# Iterative Text-based Editing of Talking-heads Using Neural Retargeting

XINWEI YAO, Stanford University

OHAD FRIED, The Interdisciplinary Center Herzliya

KAYVON FATAHALIAN, Stanford University
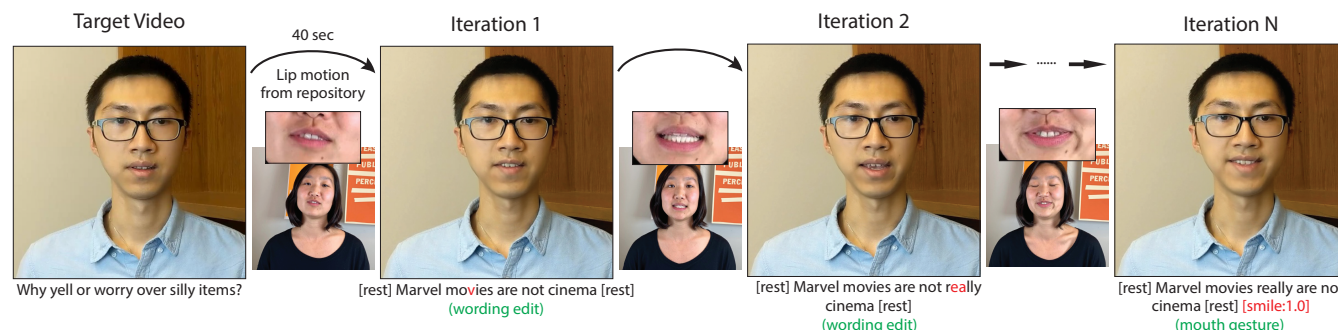
MANEESH AGRAWALA, Stanford University

Fig. 1. Our iterative text-based tool for editing talking-head video takes 2-3 minutes of a target video as input and is designed to support an iterative editing workflow. On each iteration the user might edit the wording of the speech (itr. 1 and 2), refine mouth motions if necessary to reduce artifacts, manipulate the performance by inserting mouth gestures (itr. N) or change the overall speaking style. Unlike previous techniques that require hours our tool takes about 40 seconds to generate each iteration, making it practical for users to explore a variety of different edits as they iterate. Our approach is to retarget lip motion from a repository of source actor video to the target actor. The frame shown for each iteration corresponds to the red edit text/gesture below the frame.

We present a text-based tool for editing talking-head video that enables an iterative editing workflow. On each iteration users can edit the wording of the speech, further refine mouth motions if necessary to reduce artifacts and manipulate non-verbal aspects of the performance by inserting mouth gestures (e.g. a smile) or changing the overall performance style (e.g. energetic, mumble). Our tool requires only 2-3 minutes of the target actor video and it synthesizes the video for each iteration in about 40 seconds, allowing users to quickly explore many editing possibilities as they iterate. Our approach is based on two key ideas. (1) We develop a fast phoneme search algorithm that can quickly identify phoneme-level subsequences of the source repository video that best match a desired edit. This enables our fast iteration loop. (2) We leverage a large repository of video of a source actor and develop a new self-supervised neural retargeting technique for transferring the mouth motions of the source actor to the target actor. This allows us to work with relatively short target actor videos, making our approach applicable in many real-world editing scenarios. Finally, our refinement and performance controls give users the ability to further fine-tune the synthesized results.

CCS Concepts: • **Computing methodologies** → **Motion processing**; **Computational photography**; *Reconstruction*; *Graphics systems and interfaces*.

Authors' addresses: Xinwei Yao, xinwei.yao@cs.stanford.edu, Stanford University, Department of Computer Science; Ohad Fried, ohad.fried@post.idc.ac.il, The Interdisciplinary Center Herzliya, Department of Computer Science; Kayvon Fatahalian, kayvonf@cs.stanford.edu, Stanford University, Department of Computer Science; Maneesh Agrawala, maneesh@cs.stanford.edu, Stanford University, Department of Computer Science.

Additional Key Words and Phrases: text-based video editing, talking-heads, phonemes, retargeting

## 1 INTRODUCTION

Tools for editing talking-head video using transcripts have made it possible to easily remove filler words, emphasize phrases, correct mistakes, and try different wordings of the speech [Berthouzoz et al. 2012; Fried et al. 2019; Suwajanakorn et al. 2017; Thies et al. 2020]. Many of these tools can synthesize high-quality results that closely match the appearance of the unedited video. Such tools have the potential to enable a variety of post-capture editing applications including re-phrasing dialogue in a film scene, dubbing commercials to a new language, developing dialogue for a conversational video assistant, and fixing wording errors in an online lecture.

Yet, current transcript-based video editing tools are impractical for use in many real-world editing scenarios for four main reasons.

***(1) Slow feedback loop hinders iterative editing.*** Synthesizing the edited result at high-quality is often extremely slow. For example, while viewers report that Fried et al.'s [2019] results appear very realistic, their approach takes hours to generate a few seconds of edited video. The slow feedback loop — time between specifying an edit and seeing the result — significantly hinders iterative editing (e.g. trying different phrasings of dialogue).

***(2) Require hours of target talking-head video.*** To produce realistic results, many of these tools [Fried et al. 2019; Suwajanakorn et al. 2017] require hours of video of the target talking-head actor.

Some tools further require the actor to speak a set of specialized phrases (e.g. TIMIT corpus [Garofolo et al. 1993]). In practice however, many video editing projects lack access to such large amounts of target actor video.

*(3) Missing controls for refining results.* None of these editing tools provide controls for manually refining the lip motions of the synthesized results, making it impossible to fix objectionable artifacts these editing tools sometimes generate (e.g. mouth doesn't fully close on \m, \b, \p phonemes).

*(4) Missing controls for adjusting non-verbal performance.* None of these editing tools include controls for adjusting the target actor's non-verbal performance by inserting mouth gestures (e.g. a smile) or changing the overall speaking style (e.g. mumbling, energetic).

In this work we present an iterative talking-head video editing tool that explicitly addresses all four of these issues. While our approach builds on the high-quality synthesis technique of Fried et al. [2019], we make several new contributions. We significantly reduce the time required to synthesize video (from hours to about 40 seconds for a 6 word edit) by developing a fast algorithm for searching the source repository for the desired lip motions. We lower the data requirement on the target actor video (2-3 minutes are usually enough) by leveraging a large repository of video from a source actor and use a new self-supervised neural retargeting technique to transfer their lip motions to the target actor. We provide controls to refine results by allowing users to smooth over jumpy transitions and force mouth closure on the results of the automated synthesis pipeline. Finally, we enable insertion of non-verbal mouth gestures with the same text interface, as well as controls to switch between different speaking styles by using a version of the source repository with the desired style.

As shown in Figure 1 our tool enables an iterative editing workflow. Given a short video of the target actor, the user can edit the transcript and our tool synthesizes the corresponding video in under a minute. The user can inspect the feedback and further adjust wording, refine the lip motions and/or insert mouth gestures and quickly see how the adjustment affects the synthesized video. Note that our work focuses on generating video from text; to obtain the corresponding speech audio, we rely on either having access to the actor speaking the new content (e.g. from a prerecorded library of the actor's speech or recorded by the actor in real-time during editing), text-to-speech voice synthesis [van den Oord et al. 2016] or voice cloning [Jia et al. 2018; Kumar et al. 2019].

We demonstrate a variety of iterative editing sessions facilitated by our tool and we conduct user studies which show that our synthesized results are rated as "real" for 56.2% of the sentence-long edits and for 64.9% of the phrase-long edits – slightly better than the previous state-of-the-art approach of Fried et al. [2019]. Together these results suggest that our algorithm provides the speed, data efficiency and controls necessary for a practical iterative editing workflow while maintaining high-quality synthesis results.

## 2 RELATED WORK

*Video-driven talking-head synthesis.* A common approach to synthesizing a talking-head video is to use a "driving" video from a different actor that has the desired motion, expression and speech, and transfer those elements to the primary talking-head. Early attempts used facial landmarks from a video to retrieve frames of a different person and play them back directly [Kemelmacher-Shlizerman et al. 2010] or after warping [Garrido et al. 2014]. Opting for a lower data requirement, several approaches synthesize video given only one or a few photos of the target person, either by morphing and blending [Averbuch-Elor et al. 2017] or using neural networks [Geng et al. 2018; Pumarola et al. 2019; Wiles et al. 2018; Zakharov et al. 2019]. These methods are successful in producing short expression videos, but are less convincing for full sentences. Several approaches use a tracked head model, to decouple properties (e.g. pose, identity, expressions) to produce convincing results [Garrido et al. 2015; Kim et al. 2019, 2018; Thies et al. 2016; Vlasic et al. 2005]. We similarly use a tracked head model to decouple such properties. All of these previous methods require a driving video to specify the desired output head motion and expression. In this work we specify those properties via text, which is often a simpler, lower-cost interaction.

*Voice-driven talking-head synthesis.* Another approach for talking-head synthesis is to drive it with voice. The pioneering work of Bregler et al. [1997] creates talking-heads through a combination of alignment and blending, and was improved upon in various followups [Chang and Ezzat 2005; Ezzat et al. 2002; Liu and Ostermann 2011]. Others have used human voice-driven synthesis to dub non-humans [Fried and Agrawala 2019]. Several methods synthesize a talking-head given only one or a few video frames in addition to the voice track [Chen et al. 2019; Chung et al. 2017; Song et al. 2019; Vougioukas et al. 2018, 2019; Zhou et al. 2019]. However, the result is a fixed frame with a moving inner-face region, or a tightly cropped head, and is easily distinguishable from realistic video. Suwajanakorn et al. [2017] demonstrate that using a large repository of video (17 hours) can produce convincing synthesis results. In work concurrent to ours, Thies et al. [2020] produce video from speech. We compare to their results in Section 5.3. None of these voice-drive methods provide refinement and performance controls that are essential for iterative editing.

*Text-driven talking-head synthesis.* Most related to our work are methods that perform text-based video editing and synthesis. Wang et al. [2011] synthesize a talking-head, and allow control over facial expressions, but the head is floating in space and is not part of a photorealistic video. Mattheyses et al. [2010] synthesize audio-visual speech from text, but the resulting videos have no head motion making them unrealistic. Berthouzoz et al. [2012] can edit talking-head video by cutting, copying and pasting transcript text. However, they do not allow synthesis of new words to change phrasing or fix flubbed lines. ObamaNet [Kumar et al. 2017] synthesizes both audio and video from text, using a large dataset of 17 hours of the president's speeches. While predominantly an audio-based method, Thies et al. [2020] also show text-based results by incorporating a text-to-speech system. The work of Fried et al. [2019] most closely resembles ours, but it requires over 1 hour of target video and takes
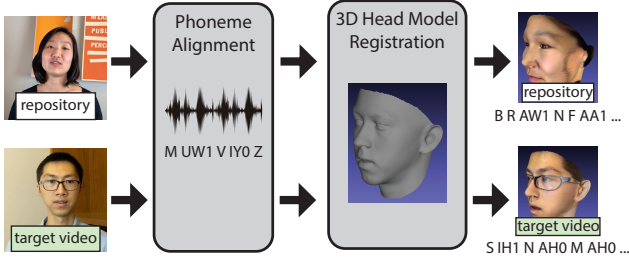
Fig. 2. Our preprocessing pipeline annotates both the source repository and the target video with phonemes and registers a 254-parameter 3D head model to each frame of each video.

hours to produce a result, while our tool requires 2–3 minutes of target video and produces results in about 40 seconds. We compare our results to both Thies et al. and Fried et al. in Sections 5.3 and 5.4 and find that the quality of our results is similar to both of these techniques. Moreover we provide refinement and performance controls that are missing in all previous text-driven talking-head synthesis tools, but are critical for a practical video editing tool.

## 3  METHOD

Given a short talking-head video of a *target actor* (often 2-3 minutes in length), and an edit of the video transcript, our system synthesizes new video of the target actor matching the edit. An edit is specified as a replacement of one continuous sequence of words in the original transcript with a new sequence of words. Since a short target actor video is unlikely to contain all the lip motions necessary to convincingly synthesize the sequence of phonemes in the edit, we leverage a large *repository* of video from a different, *source actor*. Specifically, we pre-capture an hour of a source actor speaking the TIMIT corpus [Garofolo et al. 1993] which includes the most common phoneme combinations (coarticulations) in English and we retarget their lip motions to the target actor during synthesis.

Our approach for quickly synthesizing the edited result is based on the approach of Fried et al. [2019] but involves several critical modifications. As in Fried et al., our *preprocessing pipeline* (Figure 2) annotates both the repository and target videos with phonemes and registers a parametric 3D head model to the face in each frame of each video. Our *synthesis pipeline* (Figure 3) provides a new, fast phoneme search algorithm that finds subsequences of phonemes in the source video that match the desired edit. It then stitches together the corresponding parameters of the 3D head model for the source actor across subsequence boundaries to smooth the lip motions. We introduce a new self-supervised neural retargeting step that adapts the parameters representing the lip motion of the source actor to those of the target actor and blend the resulting parameters into the target video. Finally we render photorealistic frames from the parameters using neural rendering [Tewari et al. 2020].

We briefly summarize how we adapt each step in Fried et al.'s pipelines to our problem in Sections 3.1 and 3.2. We then present the details of our new algorithms; fast phoneme search and stitching in Section 3.3 and neural retargeting algorithm in Section 3.4. In Section 3.5, we describe the iterative refinement and performance controls enabled by our approach.

### 3.1  Preprocessing Pipeline

Our preprocessing pipeline annotates each frame of the repository and the target videos in two main steps, (1) phoneme alignment and (2) registration of a parametric 3D head model. These resulting phoneme and face parameter annotations are used by our synthesis pipeline to establish correspondences between the target video and source repository. Note that the repository is only annotated *once* and the resulting *annotated repository* is then bundled as part of the system. In contrast, the target video must be annotated each time a new target video is given as input.

*Phoneme Alignment.* The phoneme alignment step takes as input a video (repository or target) paired with its text transcript, and computes the identity and timing of the phonemes in the video. Specifically, we use P2FA [Rubin et al. 2013; Yuan and Liberman 2008] to convert the transcript into phonemes and align them to the audio speech track of the video. This produces an ordered sequence $V = (v_1, \ldots, v_n)$ of phonemes, where each phoneme $v_i$ contains its name, start time and end time. If the transcript is not available, we can obtain one using a transcription service such as Google Cloud Speech-to-Text [2020a], or rev.com [2020].

*3D Head Model Registration.* We fit a parametric head model [Blanz and Vetter 1999; Thies et al. 2016] to each frame of video using a monocular head tracker [Garrido et al. 2016]. At every frame, the fitted model includes 80 parameters for 3D facial geometry, 80 for facial reflectance, 3 for head pose, 27 for scene illumination and 64 for face and lip expressions. In the fitting procedure we hold the facial geometry and reflectance parameters constant across all the frames of the same actor, but we allow the pose, illumination and expression to vary across time.

### 3.2  Synthesis Pipeline

Our synthesis technique is based on matching phonemes in the edit to phonemes in the repository. Therefore, we first convert the input text of the edit from words into a sequence of phonemes $W = (w_1, \ldots, w_m)$, where each phoneme contains its name, start time and end time. Specifically, we convert the edit into audio using either text-to-speech voice synthesis [goo 2020b; van den Oord et al. 2016] or voice cloning techniques [lyr 2020; Kumar et al. 2019], and then apply P2FA [Rubin et al. 2013; Yuan and Liberman 2008] to time-align the resulting speech to the phonemes of the edit. Note that our synthesis algorithm only uses the timing of the phonemes and does not use any other aspect of the synthesized speech audio signal. If the user has access to the audio of the target actor saying the new content (either from a prerecorded library of the actor's speech or recorded by the actor in real-time during editing) they can run our fast synthesis pipeline with the timings obtained from the real-voice recordings.

*Fast Phoneme Search and Stitching.* The fast phoneme search and stitching step is designed to quickly find the best subsequences of phonemes in the repository and then stitch together the corresponding expression parameters of the source actor 3D head model, in order to produce the the edit $W$ we wish to synthesize. Our new algorithm operates three orders of magnitude faster than Fried et al. [2019] and finds the best repository subsequences in
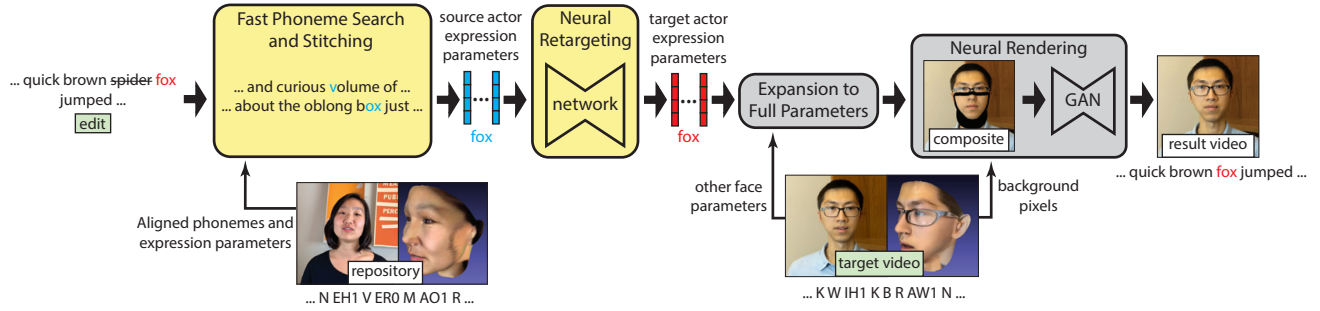
Fig. 3. Our synthesis pipeline adapts the pipeline of Fried et al. [2019] by introducing a fast phoneme search and stitching step (yellow), and a self-supervised neural retargeting step (yellow). The input to our pipeline is a target video (green) and an edit (green) – here changing the word "spider" to "fox". Our fast phoneme search finds phonemes in the repository that visually match the desired edit – here the "v" in volume and the "ox" in box. We then stitch together the corresponding facial expression parameters of the 3D head model for the source repository actor, and use a novel neural retargeting model to translate those parameters into those for the target actor. Next, we expand the retargeted expression parameters to the full face parameters (e.g. pose, illumination) for the target actor. Finally we render photorealistic frames from the parameters using a neural rendering approach that first composites the lower part of the face rendered from the 3D head model, with background pixels from the original target video and then uses a generative adversarial network (GAN) to map the composites to photorealistic frames.

seconds rather than hours. We present this fast algorithm in detail in Section 3.3.

*Neural Retargeting.* The retargeting step converts a sequence of expression parameters for the source actor into those for the target actor. We introduce a learned retargeting model, trained in a self-supervised manner from corresponding pairs of repository and target video sequences and transforms the expression parameters as detailed in Section 3.4. The result is a sequence of target actor face expression parameters corresponding to the edit.

*Expansion to Full Parameters.* Next, we combine the synthesized target actor expression parameters with geometry, reflectance, pose and illumination parameters from the input target video to produce a sequence of full face parameters for the target actor corresponding to the edit. Specifically, we take an interval of frames around the edit location in the target video, retime it to account for the duration of the edit, and use the geometry, reflectance, pose and illumination parameters from the retimed interval.

*Neural Rendering.* The neural rendering step takes the sequence of full face parameters for the target actor and first generates a composite image in which the lower face region is a rendering of the 3D head model, while the upper part of the head and the surrounding background are from the original target video, but retimed to match the length of the edit. It then uses a GAN trained on the target video to complete the image-to-image translation from composite image to photorealistic frame.

### 3.3 Fast Phoneme Search and Stitching

Our synthesis pipeline takes an edit $W$ specified as a sequence of phonemes with timings $(w_1, \ldots, w_m)$ and starts by finding matching subsequences of video in the source repository $V$, that can be combined to produce $W$. More precisely, we partition the edit $W$ into phoneme subsequences $(W_1, W_2, \ldots, W_k)$ and for each subsequence $W_i$ find its best match $V_i$ in the repository $V$. Fried et al. [2019] use a brute-force method that considers all possible partitions $\text{split}(W)$ of the edit $W$, and all possible matches with subsequences of $V$ to

find $(V_1, V_2, \ldots, V_k)$ that minimizes the objective:

$$C(W, V) = \min_{\substack{(W_1, W_2, \ldots, W_k) \in \text{split}(W) \\ (V_1, V_2, \ldots, V_k)}} \sum_{i=1}^{k} C_{\text{match}}(W_i, V_i) + C_{\text{len}}(W_i)$$

where $C_{\text{match}}$ is a custom Levenshtein edit distance [1966] between two phoneme subsequences that takes into account the phoneme label, the viseme label and the timing difference, and $C_{\text{len}}$ penalizes short subsequences. In order to find subsequences that transition well at their endpoints, during the search, we expand each subsequence $W_i$ by a single phoneme on either end. Thus, adjacent subsequences overlap by two *context* phonemes. We find that this new *context expansion* approach better captures co-articulation effects between the subsequences, than the algorithm of Fried et al. which does not use context expansion (see user studies in Section 5.4).

We further modify Fried et al.'s search algorithm in three key ways to obtain a speedup of over 3 orders of magnitude: (1) we propose a fast alternative to to the Levenshtein distance, (2) we reduce the size of the search space on $W_i$ and, (3) given $W_i$, we use a viseme-based indexing scheme to quickly find the optimal $V_i$. Finally, we stitch together source actor expression parameters corresponding to the $V_i$s to produce a single coherent sequence of expression parameters.

*(1) Fast alternative to edit distance.* The full Levenshtein edit distance allows substitution, insertion and deletion of phonemes when computing $C_{\text{match}}$. However, we have observed that when the matching subsequences between the edit $W$ and the repository video $V$ contain phoneme insertions or deletions, the final synthesized video appears out-of-sync with the audio; it either contains extraneous mouth motions due to phoneme insertion, or it misses mouth motions due to deletion. In practice we find it is beneficial to disallow insertions and deletions and only allow phoneme substitutions. Given a subsequence $W_i$ of the edit $W$, this approach forces $C_{\text{match}}$ to only consider subsequences $V_j$ of the repository $V$ that contain the same number of phonemes as $W_i$. We can therefore replace the Levenshtein distance with the sum of element-wise substitution

cost which requires linear time in the number of phonemes rather than the quadratic time required for computing the full Levenshtein distance [Wagner and Fischer 1974].

*(2) Reduce search space for partitioning.* The brute force search considers all possible partitions of $W$ into $(W_1, \ldots, W_k)$. But, an extremely long edit subsequence $W_i$ is unlikely to have a good match with a repository subsequence $V_i$. Thus, we can reduce the search space of possible partitions by capping the maximum length of the $W_i$'s to $L$. In our experience, 99% of the matches found by the brute force search are of length 6 or less, and we therefore set $L = 6$. This approach reduces the number of partitions to search. More importantly, it typically reduces the number of distinct $W_i$'s we need to consider by over an order of magnitude, especially when $W$ the full edit sequence is itself very long.

*(3) Viseme-based index to search repository.* For each edit subsequence $W_i$ we consider in our search space, we must find the optimal $V_i$ in the repository with respect to $C_{\text{match}}$. Instead of checking all possible subsequences in the repository, we impose an additional constraint on $V_i$ that restricts the set of $V_i$ we consider to only the most likely match candidates and allows us to build an index structure on the set of $V_i$ to retrieve the likely candidates quickly.

As in Fried et al., our $C_{\text{match}}$ cost function considers phonemes to match when they appear visually similar – that is, their corresponding visemes match. By imposing the restriction that $V_i$ start with the same viseme n-gram as $W_i$ we can pre-compute an index for the repository using viseme n-grams as the key and the location of the n-gram in the source repository video as the value. At search time, we look up all possible candidate $V_i$'s using this index and only compute the $C_{\text{match}}$ for them. While this indexing approach speeds up the search significantly it also reduces the space of subsequence matches the search considers. In general, the longer the n-gram key, the stronger the reduction and the more likely it is that no good match will be found. In practice, we find that using a bi-gram index best balances this trade-off between search speed and result quality.

*Stitching.* After we find the best matching phoneme subsequences $V_1, \ldots, V_k$ from the repository, we look up the expression parameters of the source actor's 3D head model corresponding to each phoneme, and linearly re-time them to match the phoneme durations specified in the edit. We then stitch together adjacent subsequences by first linearly blending the the expression parameters across the overlapping context phonemes and then applying a Gaussian filter over a window of 4 frames around the transition boundaries between the subsequences to further smooth the transition.

We have found in practice, that errors in tracking expression parameters, timing misalignments and our linear blending can sometimes fail to fully close the mouth of the 3D head model at the beginning of \m, \b, and \p phonemes. However, proper mouth closures for these phonemes is crucial for producing perceptually realistic results [Agarwal et al. 2020]. We therefore force the desired mouth closure by linearly blending in a closed-mouth expression from the repository at the beginning of all \m, \b, and \p phonemes with a default length of 2 frames. Note that the closed-mouth expressions are manually annotated once during repository preprocessing and
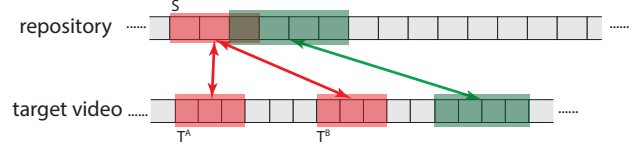


Fig. 4. To automatically build training data for our neural retargeting model, we consider a sequence of phonemes $S$ (red) in the repository and find up to two best matches $T^A$ and $T^B$ (red) in the target video. We then find matches for the next repository sequence (green) starting at the last phoneme of the previous sequence. The overlap allows us to capture the phoneme transition out of the final phoneme in the red subsequence.

we automatically use the one closest to the parameter values at the point of insertion.

Overall, the fast search strategy reduce the runtime of the phoneme search process by three orders of magnitude compared to the brute-force approach. It takes around 5 seconds to find and stitch snippets $V_1, \ldots, V_k$ for an edit $W$ of 20 phonemes in an hour-long repository.

## 3.4 Neural Retargeting

Given a stitched-together sequence of face expression parameters for the source actor in the repository, the goal of retargeting is to generate a matching sequence of expression parameters for the target actor. We have developed a self-supervised neural network model for retargeting and in Section 5.1 we show that using a neural network for retargeting produces higher-quality results than baseline methods such as directly copying source actor expression parameters to the target actor, or applying a linear retargeting model.

*Self-supervised training data.* To train our retargeting network we require sequences of expression parameters for the source and target actor that correspond to one another with respect to their mouth motions. Assuming that uttering the same sequence of visemes will produce similar mouth motions, we automatically construct corresponding pairs of training data by finding the longest matching sequences of phonemes between the source repository video and the target video, as follows.

Since we apply our retargeting model to a stitched-together subsequences of phonemes that can come from anywhere in the source repository, we would like the training sequences to cover as much of the repository as possible. Therefore, we start from the first phoneme $s_1$ in the repository, and find the longest sequence $S = (s_1, \ldots, s_k)$ for which there is at least one corresponding sequence $T = (t_1, \ldots, t_k)$ in the target video where $t_i$ and $s_i$ belongs to the same viseme group (i.e. phonemes that require the same lip expressions are in the same viseme group). We take up to two best matches in the target video with respect to $C_{\text{match}}$ score defined in Section 3.3, and call them $T^A = (t_1^A, \ldots, t_k^A)$ and $T^B = (t_1^B, \ldots, t_k^B)$. We add the pairs $(S, T^A)$ and $(S, T^B)$ to the set of phoneme sequences in correspondence, and continue the scan through the source repository at $v_k$, until we finish scanning through the entire repository for subsequence matches (Figure 4). To quickly find the best $T^A$ and $T^B$ sequences in the target video, we apply fast search techniques discussed in Section 3.3. Finally, to convert each resulting phoneme sequence
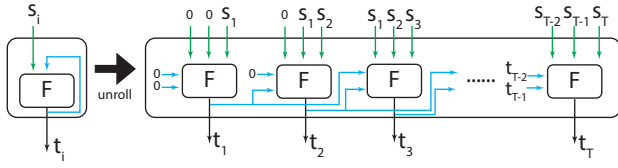
Fig. 5. The retargeting network with the unrolled recurrent unit $F$, which looks two frames back at each time step and uses the previous two outputs as its state. The diagram shows how the RNN processes the input of $T$ frames of source actor parameters $(s_1, s_2, \ldots, s_T)$ and produces the prediction for $T$ frames of target actor parameters $(t_1, t_2, \ldots, t_T)$.

pair into a parameter sequence pair, we linearly retime the target video phoneme interval to the corresponding repository interval and similarly interpolate the expression parameters. The retargeting model is trained once for each new target video.

*Neural Network Architecture.* We employ a recurrent neural network (RNN) manually unrolled for $T$ time-steps to encode the temporal dynamics of the facial expressions and regress from source parameters to target parameters (Figure 5). The resulting network takes as input $T$ frames of 64 expression parameters for the source actor, and outputs $T$ frames of 64 parameters for the target actor. At the core of the network is the recurrent unit that is made up of 3 fully connected layers of 1024 nodes with relu activation. The inputs to the recurrent unit are 64 expression parameters of the current time-step, as well as the parameters and outputs of the recurrent unit for the previous $H$ time-steps. To facilitate learning of deviations from the identity transformation, the output layer of the recurrent unit with 64 nodes and no activation produces the residual values that are added to the input source parameters element-wise to obtain the prediction for the target parameters for the current time-step. We zero-pad the unavailable inputs at the first $H$ time-steps. We empirically found that setting $H = 2$ and $T = 7$ produced high-quality retargeting results.

*Loss function.* Our loss function is a linear combination of a data term and a temporal regularizer with regularization weight $\lambda$:

$$\mathcal{L} = \frac{1}{T} \sum_{i=1}^{T} \|F(s_i) - t_i\|_1 + \lambda \|F^{(2)}(s_i)\|_2$$

where $s_i, t_i$ are 64-dimensional vectors representing the $i$th time-step of the source and target parameters respectively, $F(s_i)$ is a 64-dimensional vector of the predicted target actor parameters and $F^{(2)}(s_i)$ is a 64-dimensional vector that is the second temporal difference vector, or acceleration, of the predicted values. Empirically we found that using an $L_1$ norm for the data term significantly outperforms $L_2$ by generating more expressive motions and better preserving mouth closures. The temporal regularization term is needed to make the network predict temporally stable parameters.

*Hyperparameters and training.* Since the network takes a fixed-size input of $T$ frames, we run a sliding window on each matching parameter sequences to obtain the training examples. Experimentally we set the temporal window $T = 7$ frames. For training we set

$\lambda = 10$, and dropout rate at 25%, 50% and 25% for the three layers in the recurrent unit, respectively. To train the network we use stochastic gradient descent with the Adam solver [Kingma and Ba 2015] and set an initial learning rate of 0.0002 with an exponential decay rate of 0.5. We employ gradient clipping [Pascanu et al. 2013] to avoid exploding gradients. We train the network with minibatch size 100 and training typically converges within 100 epochs.

*Inference.* At inference time, we convert a sequence of source actor expression parameters into target actor parameters. Since our retargeting model accepts fixed-size $T$ frames of input and produces $T$ frames of output, we run a sliding window of length $T$ over the new sequence of source actor expression parameters at inference time. Each frame is covered by exactly $T$ such sliding windows. In order to obtain a more temporally stable output, at each frame we average the $T$ outputs produced by those $T$ sliding windows as the final output of the frame. The result is a synthesized sequence of target actor expression parameters that animate the face to speak the new content of the edit with the desired timing. We then proceed to expand these expression parameters into parameters for the whole head and use neural rendering to generate the video frames (Section 3.2).

Training our retargeting model typically requires 2–3 minutes of target actor video speaking arbitrary speech to produce high-quality synthesis results. Retargeting allows our tool to leverage the large repository of controlled source actor video (speech consists of TIMIT sentences) to generate the target actor lip motions and opens our tool to many practical applications where large amounts of controlled target actor video are not available.

## 3.5 User Controls

This speed of our synthesis pipeline opens the door to interactive user controls for iteratively refining the edit and further manipulating the the facial performance.

*Refinement Control: Smoothing jumpy transitions.* Our synthesis pipeline stitches together different subsequences of expression parameters from the source repository by smoothing over a window of 4 frames around the transition boundary (Section 3.3 Fast Phoneme Search and Stitching). At times however, some transitions may still appears jumpy even after this smoothing. We allow the user to inspect the result and further refine it by manually specifying the interpolation radius (in number of frames) at user-specified transition boundaries to better smooth out visibly jumpy transitions.

*Refinement Control: Adjusting mouth closure.* As noted in Section 3.3 mouth closure on \m, \b and \p phonemes is crucial for the mouth motions to appear realistic. Thus our stitching procedure automatically inserts 2 closed mouth frames at the beginnings of these three phonemes to ensure the mouth closes correctly for them. We further allow users refine any synthesized result by extending (or reducing) the length of the inserted closed mouth frames.

*Performance Control: Inserting mouth gestures.* Users can also insert non-verbal mouth gestures (e.g. a smile) into an edit. To enable such performance control we manually annotate mouth gestures including *rest, closed-mouth smile, regular teeth-showing smile, big*

*open-mouth smile, sad, scream, mouth gesturing left* and *mouth gesturing right,* in the repository video. These segments can then be retrieved by our fast phoneme search just like any other phoneme annotation.

Since the annotations are on the repository, this manual annotation only needs to be done once during repository preprocessing. Note however, that users do not label the target video with these mouth gestures and our retargeting network is never explicitly trained with corresponding pairs of mouth gesture frames between the repository and target videos. Nevertheless, we have found that our retargeting network is able to generalize to unseen expressions and produce good quality expression parameters for the target actor.

With these annotations, the user can add special *mouth gesture directives* like [smile] anywhere in their edit of the transcript, and our tool constructs a "generalized phoneme" edit sequence $W$ that contains phonemes and such directives. Any mouth gesture that appears in $W$ is given a default duration of 0.5 seconds that the user can override with an explicit duration e.g. [smile:1.5s]. We employ a special substitution cost in $C_{\text{match}}$ described in Section 3.3 for "gesture phonemes" that is set to infinity for a non-match to ensure that we retrieve the correct "phoneme" match for the gesture. When there are multiple candidates, $C_{\text{match}}$ takes duration into account and picks the gesture with duration closest to the query. The rest of the editing pipeline (Section 3.2) is otherwise unmodified.

*Performance Control: Adjusting speaking style.* Our tool allows the user to select a different speaking style for the synthesized result by using a version of the repository with the desired style. In addition to the default repository which captures a "neutral" speaking style, we have recorded an "energetic" repository of our source actor with more pronounced mouth movements, and a "mumble" repository with significantly less mouth movements. Figure 6 (third row) shows frames from these alternative repositories.

Importantly, we do not have to retrain our neural retargeting model (Section 3.4) for each additional style repository. We train this retargeting model once using only the default neutral repository. We have found that our default retargeting model can extrapolate to retarget subsequences of source actor face parameters retrieved from other speaking style repositories of the same actor. Moreover, the other repository videos can be captured at different times, with different background and the source actor can even be wearing different clothes or have a different hairstyle. Thus, our tool generates videos with different speaking styles by using one of the alternative style repositories in the fast phoneme search step, but leaves the remainder of the synthesis pipeline unchanged.

### 3.6 Implementation Details

We implemented the fast phoneme search in Python and both our neural retargeting model and the GAN renderer in Tensor-Flow [Abadi et al. 2015]. The monocular head tracker and renderer are written in C++ with shader language extensions.

In preprocessing the repository and target videos, phoneme alignment takes one third of the video time, and face registration takes 110 ms per frame. It takes 30 minutes to generate training data for our neural retargeting model and another 30 minutes to train it on

| #words | #phonemes | #frames | search (s) | render (s) | total (s) |
|--------|-----------|---------|------------|------------|-----------|
| 1 | 4 | 24 | 1.51 | 9.96 | 12.39 |
| 3 | 15 | 49 | 2.67 | 12.94 | 20.30 |
| 6 | 25 | 72 | 5.32 | 17.74 | 28.95 |
| 8 | 39 | 105 | 7.57 | 24.33 | 37.97 |
| 10 | 49 | 134 | 11.60 | 31.19 | 50.60 |

Table 1. Runtime of our tool on a variety of edit lengths. Search time scales roughly linearly with the number of phonemes, and render time scales linearly with the number of frames. Even for long edits of 10 words, our system can generate video in approximately a minute.

one NVIDIA GTX 1080Ti. Training the GAN for neural rendering takes 17 hours on one NVIDIA Tesla V100.

In our synthesis pipeline, our fast phoneme search requires 5 seconds for a typical edit of 5 words containing 20 phonemes. Retargeting inference speed is 10K fps. Composite images are rendered at 12 fps and final GAN rendering takes 7 fps on two NVIDIA GTX 1080Ti. All together, a typical 5 word edit takes around 30 seconds for the full video generation (Table 1).

It should be possible to further reduce the feedback time by parallelizing our synthesis pipeline. Phoneme search could be distributed where each worker job is responsible for searching a fraction of the repository. Both parts of the neural rendering step – forming the composite images from target actor head parameters and applying the GAN to generate photorealistic frames – are parallelizable by distributing the frames. We have performed initial experiments on parallelizing the GAN rendering which is the main bottleneck in our pipeline. Distributing the GAN rendering across a cluster of 8 NVIDIA Tesla V100s achieved a rendering rate of 24fps, a 3.4x speedup from the original 7fps for this step. Note that this speed up rate includes image compression overhead. Overall this experiment cuts the end-to-end synthesis time from 40 to 20 seconds for a typical 8-word sentence. Similarly parallelizing the other parts of the pipeline and using sufficient hardware we believe that the end-to-end video generation feedback time could be reduced significantly. Streaming the frames as they are ready could also further reduce the latency from issuing an edit to seeing the first frames of the result, enabling real-time interactive editing sessions.

## 4 RESULTS

Figures 1 and 6 show examples of iterative text-based editing sessions for a variety of target videos including recordings of graduate students, YouTube video and a take from filming a dialog scene. We encourage readers to watch the videos in our supplementary materials to see how our text-based interface facilitate the iterative editing workflow used in each of these sessions.

*Session 1: Talking-head with glasses.* Our first session works on a 2.5 minute target video (Figure 1). The editor explores ways for the actor to parody Martin Scorsese and express that Marvel movies are not to be considered real cinema. They first synthesize "Marvel movies are not cinema" from a resting pose. Feeling it is too blunt, they slightly change the wording to "Marvel movies are not really cinema", and eventually settle on the firmer statement "Marvel movies really are not cinema". They then insert a smile at the end to soften the overall tone. Our tool is able to produce results with synchronized mouth motions at each step.

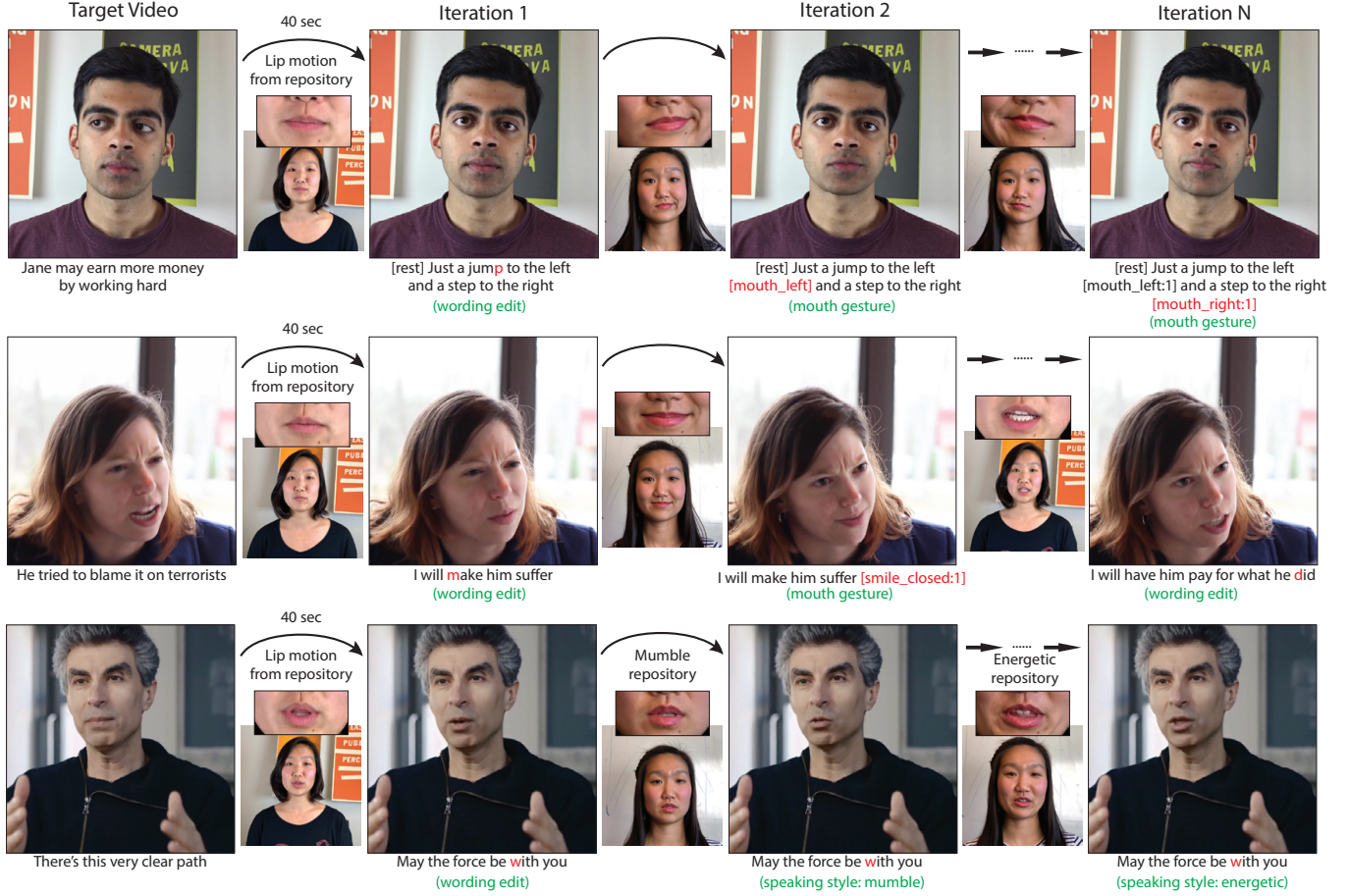Target Video    40 sec    Iteration 1    Iteration 2    Iteration N



Fig. 6. Editing sessions facilitated by our tool on target videos of 1 to 5 minutes in length. In these sessions our tool lets the editor change wording, insert mouth gestures, refine mouth motions and manipulate speaking style. Here we highlight a few of these edits in each session. Our tool finds the source actor mouth motions that conform to the specified edit operation and then retargets them to the target actor. Please see our supplemental materials and video for full video results as well as screen recordings of the edit sessions in our interface.

*Session 2: Talking-head with stubble.* Our second session works on a 3.5 minute target video (Figure 6 first row). The editor explores ways for the actor to give the instruction to start the time warp by jumping to the left then stepping to the right. They first synthesize the instruction phrase, then add gestures "[mouth_left]" and "[mouth_right]" to the corresponding location in the dialog. Next, feeling the gestures go too quickly, they lengthen the gestures to one second each. Our tool produces realistic video with mouth motions synchronized to the audio and the gesture directives.

*Session 3: Movie scene.* Our third session works on a target video of a single take from a dialogue scene (Figure 6 second row). The target video is a challenging one because it is only 1 minute long, and the actress speaks for only 30 seconds in the take. Our tool nevertheless is able to synthesize compelling results for this session. In the session, the editor prototypes ways for the actress to express her hatred towards the murderer of her sister's dog. They first try "I will make him suffer". Then for added creepiness, add a tight-lipped smile of 1 second duration at the end. Finally, they settle on a less hostile line instead. While the neural renderer struggles with the lack of data to produce images as sharp as those from the

previous two sessions, our tool is still able to produce realistic and synchronized mouth motions that give the user a good sense of how the scene would look with the alternative line and gesture.

*Session 4: Interview.* Our fourth session works on a 5 minute excerpt of a YouTube video (Figure 6 third row). The editor explores different delivery styles for the phrase "may the force be with you". They first synthesize the phrase with the default repository, then switch to a mumble style and finally to an energetic style. While the mouth movements match the audio, there is visibly less motion with the mumble style and more motion with the energetic style.

## 5 EVALUATION

To evaluate our tool we analyze the quality of the synthesized video as we vary the algorithmic methods (e.g. fast phoneme search, neural retargeting) used in our synthesis pipeline, and as we vary the data (e.g. length of target video, repository or edit) provided to the algorithm. We then compare the quality of our results to those of previous work. Finally we report on user studies that quantitatively evaluate the quality of our synthesized results. Unless otherwise

P  IY1  P  AH0

Fried et al. [2019]

Only Our Phoneme Context Expansion

Our Fast Phoneme Search and Stitching

that's all, people

Fig. 7. Comparison of phoneme search and stitching method. While the three methods usually yield visually indistinguishable frames, both the baseline method in Fried et al. [2019] and that with phoneme context expansion fail to close the mouth at the first "P" in people, whereas our fast phoneme search and stitching algorithm fully closes it.



Source Actor  Copy  Linear  Ours

Composite Image

Ground Truth

Composite Image
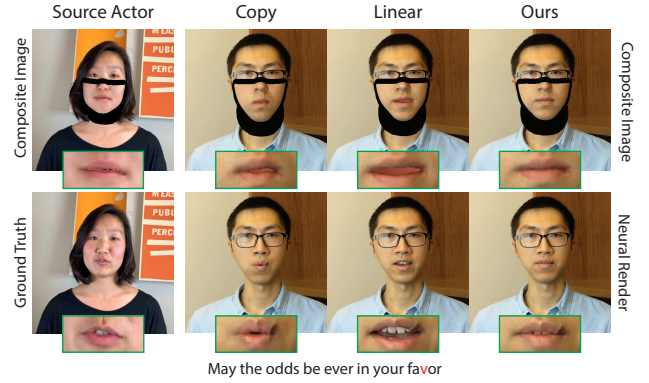
Neural Render

May the odds be ever in your fa**v**or

Fig. 8. Comparison of retargeting methods. Copying expression parameters yields an unnatural, rounded mouth shape on the target actor, while linear regression fails to close his mouth for the "v" sound in the word "favor". In contrast, our retargeting method achieves a good match in lip shape to the source actor and properly matches the shape required for the "v" sound.



30-second target  1-minute target  2.5-minute target  5-minute target

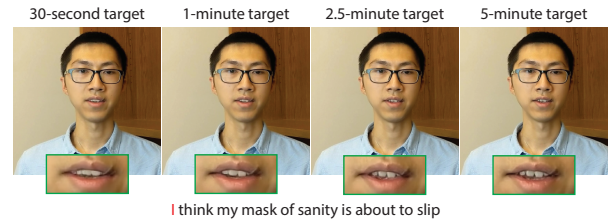I think my mask of sanity is about to slip

Fig. 9. With 2.5 minutes or more target video, our approach produces relatively sharp, realistic frames. With only 30 seconds or 1 minute of target video, the frames become noticeably blurry especially around the mouth and teeth (Zoom in at the front teeth to see the difference).

noted, the evaluations in this section use results from our automatic synthesis pipeline with no additional user refinement. Readers should refer to supplemental materials to evaluate the video results presented in this section.

## 5.1 Varying the Algorithmic Methods

*Comparison of phoneme search and stitching methods.* We compare the impact of using ablated versions of our fast phoneme search and stitching algorithm (Section 3.3) with the phoneme search method of Fried et al. [2019]. More specifically, because their synthesis pipeline does not include neural retargeting we treat Fried et al.'s approach as a baseline method and build two comparison pipelines. The first one adds only phoneme *context expansion* (Section 3.3) to the baseline stitching method. The second pipeline replaces their phoneme search and stitching method with the full version of our fast method (we call this version of the pipeline "Modified Fried et al. [2019]" in later comparisons). All three pipelines assume access to an hour of target video which serves as the repository. Figure 7 and videos in the supplemental materials show that results generated using our fast phoneme search and stitching method are often indistinguishable from the those generated by the baseline. The main differences that do appear are often subtle as our forced mouth closure on \m, \b, and \p phonemes reduces open mouth artifacts and our new context aware stitching (Section 3.3) across subsequence transitions yields smoother, less jittery lip motions. More importantly our method takes only seconds to run, which is three orders of magnitude faster than the hours required by baseline Fried et al. [2019], making it possible for the user to iterate on edits.

*Comparison of retargeting methods.* Our neural retargeting method (Section 3.4) transforms a sequence of source actor expression parameters to those of the target actor. We compare our method with two simpler baseline retargeting methods. The *copy* baseline directly copies the source actor parameters to the target actor. The *linear* baseline replaces our retargeting network with a linear model. More specifically, we manually chose 2.5 minutes of target video containing phrases that exactly matched phrases in the source repository and established a frame-to-frame correspondence by re-timing the target phonemes to match the lengths of the source phonemes. We then applied linear regression on the source and target parameter pairs to obtain the linear baseline model. Figure 8 and the supplemental materials show that our neural retargeting model produces the best results, while direct copying produces uncanny mouth shapes and the linear model often fail to close the mouth, causing out-of-sync lip motions.

## 5.2 Varying the Amount of Data

*Varying the length of the target video.* We examine the effect of different amounts of target video by applying our tool on 10 minute, 5 minute, 2.5 minute, 1 minute and 30 second subsets of a target video. More target video generally results in sharper images, higher
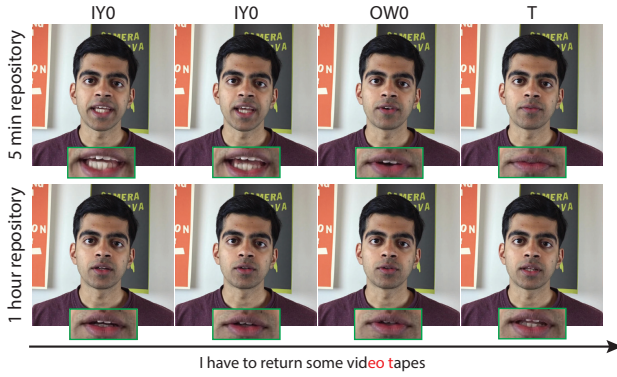
Fig. 10. Effect of small repository size. As shown here, the result from the 5-minute repository look less temporally stable as the transition from IY0 to OW0 is more drastic than the result from the full-hour repository. With 5-minute repository, the actor also has an incorrect closed-mouth at T because we do not have as good a selection of phoneme coarticulations with a small repository as we do with a full-hour one.
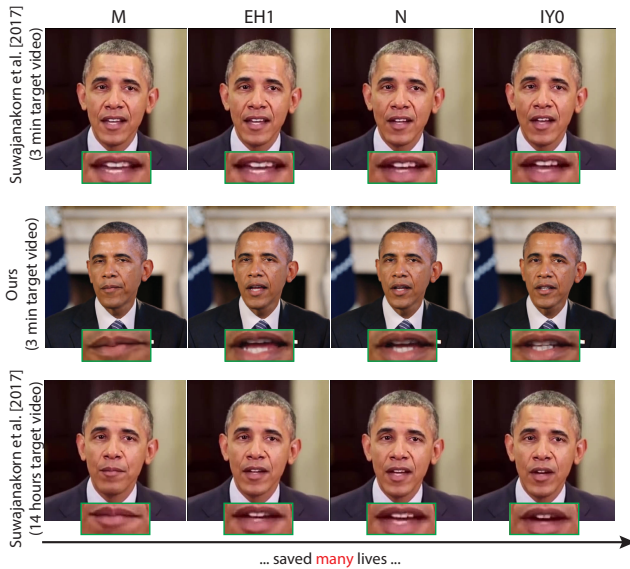


Fig. 11. Comparison to Synthesizing Obama [Suwajanakorn et al. 2017]. With 3 minutes of Obama video, Suwajanakorn et al. [2017] cannot give realistic mouth motions (top row). It produces similar frames throughout the word "many", and in particular fails to close the mouth at \m. Our method produces well-synchronized mouth shapes and closures with only 3 minute of target video (middle row), while Suwajanakorn et al. [2017] needs 14 hours of target video to produce a good result (bottom row).

quality mouth interiors and smoother mouth motions. But the difference is subtle when target video exceeds 5 minutes, and at 2.5 minutes the results remain plausible. With a 30 second target video however, although the results still have well-synchronized mouth motions, our neural renderer struggles and produces noticeably blurrier images, as shown in Figure 9 and videos in supplemental materials. Please zoom in on the front teeth to see the difference.

*Varying the amount of repository video.* We examine the effect of different amounts of repository video by applying our tool on



Fig. 12. Comparison to Fried et al. [2019]. Mouth motions from Fried et al. [2019] are less temporally coherent as shown from phoneme 'JH' to 'EY1', both with 3.5 minutes (top row) and with 1 hour (middle row) of target video. With 3.5 minutes, Fried et al. [2019] also produced an incorrect mouth closure during the second 'JH' frame. Our results (bottom row) have smoother mouth motions, and we use only 3.5 minutes of target video.

a 3.5 minute target video with 60 minute, 30 minute, 10 minute and 5 minute subsets of repository data. Generally a larger repository leads to better results. Figure 10 and videos in supplemental materials show that as the repository shrinks, mouth motions become choppier as it becomes harder to find long matching phoneme subsequences in the repository and our tool has to introduce more transitions. However, the quality degrades gracefully.

*Varying the length of the edit.* We examine the effect of synthesizing edits of different lengths by synthesizing increasingly longer portions of the sentence "only the most accomplished artists obtain popularity"; starting by only synthesizing the first word and sequentially adding words until the full sentence is synthesized by our pipeline. Videos in supplemental materials show that, while the full sentence synthesis still has good mouth motions and looks plausible, it generally produce results that contain more artifacts than shorter edits, since with more phonemes to synthesize there are more opportunities for artifacts to emerge. Our user studies (Section 5.4) also found that results from short edits are rated more real than full-sentence syntheses.

## 5.3 Comparisons to Other Methods

*Comparisons with Synthesizing Obama [Suwajanakorn et al. 2017].* Suwajanakorn et al. [2017] have presented a technique for taking audio speech of Obama as input and synthesizing a video of Obama saying the speech. We compare our results to theirs using only 3 minutes of video of Obama. As shown in Figure 11, our approach gives more plausible and synchronized mouth motions compared to their method which can fail to close the mouth. While our approach

Fig. 13. Comparison to Neural Voice Puppetry (NVP) [Thies et al. 2020]. The actor fails to close his mouth with Neural Voice Puppetry at \m and \p, whereas our approach with no user refinement fully closes it. We refer to the supplemental materials for video comparison.
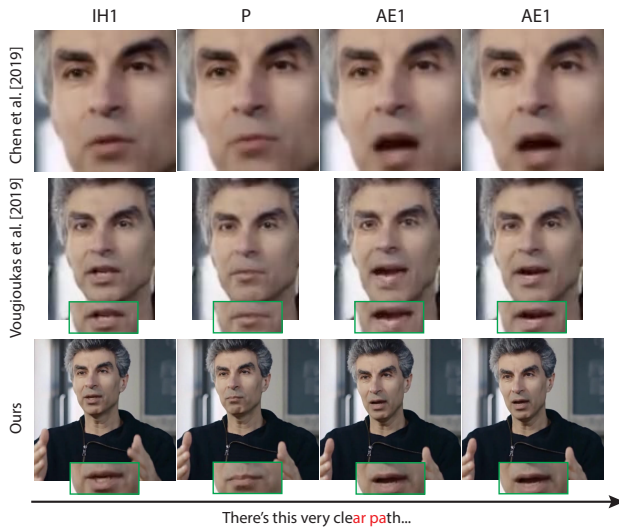


Fig. 14. Comparison to Chen et al. [2019] and Vougioukas et al. [2019]. All three methods produce good lip-sync, but our result has better quality due to sharper images and natural head motion.

produces good results with only 3 minutes of Obama video, their method needs 14 hours.

*Comparisons with Fried et al. [2019].* Our work builds on Fried et al.'s [2019] synthesis pipeline. However, they require over an hour of target video to produce high-quality results. We compare our method against theirs using 3 minutes of target video, as well as theirs using an hour of target video. Figure 12 and videos in supplemental materials show that our result has more temporally coherent mouth motions than Fried et al. [2019], both when they use the same 3 minutes of target video, and when they use 1 hour of target video. We further compare our results to those for Fried et al. [2019] in the user studies in Section 5.4.

*Comparisons with Neural Voice Puppetry [Thies et al. 2020].* Concurrent to our work, Neural Voice Puppetry (NVP) can synthesize talking head videos from audio signal input, or from text by using a

synthetic voice to obtain the audio signal. Given an audio speech track, we compare our result to theirs by applying our method on the phoneme timings extracted from the audio. Figure 13 and videos in supplemental materials show that while both approaches generate mouth motions that synchronize with the audio, our fully automatic result (without user refinement) generates closed-mouth frames at the desired phonemes (\m, \b, \p), whereas Neural Voice Puppetry leaves the mouth open for many of these phonemes. In addition, unlike NVP, our tool allows the user to iteratively refine the automatic results and adjust the performance. We further compare our results to NVP in the user studies in Section 5.4.

*Comparisons with Chen et al. [2019] and Vougioukas et al. [2019].* Both Chen et al. [2019] and Vougioukas et al. [2019] generate talking head videos from audio input and a single frame of the target actor. Given an audio speech track, we compare our result to theirs by applying our method on the phoneme timings extracted from the audio. Figure 14 and videos in supplemental materials show that while all three approaches produce good lip-sync with proper mouth closures, both Chen et al. [2019] and Vougioukas et al. [2019] produce videos of less resolution than our result. In addition, their results do not have the natural head motion in our result, and by always centering the video around the cropped head, their results can contain warping artifacts in the background, making them ill-suited for incorporation into a video-editing workflow.

## 5.4 User Studies and Automatic Metrics

We use both user studies and automatic metrics to quantitatively evaluate the quality of the video generated by our editing tool. In the user studies, we investigate both short and long edits, while ablating the target video length and the neural retargeting step. We compare to previous work [Fried et al. 2019; Thies et al. 2020] and to ground-truth video. We follow the study design used in previous talking-head synthesis research [Fried et al. 2019; Kim et al. 2018]. Specifically, participants see one video at a time in randomized order and are asked to rate the statement "This video clip looks real to me" on a 5-point Likert scale ranging from strongly disagree (1) to strongly agree (5). All videos used in our studies are available in the supplemental material.

*User study 1: Short Phrases (1 − 4 words).* Short phrases are the main type of result shown in Fried et al. [2019]. Such edits are useful for minor fix-ups on existing sentences. In this study we compare our automatic synthesis results ("Ours") to 3 versions of Fried et al. [2019]. (1) Their method with the same amount of target video as used by our tool, which is less than 5 minutes in all cases. (2) Their method with 1 hour of target video, which is their recommended amount, and more than 12 times the amount we use. (3) A version of their method with our fast phoneme search and stitching algorithm (Section 3.3), but with 1 hour of target video ("Modified Fried et al. [2019]" in Section 5.1), to evaluate the effect of ablating our neural retargeting step. We also compare to ground truth video recordings. We recruited 110 participants to view 25 videos each (5 conditions for each of 5 edits). We report Likert scale responses in Table 2 ("Short Phrase"). The differences between all pairs, except "Ours" vs. "Modified Fried", are statistically significant. All p-values have

been adjusted for multiple testing and are reported in supplemental materials.

Our tool outperforms Fried et al. [2019] both when using 5 minutes of data and 1 hour of data. We believe this is due to our results having more accurate mouth closures and better temporal coherency in mouth motions. Results are similar for our tool and Modified Fried, indicating that our neural retargeting step does not have much negative effect on result quality. Together these results also suggest that our fast phoneme search with stitching that forces closed mouths for \m, \b, and \p phonemes leads to higher-quality synthesis than the slow phoneme search and stitching approach used originally by Fried et al. [2019]. Although a gap still remains between our synthesized results and ground-truth videos, our results for short edits are rated as real almost two thirds of the time.

*User study 2: Full Sentences (6 − 9 words).* Full sentence synthesis is more challenging, since longer synthesis equates to a larger chance of inaccurate matches and synthesis artifacts. However, synthesizing full sentences as opposed to short phrases opens up more use cases (Section 4). Investigating full-sentence synthesis also emphasizes the quality differences between methods. The conditions in this user study are the same as for user study 1. We recruited 153 participants to view 25 videos each (5 conditions for each of 5 sentences). We report Likert scale responses in Table 2 ("Full Sentence"). The differences between all pairs, except "Fried < 5 min" vs. "Fried > 1 hr", are statistically significant. All p-values have been adjusted for multiple testing and are reported in supplemental materials.

Our tool produces the highest-quality results, followed by Modified Fried with over 1 hour of data, then by Fried et al. [2019]. Similar to user study 1, here our results have better mouth closures and smoother mouth motions than Fried et al. [2019]. It is worth noting that our results are even better than Modified Fried. We believe this is because our tool has a higher-quality source repository which becomes more salient when the edits are long. It shows the advantage of our approach to decouple source repository from target video, as data quality improvements to the repository can benefit many different target videos. The one-time cost of building a high-quality repository amortizes across all the edits that use it.

*User study 3: our tool vs Neural Voice Puppetry [Thies et al. 2020].* The third user study compares our results to those of Neural Voice Puppetry [Thies et al. 2020], where we show viewers videos generated by the two methods from the same audio speech track. We recruited 90 participants to view 8 videos each (4 from each of the two methods, Table 3). The audio tracks used in the videos are not the actor's real voice, and we believe this is the predominant reason for overall lower scores (for both methods). The difference between conditions is not statistically significant, and our results have similar mean scores to those of NVP. Nevertheless, as mentioned in Section 5.3, closely examining the videos generated by the two approaches, we find that our method often does a better job of closing the mouth on \m, \b, and \p phonemes. We also note that while our user studies evaluate our *automatic* results, unlike NVP, our tool also provides refinement and performance controls that can be used to improve results over the course of an interactive editing session.

| | Condition | Length | Likert response (%) | | | | | Mean | 'Real' |
|---|---|---|---|---|---|---|---|---|---|
| | | | 5 | 4 | 3 | 2 | 1 | | |
| Short Phrase | Fried et al. [2019] | < 5 min | 19.5 | 28.0 | 11.3 | 22.1 | 19.0 | 3.1 | 47.6% |
| | Fried et al. [2019] | > 1 hr | 24.1 | 31.5 | 13.7 | 20.0 | 10.7 | 3.4 | 55.6% |
| | Modified Fried | > 1 hr | 27.5 | 39.9 | 13.3 | 13.9 | 5.4 | **3.7** | **67.4%** |
| | Ours | < 5 min | 30.37 | 34.2 | 14.1 | 15.2 | 5.7 | **3.7** | 64.9% |
| | Ground truth | n/a | 40.1 | 37.8 | 11.5 | 9.8 | 0.9 | 4.1 | 77.8% |
| Full Sentence | Fried et al. [2019] | < 5 min | 14.7 | 22.9 | 11.1 | 25.0 | 26.3 | 2.7 | 37.6% |
| | Fried et al. [2019] | > 1 hr | 14.6 | 22.5 | 12.2 | 26.3 | 24.5 | 2.8 | 37.1% |
| | Modified Fried | > 1 hr | 16.5 | 31.5 | 14.1 | 20.7 | 17.2 | **3.1** | **48.0%** |
| | Ours | < 5 min | 17.9 | 38.2 | 16.0 | 20.1 | 7.7 | **3.4** | 56.2% |
| | Ground truth | n/a | 39.0 | 42.1 | 7.3 | 8.8 | 2.9 | 4.1 | 81.1% |

Table 2. Results from user studies on short phrases and full sentences. The "Length" column shows the length of the input target video for each method. We compare to previous work [Fried et al. 2019] and to ground-truth recordings. We report percentage of each answer on a 5-Point Likert scale, the mean score, and percent of videos that received a score of 4 or 5 ('real'). The difference between conditions is significant in both studies (Kruskal-Wallis test, $p < 10^{-20}$ each). A followup Tukey's range test shows that all pairwise comparisons are statistically significant ($p < 0.008$ each) except for "Ours" vs. "Modified Fried" for short edits and "Fried < 5 min" vs. "Fried > 1 hr" for full sentences. Note that Tukey's procedure adjusts the p-values for multiple comparisons. We report all adjusted p-values in the supplemental materials. Using our fast phoneme search and stitching algorithm improves results from Fried et al. [2019]. Our tool outperforms the method of Fried et al. [2019] when they use the same amount of target video data, and when they use 12x the amount of data. For full-sentence synthesis, our tool also outperforms Modified Fried even while they use 12x the amount of data.

| Condition | Length | Likert response (%) | | | | | Mean | 'Real' |
|---|---|---|---|---|---|---|---|---|
| | | 5 | 4 | 3 | 2 | 1 | | |
| Ours | < 5 min | 6.0 | 9.6 | 10.6 | 26.2 | 47.5 | 2.0 | 15.6% |
| NVP | < 5 min | 6.0 | 10.6 | 11.3 | 24.8 | 47.4 | 2.0 | 16.6% |

Table 3. Results from user study on our tool and Neural Voice Puppetry (NVP) [Thies et al. 2020]. We compare to NVP and report percentage of each answer on a 5-Point Likert scale, as well as mean score and percent of videos that received a score of 4 or 5 ('real'). Our tool and NVP received similar mean scores and the difference is not statistically significant (Kruskal-Wallis test, $p = 0.84$).

*Automatic Metrics.* For videos in user study 1 and 2 where we have ground-truth recordings, we further evaluate the results using automatic metrics against grouth truth videos. To measure reconstruction quality, we compute structural similarity index (SSIM [Zhou Wang et al. 2004]) and peak signal-to-noise ratio (PSNR). To measure lip motion quality, we compute the Landmarks Distances (LMD) [Chen et al. 2018]. We report the results in Table 4. Although our methods achieve favorable scores on many of these metrics, the score differences are quite small and visual differences can be subtle. We believe the user studies provide a better and more trustworthy measure of video quality.

| | Condition | Length | SSIM | PSNR | LMD |
|---|---|---|---|---|---|
| Short Phrase | Fried et al. [2019] | < 5 min | 0.89929 | 25.08146 | 4.57492 |
| | Fried et al. [2019] | > 1 hr | 0.89925 | 25.09057 | 4.47108 |
| | Modified Fried | > 1 hr | **0.89934** | **25.09670** | **4.38139** |
| | Ours | < 5 min | 0.89921 | 25.09349 | 4.57060 |
| Full Sentence | Fried et al. [2019] | < 5 min | 0.96339 | 32.74630 | 3.75896 |
| | Fried et al. [2019] | > 1 hr | 0.96552 | 32.94852 | 3.67367 |
| | Modified Fried | > 1 hr | 0.97578 | 35.00073 | **2.90959** |
| | Ours | < 5 min | **0.97630** | **35.12082** | 3.18670 |

Table 4. Results of automatic metrics. We compare our results to 3 versions of Fried et al. [2019] for both short phrase edits and full sentence syntheses. Best score is bolded. Our tool tops SSIM and PSNR for full sentence syntheses and ranks second after Modified Fried for LMD on full sentence and PSNR on short edits.

## 6 LIMITATIONS AND FUTURE WORK

We have demonstrated an iterative text-based tool for editing talking-head dialogue and performance that can be applied to many real-world editing scenarios in which only a few minutes of target actor video is available. However, our approach does have several limitations that could be addressed in future work.

*Further reduce feedback loop time.* Our tool currently requires about 30 seconds to synthesize a typical 5 word edit. While this feedback loop time allows users to try a variety of edits and refinements, seeing a synthesized result immediately (in real-time) would allow even more iteration and exploration of design space. As noted in Section 3.6, parallelization of our fast phoneme search and neural rendering steps as well as streaming playback of the synthesized video could reduce the feedback loop time significantly.

*Improve quality of synthesis results.* Although our method compares favorably in quality with previous talking-head synthesis techniques (Section 5.4), there is still a gap in realism between our results and ground-truth videos. As our method relies on a rich repository of source video to provide mouth motions for phoneme coarticulations, it may be possible to improve synthesis by developing higher-quality repositories. One approach may be to leverage existing work in text-driven 3D human mouth animation [Edwards et al. 2016] to render unlimited amounts of mouth motions to serve as the repository. Another direction is to build multiple repositories of many different source actors and then given a target video, develop techniques to pick the best source actor for the target.

*Performance controls over full face.* Our current approach focuses on synthesizing lip motions that match the target edit. While our tool offers controls for inserting mouth gestures and changing the speaking style, the effects of these controls are limited to the lower part of the face. Others have demonstrated techniques for controlling more of the head, including the ability to change head pose, gaze direction and whole face expressions [Kim et al. 2018] However, these techniques often introduce artifacts in the hair and with the clothes at the neckline. Adding such full face controls in an artifact-free manner remains an open research direction.

*Previsualizing dialogue using existing film scenes.* When writing dialogue, scriptwriters have to imagine the sound and appearance of the scene. Using our video editing tool with a catalog of video from existing film scenes might allow such scriptwriters to quickly visualize the dialogue in different settings and with different actors. Users might search for scenes based on their settings and actors using a tool like SceneSkim [Pavel et al. 2015] and our tool could be used to insert the new dialogue.

## 7 ETHICAL CONSIDERATIONS

Our editing tool is designed to enable an iterative workflow for removing filler words, adjusting phrasing, or correcting mistakes in a talking-head video. While such tools can facilitate content creation and storytelling, tools like ours, that let users manipulate what a target actor is saying, can also be misused. We follow the guidelines suggested by Fried et al. [2019] for ethically using such tools. (1) Video generated by our tool should be *transparent* about the fact that it has been manipulated. (2) Actors must give *consent* to any manipulation before a resulting video is shared widely.

We also recognize that these guidelines alone will not stop bad actors from using tools like ours to create false statements and slander others. Therefore, it is also critical for researchers to continue developing tools for detecting, fingerprinting, and verifying such video manipulation. Openly publishing the technical details of our tool can increase public awareness and help detection efforts. Ultimately these issues may also require regulations and laws that penalize misuse while allowing creative and consensual use cases.

## 8 CONCLUSION

Iterative editing is central to many content-creation tasks and is especially common in video editing. We have shown how to enable such iterative editing in the context of editing talking-head video using a text-based interface that allows changes to wording and facial performance while providing refinement controls. Whether an editor trying different ways to phrase the dialogue in a film, developing dialogue for a conversational agent, or correcting a mistake in a lecture, such iteration is often essential for finding the most appropriate result. We believe such tools that facilitate video editing can democratize content-creation and enable many more people to tell their stories.

## REFERENCES

2020a. Google Cloud Speech to Text API. https://cloud.google.com/speech-to-text
2020b. Google Cloud Text to Speech API. https://cloud.google.com/text-to-speech
2020. Lyrebird AI. https://www.descript.com/lyrebird-ai
2020. Rev. https://rev.com
Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dandelion Mané, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. 2015. TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems. https://www.tensorflow.org/ Software available from tensorflow.org.
Shruti Agarwal, Hany Farid, Ohad Fried, and Maneesh Agrawala. 2020. Detecting Deep-Fake Videos from Phoneme-Viseme Mismatches. In *Workshop on Media Forensics at CVPR*. Seattle, WA, USA.

Hadar Averbuch-Elor, Daniel Cohen-Or, Johannes Kopf, and Michael F. Cohen. 2017. Bringing Portraits to Life. 36, 6 (Nov. 2017), 196:1–13. https://doi.org/10.1145/3130800.3130818

Floraine Berthouzoz, Wilmot Li, and Maneesh Agrawala. 2012. Tools for Placing Cuts and Transitions in Interview Video. ACM Trans. Graph. 31, 4, Article 67 (July 2012), 8 pages. https://doi.org/10.1145/2185520.2185563

Volker Blanz and Thomas Vetter. 1999. A Morphable Model for the Synthesis of 3D Faces. In Proceedings of the 26th Annual Conference on Computer Graphics and Interactive Techniques (SIGGRAPH '99). ACM Press/Addison-Wesley Publishing Co., USA, 187–194. https://doi.org/10.1145/311535.311556

Christoph Bregler, Michele Covell, and Malcolm Slaney. 1997. Video Rewrite: Driving Visual Speech with Audio. In Proceedings of the 24th Annual Conference on Computer Graphics and Interactive Techniques (SIGGRAPH '97). ACM Press/Addison-Wesley Publishing Co., USA, 353–360. https://doi.org/10.1145/258734.258880

Yao-Jen Chang and Tony Ezzat. 2005. Transferable Videorealistic Speech Animation. 143–151. https://doi.org/10.1145/1073368.1073388

Lele Chen, Zhiheng Li, Ross K Maddox, Zhiyao Duan, and Chenliang Xu. 2018. Lip movements generation at a glance. In Proceedings of the European Conference on Computer Vision (ECCV). 520–535.

Lele Chen, Ross K Maddox, Zhiyao Duan, and Chenliang Xu. 2019. Hierarchical cross-modal talking face generation with dynamic pixel-wise loss. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 7832–7841.

Joon Son Chung, Amir Jamaludin, and Andrew Zisserman. 2017. You said that?. In British Machine Vision Conference.

Pif Edwards, Chris Landreth, Eugene Fiume, and Karan Singh. 2016. JALI: an animator-centric viseme model for expressive lip synchronization. ACM Transactions on Graphics (TOG) 35, 4 (2016), 1–11.

Tony Ezzat, Gadi Geiger, and Tomaso Poggio. 2002. Trainable Videorealistic Speech Animation. 21, 3 (July 2002), 388–398. https://doi.org/10.1145/566654.566594

Ohad Fried and Maneesh Agrawala. 2019. Puppet Dubbing. In Eurographics Symposium on Rendering - DL-only and Industry Track, Tamy Boubekeur and Pradeep Sen (Eds.). The Eurographics Association. https://doi.org/10.2312/sr.20191220

Ohad Fried, Ayush Tewari, Michael Zollhöfer, Adam Finkelstein, Eli Shechtman, Dan B Goldman, Kyle Genova, Zeyu Jin, Christian Theobalt, and Maneesh Agrawala. 2019. Text-Based Editing of Talking-Head Video. ACM Trans. Graph. 38, 4, Article 68 (July 2019), 14 pages. https://doi.org/10.1145/3306346.3323028

John S Garofolo, Lori F Lamel, William M Fisher, Jonathan G Fiscus, and David S Pallett. 1993. DARPA TIMIT acoustic-phonetic continous speech corpus CD-ROM. NIST speech disc 1-1.1. NASA STI/Recon technical report n 93 (1993).

Pablo Garrido, Levi Valgaerts, Ole Rehmsen, Thorsten Thormaehlen, Patrick Pérez, and Christian Theobalt. 2014. Automatic Face Reenactment. 4217–4224. https://doi.org/10.1109/CVPR.2014.537

Pablo Garrido, Levi Valgaerts, Hamid Sarmadi, Ingmar Steiner, Kiran Varanasi, Patrick Pérez, and Christian Theobalt. 2015. VDub: Modifying Face Video of Actors for Plausible Visual Alignment to a Dubbed Audio Track. 34, 2 (May 2015), 193–204. https://doi.org/10.1111/cgf.12552

Pablo Garrido, Michael Zollhöfer, Dan Casas, Levi Valgaerts, Kiran Varanasi, Patrick Pérez, and Christian Theobalt. 2016. Reconstruction of Personalized 3D Face Rigs from Monocular Video. ACM Trans. Graph. 35, 3, Article 28 (May 2016), 15 pages. https://doi.org/10.1145/2890493

Jiahao Geng, Tianjia Shao, Youyi Zheng, Yanlin Weng, and Kun Zhou. 2018. Warp-guided GANs for Single-photo Facial Animation. In SIGGRAPH Asia 2018 Technical Papers (Tokyo, Japan) (SIGGRAPH Asia '18). ACM, New York, NY, USA, Article 231, 231:1–231:12 pages. http://doi.acm.org/10.1145/3272127.3275043

Ye Jia, Yu Zhang, Ron Weiss, Quan Wang, Jonathan Shen, Fei Ren, Patrick Nguyen, Ruoming Pang, Ignacio Lopez Moreno, Yonghui Wu, et al. 2018. Transfer learning from speaker verification to multispeaker text-to-speech synthesis. In Advances in neural information processing systems. 4480–4490.

Ira Kemelmacher-Shlizerman, Aditya Sankar, Eli Shechtman, and Steven M. Seitz. 2010. Being John Malkovich. 341–353. https://doi.org/10.1007/978-3-642-15549-9_25

Hyeongwoo Kim, Mohamed Elgharib, Michael Zollhöfer, Hans-Peter Seidel, Thabo Beeler, Christian Richardt, and Christian Theobalt. 2019. Neural Style-Preserving Visual Dubbing. ACM Trans. Graph. 38, 6, Article 178 (Nov. 2019), 13 pages. https://doi.org/10.1145/3355089.3356500

H. Kim, P. Garrido, A. Tewari, W. Xu, J. Thies, N. Nießner, P. Pérez, C. Richardt, M. Zollhöfer, and C. Theobalt. 2018. Deep Video Portraits. ACM Transactions on Graphics 2018 (TOG) (2018).

Diederik P. Kingma and Jimmy Ba. 2015. Adam: A Method for Stochastic Optimization. In 3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings, Yoshua Bengio and Yann LeCun (Eds.). http://arxiv.org/abs/1412.6980

Kundan Kumar, Rithesh Kumar, Thibault de Boissiere, Lucas Gestin, Wei Zhen Teoh, Jose Sotelo, Alexandre de Brebisson, Yoshua Bengio, and Aaron Courville. 2019. MelGAN: Generative Adversarial Networks for Conditional Waveform Synthesis. arXiv:1910.06711 [eess.AS]

Rithesh Kumar, Jose Sotelo, Kundan Kumar, Alexandre de Brebisson, and Yoshua Bengio. 2017. ObamaNet: Photo-realistic lip-sync from text. arXiv:1801.01442 [cs.CV]

VI Levenshtein. 1966. Binary Codes Capable of Correcting Deletions, Insertions and Reversals. In Soviet Physics Doklady, Vol. 10. 707.

Kang Liu and Joern Ostermann. 2011. Realistic facial expression synthesis for an image-based talking head. https://doi.org/10.1109/ICME.2011.6011835

Wesley Mattheyses, Lukas Latacz, and Werner Verhelst. 2010. Optimized photorealistic audiovisual speech synthesis using active appearance modeling. In Auditory-Visual Speech Processing. 8–1.

Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. 2013. On the difficulty of training recurrent neural networks. In International conference on machine learning. 1310–1318.

Amy Pavel, Dan B Goldman, Björn Hartmann, and Maneesh Agrawala. 2015. Sceneskim: Searching and browsing movies using synchronized captions, scripts and plot summaries. In Proceedings of the 28th Annual ACM Symposium on User Interface Software & Technology. ACM, 181–190.

A. Pumarola, A. Agudo, A.M. Martinez, A. Sanfeliu, and F. Moreno-Noguer. 2019. GANimation: One-Shot Anatomically Consistent Facial Animation. (2019).

Steve Rubin, Floraine Berthouzoz, Gautham J Mysore, Wilmot Li, and Maneesh Agrawala. 2013. Content-based tools for editing audio stories. In Proceedings of the 26th annual ACM symposium on User interface software and technology. ACM, 113–122.

Yang Song, Jingwen Zhu, Dawei Li, Andy Wang, and Hairong Qi. 2019. Talking Face Generation by Conditional Recurrent Adversarial Network. In Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI-19. International Joint Conferences on Artificial Intelligence Organization, 919–925. https://doi.org/10.24963/ijcai.2019/129

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

Ayush Tewari, Ohad Fried, Justus Thies, Vincent Sitzmann, Stephen Lombardi, Kalyan Sunkavalli, Ricardo Martin-Brualla, Tomas Simon, Jason Saragih, Matthias Nießner, Rohit Pandey, Sean Fanello, Gordon Wetzstein, Jun-Yan Zhu, Christian Theobalt, Maneesh Agrawala, Eli Shechtman, Dan B. Goldman, and Michael Zollhöfer. 2020. State of the Art on Neural Rendering. Computer Graphics Forum (2020). https://doi.org/10.1111/cgf.14022

Justus Thies, Mohamed Elgharib, Ayush Tewari, Christian Theobalt, and Matthias Nießner. 2020. Neural Voice Puppetry: Audio-driven Facial Reenactment. In European Conference on Computer Vision. Springer, 716–731.

Justus Thies, Michael Zollhöfer, Marc Stamminger, Christian Theobalt, and Matthias Nießner. 2016. Face2Face: Real-Time Face Capture and Reenactment of RGB Videos. 2387–2395. https://doi.org/10.1109/CVPR.2016.262

Aäron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alexander Graves, Nal Kalchbrenner, Andrew Senior, and Koray Kavukcuoglu. 2016. WaveNet: A Generative Model for Raw Audio. In Arxiv. https://arxiv.org/abs/1609.03499

Daniel Vlasic, Matthew Brand, Hanspeter Pfister, and Jovan Popović. 2005. Face Transfer with Multilinear Models. 24, 3 (July 2005), 426–433. https://doi.org/10.1145/1073204.1073209

Konstantinos Vougioukas, Stavros Petridis, and Maja Pantic. 2018. End-to-End Speech-Driven Facial Animation with Temporal GANs. arXiv:1805.09313 [eess.AS]

Konstantinos Vougioukas, Stavros Petridis, and Maja Pantic. 2019. Realistic Speech-Driven Facial Animation with GANs. International Journal of Computer Vision (13 Oct 2019). https://doi.org/10.1007/s11263-019-01251-8

Robert A Wagner and Michael J Fischer. 1974. The string-to-string correction problem. Journal of the ACM (JACM) 21, 1 (1974), 168–173.

Lijuan Wang, Wei Han, Frank Soong, and Qiang Huo. 2011. Text-driven 3D Photo-Realistic Talking Head. In INTERSPEECH 2011 (interspeech 2011 ed.). International Speech Communication Association. https://www.microsoft.com/en-us/research/publication/text-driven-3d-photo-realistic-talking-head/

O. Wiles, A.S. Koepke, and A. Zisserman. 2018. X2Face: A network for controlling face generation by using images, audio, and pose codes. In European Conference on Computer Vision.

Jiahong Yuan and Mark Liberman. 2008. Speaker identification on the SCOTUS corpus. In In Proceedings of Acoustics 2008. Citeseer.

Egor Zakharov, Aliaksandra Shysheya, Egor Burkov, and Victor Lempitsky. 2019. Few-Shot Adversarial Learning of Realistic Neural Talking Head Models. arXiv:1905.08233 [cs.CV]

Hang Zhou, Yu Liu, Ziwei Liu, Ping Luo, and Xiaogang Wang. 2019. Talking Face Generation by Adversarially Disentangled Audio-Visual Representation. In AAAI Conference on Artificial Intelligence (AAAI).

Zhou Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli. 2004. Image quality assessment: from error visibility to structural similarity. IEEE Transactions on Image Processing 13, 4 (2004), 600–612.