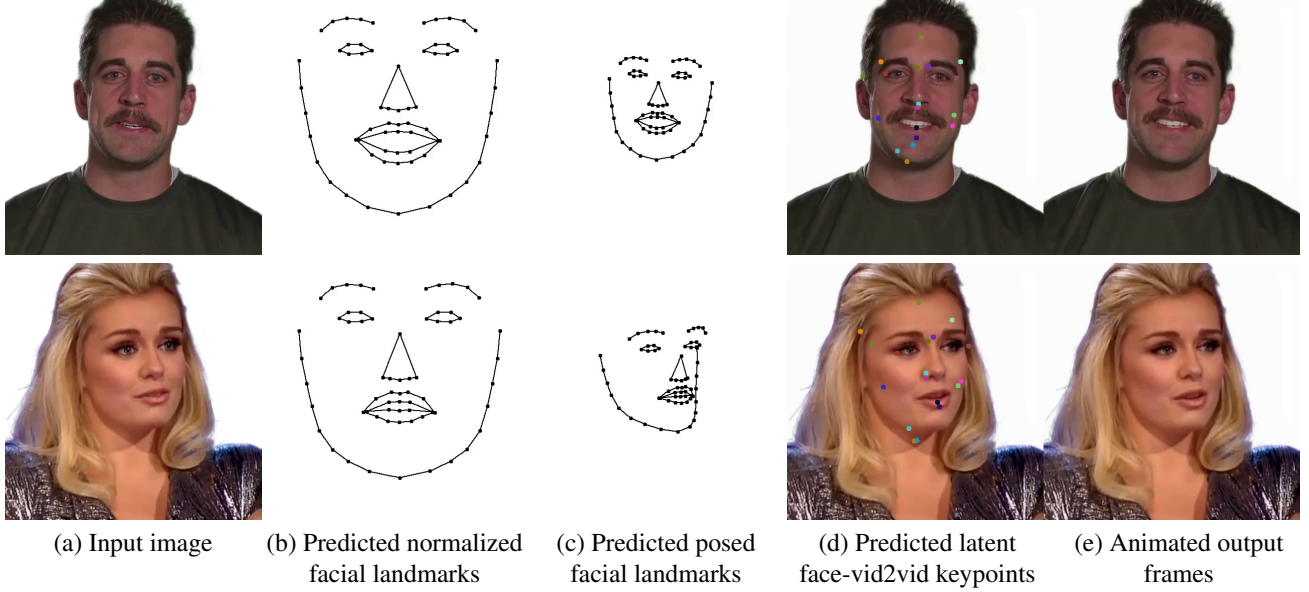


Fig. 4. **Inputs, predicted intermediate and final outputs, for a fixed pose.** In the first Speech2Landmarks (Sec. 3.2) stage, given an input image (a) and speech, our method predicts facial landmark motions in the normalized space (b). These predictions are scaled, rotated, and translated to the required pose by PoseGen (Sec. 3.3). In these samples, we transform them back to the input image space (c). In the Landmarks2Latents stage (Sec. 3.4), these facial landmarks are used to predict corresponding face-vid2vid latent keypoints (d) (which are shown overlayed on the final output image). In the final Video Synthesis stage (Sec. 3.6), the face-vid2vid generator produces the final outputs (e) using the latent keypoints and the input image. *Click any result to play video.*

3.2. Speech2Landmarks

The Speech2Landmarks (S2L) network is the first step of our framework, as shown in Fig. 2. It uses a Long Short-Term Memory (LSTM) regressor for predicting the normalized facial landmarks given the input audio MFCCs and the normalized facial landmarks of the input image. This is represented by

$$\text{face}_t = \text{LSTM}_{\text{S2L}}(\text{face}_{t-1}, a_t), \quad (1)$$

where face_t and a_t are the facial landmarks and audio features at time t , respectively. We use a convolutional neural network (CNN) to encode the audio and a multi-layer perceptron (MLP) to encode the facial landmarks. We train this network using an \mathcal{L}_1 loss with the ground truth normalized facial landmarks, with a higher loss scale on the y-axis to more heavily penalize vertical motion errors. We additionally use a velocity loss, an \mathcal{L}_1 loss between the first-order temporal difference of ground truth and predicted landmarks. Given source images shown in Fig. 4 (a) and the associated speech input, outputs from our S2L model are visualized in Fig. 4 (b).

Using facial landmarks as an intermediate representation is advantageous as it allows for the explicit manipulation of facial features. For example, we can add eye blinks by manipulating the eye landmarks. We found making predictions in the normalized space to be very important to simplify the

mapping between phonemes and lip motions.

3.3. Pose generation

For a given input speech, many different head pose sequences are valid. To model this variation, we use the conditional variational autoencoder architecture of Greenwood *et al.* [10]. We train an autoregressive decoder that predicts the rotation R and translation T for the facial landmarks. R is represented by three angles: yaw, pitch, and roll. T is the displacement of the 3D landmarks in the x-, y-, and z-axes. Our network, PoseGen, at time step t is expressed as

$$R_t, T_t = \text{LSTM}_{\text{pose}}(a_t, z), \quad (2)$$

where z is a noise sampled from $\mathcal{N}(0, 1)$ during inference, and $\mathcal{N}(\mu, \sigma)$ during training. μ, σ are predicted by a bidirectional LSTM encoder that jointly encodes audio and pose.

The poses, whether predicted or extracted from a reference video, are applied to the frontal normalized landmarks predicted by our S2L module. This transforms the normalized landmarks back to the image space after an appropriate scaling factor is applied. Examples of the predicted facial landmarks (transformed to a fixed pose for better visualization) are shown in Fig. 4 (c). Time-varying pose outputs using our generated poses are shown in Fig. 1 and the supplementary material.

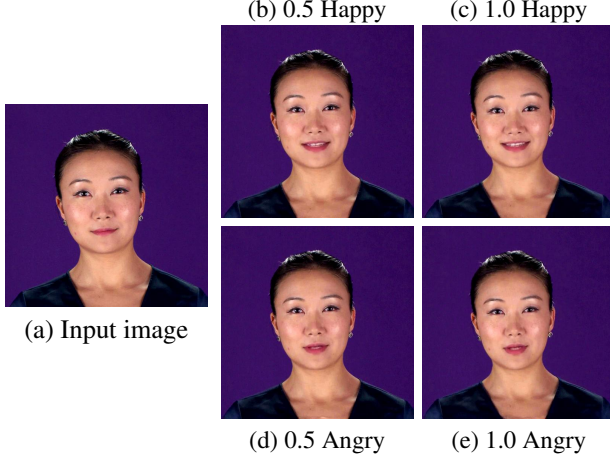


Fig. 5. **Emotion label and intensity control.** We can condition outputs on both the emotion label and its intensity. *Click any result to play video.*

3.4. Landmarks2Latents

The second step, Landmarks2Latents (L2L), takes in posed facial landmarks, and generates corresponding latent face-vid2vid keypoints, as shown in Fig. 2. An LSTM regressor predicts latent keypoints given encoded audio features, predicted posed landmarks, latent keypoints of the input image, and the previous prediction. L2L can be expressed as

$$\text{kp}_t = \text{LSTM}_{\text{L2L}}(R_t \cdot \text{face}_t + T_t, \text{a}_t, \text{kp}_s, \text{kp}_{t-1}), \quad (3)$$

where face_t is predicted by S2L (Eq. 1), R_t and T_t are predicted by PoseGen (Eq. 2), and kp_s is predicted from the input image by the face-vid2vid encoder. For training, we use the (facial landmarks, latent keypoints) pairs generated as described in Sec. 3.1. We train this network using an \mathcal{L}_1 loss with the ground truth latent keypoints. Examples of the predicted latent keypoints are shown in Fig. 4 (d), overlaid on the final outputs.

3.5. Emotion control

To provide additional user control on the generated video, we condition both S2L and L2L on the emotion of the video frames. This is achieved using Feature-wise Linear Modulation (FiLM) layers [22]. For the S2L network, we use FiLM to modulate the audio features and the initial landmark input. For L2L, we apply FiLM on the audio, landmarks, and the initial latent keypoint input. As mentioned earlier in Sec. 3.1, we use one-hot emotion labels when provided by the dataset, and a predicted per-frame probability distribution over the emotions [27] when the ground-truth emotions are unavailable. At inference, we can provide the desired combination of emotion labels and their intensities as input.

Even though we do not impose any constraints to disentangle the emotion label from the input speech, unlike prior

work [16, 17, 36], we find that our model is able to change facial emotions, as shown in Fig. 5. We believe there are two main reasons this simple modulation works – 1) Our audio encoder is shallow and unidirectional. Due to a small lookahead, various augmentations, and lack of bidirectionality, it is unable to extract emotions from the audio; 2) Thus, relying on emotion-conditioned modulations of intermediate landmarks and latent keypoints, which directly control facial expressions, becomes the most effective way to produce the required expressions.

3.6. Video synthesis

The third and final step, as shown in Fig. 2, is the talking-head video generation by the pretrained face-vid2vid [39] generator. The generator uses the input image and the per-frame latent keypoints kp_t predicted by the previous step, L2L (Eq. 3), and outputs video frames at 512×512 pixel resolution. Examples of the final outputs and intermediate predictions are shown in Fig. 4.

To summarize, given an input image and audio, our method learns to produce valid sequences of face-vid2vid latent keypoints. These latent keypoints are used to animate the input image using the pretrained face-vid2vid generator. In order to provide additional user control, such as blinking and pose change, we use facial landmarks as an interpretable intermediate prediction. Lastly, the emotion conditioning helps change the emotion and intensity of output expressions. Additional architectural and implementation details are available in the supplementary material.

4. Results

In this section, we discuss the datasets, metrics, and various experiments we conduct to validate our method.

Datasets. We train our models using three different datasets: VoxCeleb2 [4], RAVDESS [20], and MEAD [36]. VoxCeleb2 is made up of a large collection of YouTube videos of celebrities speaking in various settings, such as interviews, speeches, and podcasts. It consists of a diverse set of speakers with a wide range of languages and accents. RAVDESS is a smaller-scale dataset compared to VoxCeleb2, consisting of approximately 7k clean video recordings of various actors speaking two sentences, with 8 different emotions in two levels of emotional intensity. MEAD is a larger-scale emotional dataset consisting of 40 hours of audio-visual recordings from 60 actors (48 of which are made publicly available). It consists of multi-camera recordings of actors speaking thirty sentences with 8 emotions in three levels of intensity.

After combining, preprocessing, and filtering the three datasets as described in Sec. 3.1, we obtain $\sim 270\text{k}$ training videos. For evaluation, we manually select and set aside 350 videos of unseen identities. This subset is selected to ensure

it captures a diverse set of speakers with different initial poses. More details of the datasets are in the supplementary.

Baselines. We compare SPACE to 3 prior works on speech-driven animation – 1) **Wav2Lip** [23], 2) **MakeItTalk** [43], and 3) **Talking Face-PC-AVS (PC-AVS)** [42]. We chose these baselines for two main reasons. First, the source code and pretrained model weights are available. Second, they represent the state-of-the-art talking-head video generation using arbitrary identities. These methods, like ours, do not need to be fine-tuned or trained separately for a new identity that was not seen during training. The publicly available code for MEAD [36], which has emotion control, does not support arbitrary identities for inference. Note that Wav2Lip and PC-AVS both require an input video to determine the pose of the generated output. However, we provide all models with only two inputs in our setting: an initial photo and the driving speech. Therefore, these two methods will naturally only generate videos with a single pose. Thus, we use a fixed pose with SPACE for a fair comparison.

Metrics. We use the following metrics to measure facial landmark motion and photorealism:

- **Lip sync quality:** We use the lip sync confidence score output by SyncNet [5] as a proxy for lip sync quality.
- **Landmark accuracy:** We extract facial landmarks from outputs, frontalize, and normalize them such that the left and right edge of the mouth landmarks are at $(-1, 0)$ and $(1, 0)$, respectively. We measure the mean absolute error (MAE) between the predicted and ground truth mouth landmarks positions (M-P) and velocities (M-V). We compute errors in face geometry (F-P) and velocity (F-V), using the normalization from Sec. 3.1.
- **Photorealism:** We measure the Frechet Inception Distance (FID) [14] of outputs with the ground truth.
- **Human evaluation:** We perform a forced preference A/B test with the four baseline models for the same input image and speech audio. Each pair is rated by 3 users, and we report the average preference score.

Additional details about the metrics, including the human evaluation interface, are available in the supplementary.

Quantitative results. As shown in Table 2, SPACE achieves the lowest FID, indicating that our method produces the best image quality. We also outperform the baselines on the normalized landmark distance metrics. Despite obtaining better output quality and lower temporal jittering, our method achieves poorer SyncNet scores. We found that SyncNet scores are very sensitive to the input crop and suspect that our larger-crop outputs may have resulted in worse scores. Surprisingly, the mean SyncNet score on real videos is 4.85, lower than Wav2Lip at 5.08 and PC-AVS at 7.00.

Qualitative results. Fig. 6 shows examples of videos generated by various methods. While previous methods work on

Table 2. **Quantitative evaluation results.** ↓ indicates lower is better. Although PC-AVS obtains a better SyncNet score, we found the score to be highly sensitive to the input crop. In fact, the mean score on real videos is 4.85.

	FID ↓	M-P ↓	M-V ↓	F-P ↓	F-V ↓	Sync
PC-AVS [41]	66.68	0.032	0.013	0.024	0.004	7.00
MakeItTalk [43]	30.68	0.035	0.012	0.020	0.004	2.61
Wav2Lip [23]	15.67	0.030	0.013	0.017	0.004	5.08
Ablations						
S2Latent	11.82	0.021	0.005	0.013	0.003	2.94
L2L-landmarks	15.63	0.026	0.007	0.008	0.002	2.78
SPACE (ours)	11.68	0.025	0.007	0.008	0.002	3.61

Table 3. **User preference scores.** We perform a forced preference test where users are presented with two videos with audio. They are asked to select which one of the two is more realistic, with a focus on the face and mouth region. We obtain ~900 ratings per pair of models. Users overwhelmingly prefer SPACE.

	PC-AVS [41]	MakeItTalk [43]	Wav2Lip [23]
SPACE (ours)	89.1%	87.8%	72.6%

frontal-facing and closely-cropped input images, they suffer from degraded quality or fail for arbitrary poses and larger crops. Since we predict mouth motions in the normalized landmark space and then apply pose transformations, we are able to handle poses from face images in the wild. In addition, SPACE is also able to generate missing details such as teeth, while other methods either fail or introduce artifacts. We also found that our method adds realistic head and shoulder motions when the audio has a breathing sound. We find that SPACE has smoother lip motion compared to Wav2Lip and PC-AVS, where the motions tend to be more exaggerated. We suspect that this is another reason for poorer SyncNet confidence scores for our method.

As seen in Table 3, users overwhelmingly prefer SPACE over other methods. This study asks users to rate which of the two videos is more realistic in a two-choice test.

Ablations. To isolate the contributions of various design choices, we perform ablation studies with three alternatives – 1) S2Latent: directly predicting latent keypoints from speech instead of the proposed multi-stage system, 2) L2L-landmarks: S2L uses audio, but L2L only uses landmarks, and 3) S2L-raw: using raw landmarks to train S2L.

From Table 2, it can be seen that S2Latent achieves strong performance on the chosen metrics; however, since it is not informed of the pose, there tends to have pose drift in the generated videos, leading to poor overall motion quality. Additionally, we are unable to control aspects such as eye blinks and mouth opening/closing in this case. On the other hand, for L2L-landmarks we find the fine-grained motions to be missing. This can be attributed to the fact that facial

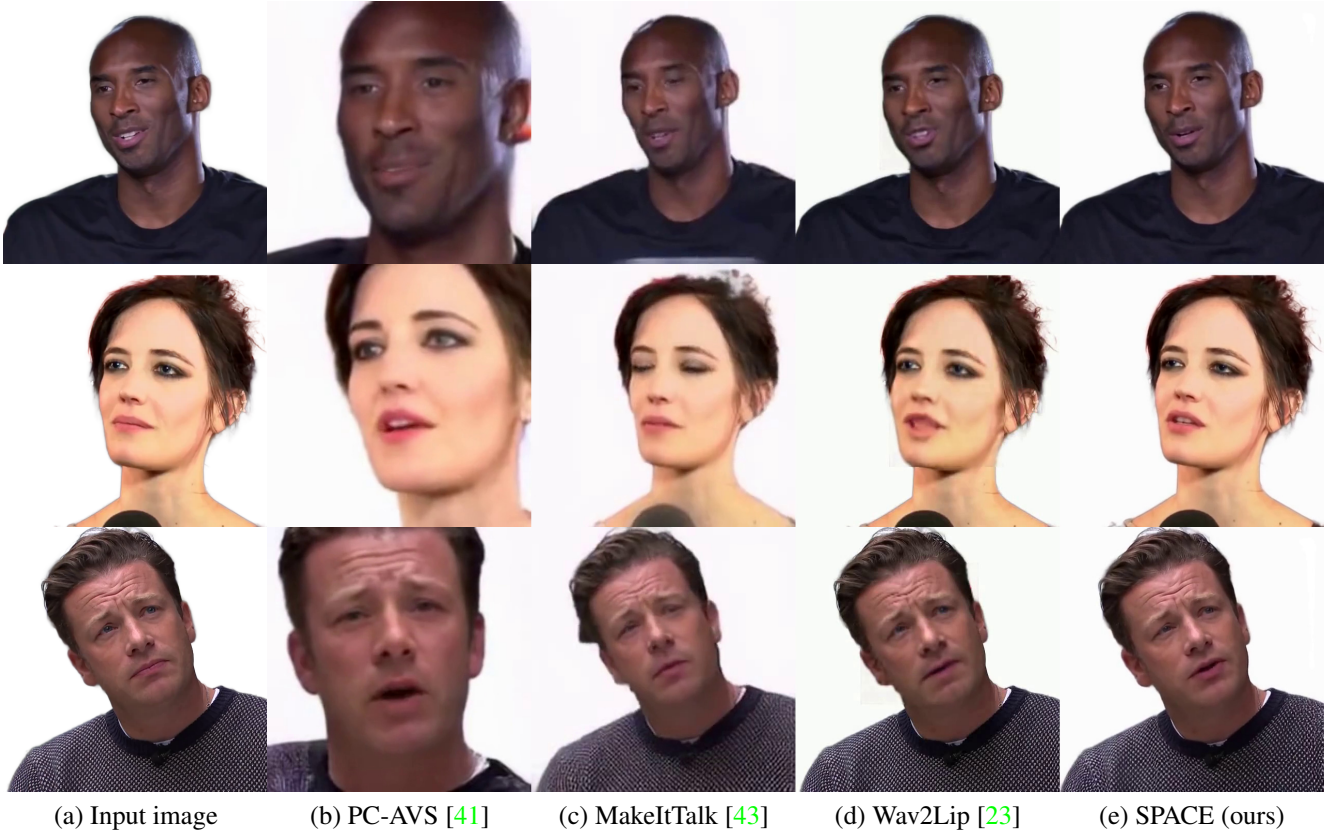


Fig. 6. **Comparisons with prior works.** SPACE is able to animate faces in any pose, with better temporal stability and per-frame quality. Note the temporal jittering in (b), deformed face shapes in (c), and unnatural mouth motions in (d). *Click any result to play video.*

landmarks do not accurately capture fine-grained mouth motions and expressions, while our full model can rely on audio for the missing details. Finally, S2L-raw suffers from both of the issues, i.e., failing to predict accurate mouth motions, and pose drift due to missing pose information. Hence we exclude it from the quantitative evaluations. For visual results, please see the supplementary material.

Novel applications. SPACE can be used in video conferencing as sending only audio and a still image saves substantial bandwidth over sending the full video. Furthermore, due to its pose controllability, it can be combined with existing approaches so that the user can freely switch between using video or audio inputs. In low bandwidth scenarios, the video conferencing system can fall back to the audio-driven mode and still generate realistic output videos. We show such an application with seamless video-to-audio-driven switching in the supplementary material.

5. Discussion

In this work, we have presented SPACE, a speech-driven face animation framework that produces state-of-the-art realistic and high resolution videos. Due to our novel inter-

mediate representations, we also gain control over the head pose, eye blinking, and gaze directions. Our method is also able to utilize emotion labels and their intensities as input to produce diverse outputs with different emotions. Moreover, SPACE also possesses the ability to generate poses from audio or use poses from a different driving video. The capabilities of SPACE open new avenues for applications in video conferencing, gaming, and media synthesis.

Limitations and Societal impact. Our method might fail on underrepresented emotions, phonemes, and visemes. A diverse training dataset should help mitigate such issues. Extreme head poses can cause inaccuracies in 3DDFA predictions, and thus our method. Our method has the potential for negative impact if used to create *deepfakes*. Using voice cloning, a malicious actor can create videos enabling identity theft or dissemination of fake news. However, in controlled settings, SPACE can be used for positive creative purposes.

References

- [1] Matthew Brand. Voice puppetry. In *PACMGIT*, 1999. 2
- [2] Egor Burkov, Igor Pasechnik, Artur Grigorev, and Victor Lempitsky. Neural head reenactment with latent pose descriptors.

- In *CVPR*, 2020. 3
- [3] Lele Chen, Ross K Maddox, Zhiyao Duan, and Chenliang Xu. Hierarchical cross-modal talking face generation with dynamic pixel-wise loss. In *CVPR*, 2019. 3
- [4] Joon Son Chung, Arsha Nagrani, and Andrew Zisserman. VoxCeleb2: Deep speaker recognition. In *INTERSPEECH*, 2018. 6
- [5] Joon Son Chung and Andrew Zisserman. Out of time: Automated lip sync in the wild. In *ACCV Workshops*, 2016. 7
- [6] Amanda Cardoso Duarte, Francisco Roldan, Miquel Tubau, Janna Escur, Santiago Pascual, Amaia Salvador, Eva Mo-hedano, Kevin McGuinness, Jordi Torres, and Xavier Giro-i Nieto. WAV2PIX: Speech-conditioned face generation using generative adversarial networks. In *ICASSP*, 2019. 2
- [7] Yingruo Fan, Zhaojiang Lin, Jun Saito, Wenping Wang, and Taku Komura. FaceFormer: Speech-driven 3D facial animation with transformers. In *CVPR*, 2022. 2
- [8] Guy Gafni, Justus Thies, Michael Zollhofer, and Matthias Nießner. Dynamic neural radiance fields for monocular 4D facial avatar reconstruction. In *CVPR*, 2021. 3
- [9] Jiahao Geng, Tianjia Shao, Youyi Zheng, Yanlin Weng, and Kun Zhou. Warp-guided GANs for single-photo facial animation. *TOG*, 2018. 3
- [10] David Greenwood, Stephen Laycock, and Iain Matthews. Predicting head pose from speech with a conditional variational autoencoder. In *INTERSPEECH*, 2017. 5
- [11] Jianzhu Guo, Xiangyu Zhu, and Zhen Lei. 3DDFA. <https://github.com/cleardusk/3DDFA>, 2018. 4
- [12] Jianzhu Guo, Xiangyu Zhu, Yang Yang, Fan Yang, Zhen Lei, and Stan Z Li. Towards fast, accurate and stable 3D dense face alignment. In *ECCV*, 2020. 3, 4
- [13] Sungjoo Ha, Martin Kersner, Beomsu Kim, Seokjun Seo, and Dongyoung Kim. MarioNETte: Few-shot face reenactment preserving identity of unseen targets. In *AAAI*, 2020. 3
- [14] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. GANs trained by a two time-scale update rule converge to a local Nash equilibrium. In *NeurIPS*, 2017. 7
- [15] Amir Jamaludin, Joon Son Chung, and Andrew Zisserman. You said that?: Synthesising talking faces from audio. *IJCV*, 2019. 3
- [16] Xinya Ji, Hang Zhou, Kaisiyuan Wang, Qianyi Wu, Wayne Wu, Feng Xu, and Xun Cao. EAMM: One-shot emotional talking face via audio-based emotion-aware motion model. *TOG*, 2022. 2, 3, 6
- [17] Xinya Ji, Hang Zhou, Kaisiyuan Wang, Wayne Wu, Chen Change Loy, Xun Cao, and Feng Xu. Audio-driven emotional video portraits. In *CVPR*, 2021. 3, 6
- [18] Tero Karras, Timo Aila, Samuli Laine, Antti Herva, and Jaakko Lehtinen. Audio-driven facial animation by joint end-to-end learning of pose and emotion. *TOG*, 2017. 2
- [19] Borong Liang, Yan Pan, Zhizhi Guo, Hang Zhou, Zhibin Hong, Xiaoguang Han, Junyu Han, Jingtuo Liu, Errui Ding, and Jingdong Wang. Expressive talking head generation with granular audio-visual control. In *CVPR*, 2022. 2
- [20] Steven R Livingstone and Frank A Russo. The Ryerson Audio-Visual Database of Emotional Speech and Song (RAVDESS): A dynamic, multimodal set of facial and vocal expressions in North American English. *PLOS One*, 2018. 6
- [21] Arun Mallya, Ting-Chun Wang, and Ming-Yu Liu. Implicit Warping for Animation with Image Sets. In *NeurIPS*, 2022. 3
- [22] Ethan Perez, Florian Strub, Harm De Vries, Vincent Dumoulin, and Aaron Courville. FiLM: Visual reasoning with a general conditioning layer. In *AAAI*, 2018. 6
- [23] 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 *ACM MM*, 2020. 2, 3, 7, 8
- [24] Albert Pumarola, Antonio Agudo, Aleix M Martinez, Alberto Sanfeliu, and Francesc Moreno-Noguer. GANimation: Anatomically-aware facial animation from a single image. In *ECCV*, 2018. 3
- [25] Aliaksandr Siarohin, Stéphane Lathuilière, Sergey Tulyakov, Elisa Ricci, and Nicu Sebe. First order motion model for image animation. In *NeurIPS*, 2019. 2, 3
- [26] Aliaksandr Siarohin, Oliver J Woodford, Jian Ren, Menglei Chai, and Sergey Tulyakov. Motion representations for articulated animation. In *CVPR*, 2021. 2
- [27] Henrique Siqueira, Sven Magg, and Stefan Wermter. Efficient facial feature learning with wide ensemble-based convolutional neural networks. In *AAAI*, 2020. 4, 6
- [28] Yang Song, Jingwen Zhu, Dawei Li, Andy Wang, and Hairong Qi. Talking face generation by conditional recurrent adversarial network. In *IJCAI*, 2019. 3
- [29] Supasorn Suwajanakorn, Steven M Seitz, and Ira Kemelmacher-Shlizerman. Synthesizing Obama: learning lip sync from audio. *TOG*, 2017. 3
- [30] 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. *TOG*, 2017. 2
- [31] Justus Thies, Michael Zollhöfer, and Matthias Nießner. Deferred neural rendering: Image synthesis using neural textures. *TOG*, 2019. 3
- [32] Justus Thies, Michael Zollhöfer, Matthias Nießner, Levi Valgaerts, Marc Stamminger, and Christian Theobalt. Real-time expression transfer for facial reenactment. *TOG*, 2015. 3
- [33] Justus Thies, Michael Zollhofer, Marc Stamminger, Christian Theobalt, and Matthias Nießner. Face2Face: Real-time face capture and reenactment of RGB videos. In *CVPR*, 2016. 3
- [34] Konstantinos Vougioukas, Stavros Petridis, and Maja Pantic. End-to-end speech-driven facial animation with temporal GANs. In *CVPR Workshop*, 2019. 3
- [35] Konstantinos Vougioukas, Stavros Petridis, and Maja Pantic. Realistic speech-driven facial animation with GANs. *IJCV*, 2019. 3
- [36] 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 *ECCV*, 2020. 2, 3, 6, 7
- [37] Suzhen Wang, Lincheng Li, Yu Ding, Changjie Fan, and Xin Yu. Audio2head: Audio-driven one-shot talking-head generation with natural head motion. *IJCAI*, 2021. 2

- [38] Ting-Chun Wang, Ming-Yu Liu, Andrew Tao, Guilin Liu, Jan Kautz, and Bryan Catanzaro. Few-shot video-to-video synthesis. In *NeurIPS*, 2019. 3
- [39] Ting-Chun Wang, Arun Mallya, and Ming-Yu Liu. One-shot free-view neural talking-head synthesis for video conferencing. In *CVPR*, 2021. 2, 3, 4, 6
- [40] Kaipeng Zhang, Zhanpeng Zhang, Zhifeng Li, and Yu Qiao. Joint face detection and alignment using multitask cascaded convolutional networks. *IEEE Signal Processing Letters*, 2016. 3, 4
- [41] Hang Zhou, Yu Liu, Ziwei Liu, Ping Luo, and Xiaogang Wang. Talking face generation by adversarially disentangled audio-visual representation. In *AAAI*, 2019. 3, 7, 8
- [42] Hang Zhou, Yasheng Sun, Wayne Wu, Chen Change Loy, Xiaogang Wang, and Ziwei Liu. Pose-controllable talking face generation by implicitly modularized audio-visual representation. In *CVPR*, 2021. 2, 3, 7
- [43] Yang Zhou, Xintong Han, Eli Shechtman, Jose Echevarria, Evangelos Kalogerakis, and Dingzeyu Li. MakeItTalk: Speaker-aware talking-head animation. *TOG*, 2020. 2, 3, 7, 8
- [44] Yang Zhou, Zhan Xu, Chris Landreth, Evangelos Kalogerakis, Subhransu Maji, and Karan Singh. VisemeNet: Audio-driven animator-centric speech animation. *TOG*, 2018. 2
- [45] Xiangyu Zhu, Xiaoming Liu, Zhen Lei, and Stan Z Li. Face alignment in full pose range: A 3D total solution. *TPAMI*, 2017. 3, 4

Appendix

A. Network architecture and training

A.1. Speech2Landmarks

Fig. 7 shows the high-level architecture of the Speech2Landmarks (S2L) model. The main components include an audio encoder, landmark encoder, and the LSTM regressor that autoregressively predicts the face landmarks frame-by-frame.

The audio encoder is a 12-layer convolutional network with 1024 1D filters. Each layer uses leaky ReLU activations and batch-norm. The first 4 layers use symmetrical filters spanning 3 time-steps. The last 8 layers use causal filters spanning 5 time-steps. The goal of this asymmetrical padding is to have the audio encoder focus more on the past and have fewer time-steps look-ahead, thus enabling use in real-time applications.

There are two landmark encoders. For the 68 3DDFA landmarks we use an 8 layer fully-connected network with 1024 hidden units with leaky ReLU activation and batch norm. For the 52 MTCNN eye landmarks we use an 4 layer fully-connected network with 1024 hidden units with leaky ReLU activation and batch norm. Both the encoders are multi-layer fully connected networks.

The regressor is a 2-layer LSTM with 1024 hidden units. We also found that using a learned initialization for the hidden state leads to better initial outputs. For this purpose we use a 4-layer fully connected network which takes the initial time-step of all inputs:

$$h_0 = \text{MLP}(\text{face}_0, a_0) \quad (4)$$

A.2. Pose Generation

Fig. 8 depicts the high-level architecture of our Pose Generation network (PoseGen). The encoder is a 2-layer bi-LSTM with 256 hidden units. The latent z is 64 dimensional. The decoder is a 2-layer LSTM with 256 hidden units.

A.3. Landmarks2Latents

Fig. 9 shows the high-level architecture of the Landmarks2Latents (L2L) model. The main components include an audio encoder, face landmark encoder, latent keypoint encoder, and the LSTM regressor that autoregressively predicts the latent landmarks frame-by-frame.

The audio encoder, face landmarks encoders, and the regressor follow the same architecture as Speech2Landmarks.

The source encoder encodes the latent keypoints from the source image. It is a 4-layer fully connected network with 1024 hidden units with leaky ReLU activation and batch norm. The latents encoder encodes the latent keypoints of the

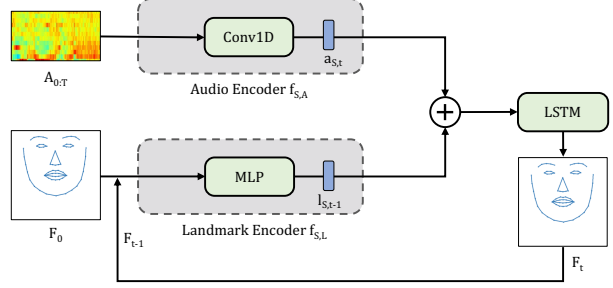


Fig. 7. **Speech2Landmarks architecture.** Speech2Landmarks takes as input the normalized face landmarks F_0 , the mouth landmarks sequence $M_{1:T}$, and the MFCC features of the speech recording $A_{0:T}$. It is an autoregressive model that takes as input the predicted face landmarks from previous time step, and audio features and mouth landmarks for the current time step. It consists of a face landmarks encoder $f_{F,L}$, a mouth landmarks encoder $f_{F,E}$, an audio encoder $f_{F,A}$, and an LSTM network.

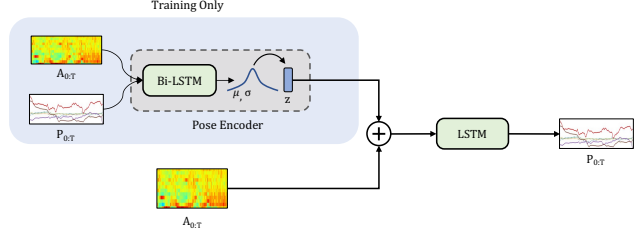


Fig. 8. **PoseGen architecture.** PoseGen is a variational autoencoder which jointly encodes audio and pose sequence. The decoder subsequently reconstructs the pose sequence conditioned on the input audio.

previous time-step. It is an 8-layer fully connected network with 1024 hidden units with leaky ReLU activation and batch norm.

A.4. Video Generation

For details about face-vid2vid, please visit the project page: <https://nvlabs.github.io/face-vid2vid/>.

A.5. Training details

Each network is trained using the Adam optimizer with a batch size of 256. We use a learning schedule which warms up the learning rate from $1e-5$ to $5e-4$ over 10000 steps. Then we use a cosine schedule to decay the learning rate to 0 over 1 million steps. We find our networks converge within 200k steps.

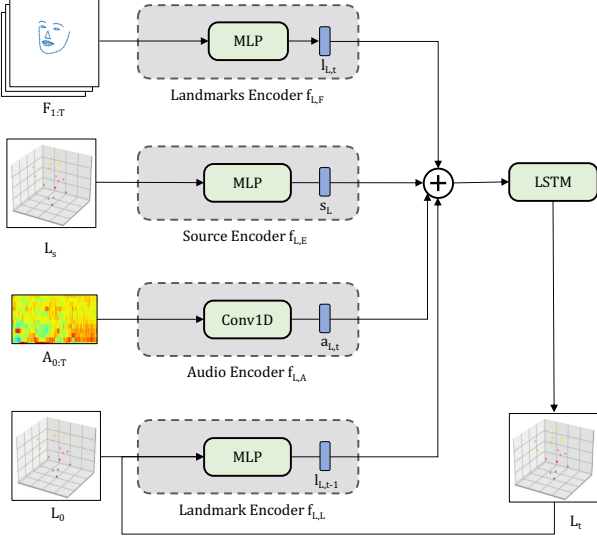


Fig. 9. **Landmarks2Latents architecture.** Landmarks2Latents is the final autoregressive component of SPACE. It is trained to transform facial landmarks to self-supervised keypoints learned by the talking head synthesis model. Landmarks2Latents takes as input the full face landmarks aligned with the image, which in turn are aligned with the keypoints. It also uses the input speech as an additional source of information. At a particular time step, the model takes the previous keypoints encodings, and the current landmarks and audio features.

B. Data preprocessing

Since VoxCeleb1 and VoxCeleb2 datasets were captured in the wild, there is a large variation in the quality of videos due to different recording devices, camera viewpoints, head poses, occlusions, *etc.* We found it essential to filter the data to remove outliers and obtain a clean subset for training our models. First, we remove any video where a state-of-the-art face detection algorithm fails to detect a face in any of its frames. This is because we assume facial landmarks are available for all the frames in a video. If a face detector cannot detect a face in a frame, then we will either fail to extract the facial landmarks from the frame or have unreliable facial landmarks from the frame. Neither is desirable. We also remove a video with a large temporal difference between facial landmarks in any pair of its neighboring frames. Such a case often occurs in an edited video where the view transitions from one subject to another. In addition, we remove videos with large head rotation and camera zoom motion. We track the scale of the face landmarks and use the difference between the minimum and maximum scales as the criterion for filtering. We also remove videos with head rotation greater than 45 degrees along any axis. We use a pre-trained hand tracker to remove any videos where hands are detected.

For the Ryerson and MEAD datasets, we use a face detector to crop and resize the full size videos to 512×512 .

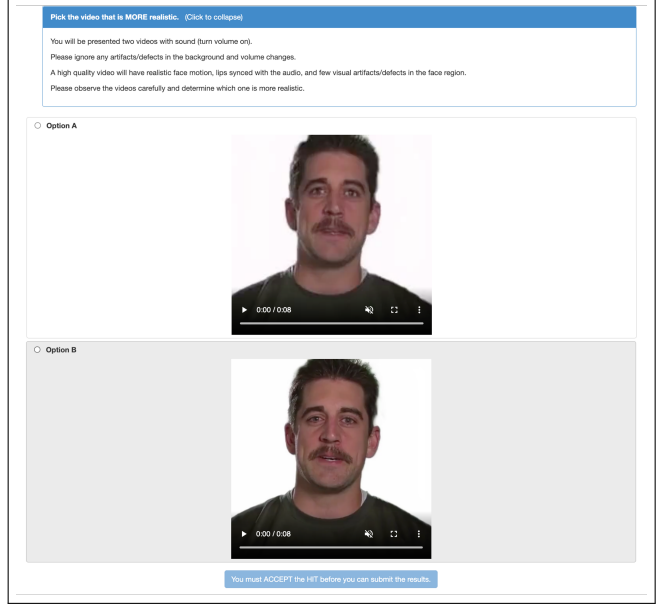


Fig. 10. **Amazon Mechanical Turk interface.** We ask users to choose the more realistic video, with a focus on the face and mouth region.

MEAD contains videos from various camera perspectives: front, up, down, left 30 degrees, right 30 degrees, left 60 degrees and right 60 degree. We discard the left 60 degrees and right 60 degrees videos.

B.1. Test Set details

Our test set was chosen to have a diverse set of initial poses. Using head pose estimation, we found that our test set has $\sim 45\%$ data with rotations less than 15 degrees, $\sim 35\%$ data with rotations between 15 and 30 degrees, and $\sim 20\%$ data with rotations between 30 and 45 degrees.

C. Evaluation

Fig. 10 shows the interface of the user study. As described in the text, a user is presented with a pair of videos from randomized models. Each pair consists of videos generated using same reference image and audio with different models. The user is asked to pick the video they prefer in terms of face and mouth motion.