

End-to-End Video-To-Speech Synthesis using Generative Adversarial Networks

Rodrigo Mira, Konstantinos Vougioukas, Pingchuan Ma, Stavros Petridis *Member, IEEE*
 Björn W. Schuller *Fellow, IEEE*, Maja Pantic *Fellow, IEEE*

Video-to-speech is the process of reconstructing the audio speech from a video of a spoken utterance. Previous approaches to this task have relied on a two-step process where an intermediate representation is inferred from the video, and is then decoded into waveform audio using a vocoder or a waveform reconstruction algorithm. In this work, we propose a new end-to-end video-to-speech model based on Generative Adversarial Networks (GANs) which translates spoken video to waveform end-to-end without using any intermediate representation or separate waveform synthesis algorithm. Our model consists of an encoder-decoder architecture that receives raw video as input and generates speech, which is then fed to a waveform critic and a power critic. The use of an adversarial loss based on these two critics enables the direct synthesis of raw audio waveform and ensures its realism. In addition, the use of our three comparative losses helps establish direct correspondence between the generated audio and the input video. We show that this model is able to reconstruct speech with remarkable realism for constrained datasets such as GRID, and that it is the first end-to-end model to produce intelligible speech for LRW (Lip Reading in the Wild), featuring hundreds of speakers recorded entirely ‘in the wild’. We evaluate the generated samples in two different scenarios – seen and unseen speakers – using four objective metrics which measure the quality and intelligibility of artificial speech. We demonstrate that the proposed approach outperforms all previous works in most metrics on GRID and LRW.

I. INTRODUCTION

AUTOMATIC speech recognition (ASR) is a well established field with diverse applications including captioning voiced speech and recognizing voice commands. Deep learning has revolutionised this task in the past years, to the point where state of the art models are able to achieve very low word error rates (WER) [29]. Although these models are reliable for clean audio, they struggle under noisy conditions [32, 49], and they are not effective when gaps are found in the audio stream [58]. The recurrence of these edge cases has driven researchers towards Visual Speech Recognition (VSR), also known as lipreading, which performs speech recognition based on video only.

Although the translation from video-to-text can now be achieved with remarkable consistency, there are various applications that would benefit from a video-to-audio model, such as videoconferencing in noisy conditions; speech inpainting [58], i. e., filling in audio gaps from video in an audiovisual stream; or generating an artificial voice for people suffering from aphonia (i. e., people who are unable to produce voiced sound). One approach for this task would be to simply combine a lipreading model (which outputs text) with a text-to-speech (TTS) model (which outputs audio). This approach is especially attractive since state-of-the-art TTS models can now produce realistic speech with considerable efficacy [36, 46].

This work has been submitted to the IEEE for possible publication. Copyright may be transferred without notice, after which this version may no longer be accessible. Manuscript received Month Day, Year; revised Month Day, Year. Corresponding author: R. Mira (email: rs2517@ic.ac.uk). Rodrigo Mira would like to thank Samsung for their continued support of his work on this project.

Rodrigo Mira, Konstantinos Vougioukas, Pingchuan Ma, Stavros Petridis, Björn W. Schuller, Maja Pantic are with the IBUG Group, Department of Computing, Imperial College London, UK

Stavros Petridis is with the Samsung AI Centre Cambridge, UK.

Björn W. Schuller is with the Chair of Embedded Intelligence for Health Care and Wellbeing, University of Augsburg, Germany.

Maja Pantic is with Facebook London, UK.

Combining video-to-text and text-to-speech models to perform video-to-speech has, however, some disadvantages. Firstly, these models require large transcribed datasets, since they are trained with text supervision. This is a sizeable constraint given that generating transcripts is a time consuming and expensive process. Secondly, generation can only happen as each word is recognized, which imposes a delay on the throughput of the model, jeopardizing the viability of real-time synthesis. Lastly, using text as an intermediate representation removes any intonation and emotion from the spoken statement, which are fundamental for natural sounding speech.

Given these constraints, some authors have developed end-to-end video-to-speech models which circumvent these issues. The first of these models [9] used visual features based on discrete cosine transform (DCT) and active appearance models (AAM) to predict linear predictive coding (LPC) coefficients and mel-filterbank amplitudes. Following works have mostly focused on predicting spectrograms [1, 13, 43], which is also a common practice in text-to-speech works [46]. These models achieve intelligible results, but are only applied to seen speakers, i. e., there is exact correspondence between the speakers in the training, validation and test sets, or choose to focus on single speaker speech reconstruction [43]. Recently, [33] has proposed an alternative approach based on predicting WORLD vocoder parameters [34] which generates clear speech for unseen speakers as well. However, the reconstructed speech is still not realistic.

It is clear that previous works have avoided synthesising raw audio, likely due to the lack of a suitable loss function, and have focused on generating intermediate representations which are then used for reconstructing speech. To the best of our knowledge, the only work which directly synthesises the raw audio waveform from video is [53]. This work introduces the use of GANs [3, 15], and thanks to the adversarial loss, it is able to directly reconstruct the audio waveform. This approach also produces realistic utterances for seen speakers, and is the

first to produce intelligible speech for unseen speakers.

Our work builds upon the model presented in [53] by proposing architectural changes to the model, and to the training procedure. Firstly, we replace the original encoder composed of five stacked convolutional layers with a ResNet-18 [20] composed of a front end 3D convolutional layer (followed by a max pooling layer), four blocks containing four convolutional layers each and an average pooling layer. Additionally, we replace the GRU (Gated Recurrent Unit) layer following the encoder with two bidirectional GRU layers, increasing the capacity of our temporal model. The adversarial methodology was a major factor towards generating intelligible waveform audio in [53]. Hence, our approach is also based on the Wasserstein GAN [3], but we propose a new critic adapted from [25]. We also propose an additional critic which discriminates real from synthesized spectrograms.

Furthermore, we revise the loss configuration presented in [53]. Firstly, we decide to forego the use of the total variation loss and the L1 loss, as their benefit was minimal. Secondly, we use the recently proposed PASE (Problem Agnostic Speech Encoder) [38] as a perceptual feature extractor. Finally, we propose two additional losses, the power loss and the MFCC loss. The power loss is an L1 loss between the (log-scaled) spectrograms of the real and generated waveforms. The MFCC loss is an L1 Loss between the MFCCs (Mel Frequency Cepstral Coefficients) of the real and generated waveforms.

Our contributions for this work are described as follows: **1)** We propose a new approach for reconstructing waveform speech directly from video based on GANs without using any intermediate representations. We use two separate critics to discriminate real from synthesized waveforms and spectrograms respectively, and apply three comparative losses to improve the quality of outputs. **2)** We include a detailed ablation study where we measure the effect of each component on the final model. We also investigate how the type of visual input, size of training set and range of vocabulary affect the performance. **3)** We show results on two different datasets (GRID [8] and TCD-TIMIT [19]) for seen speakers. We find that our model substantially outperforms the state-of-the-art for GRID and adapts well to a larger pool of speakers. **4)** We also include results for unseen speakers on two datasets (GRID and LRW [6]). We show that our model achieves intelligible results, even when applied to utterances recorded ‘in the wild’, and outperforms the state-of-the-art for the corpora we present. **5)** Finally, we study our model’s ability to generalize for videos of silent speakers, and discuss our findings.

II. RELATED WORK

Video-driven speech reconstruction is effectively the combination of two tasks: lipreading and speech synthesis. As such, we begin by briefly describing the main works in each field, and then go on to describe existing approaches for video-to-speech.

A. Lipreading

Traditional lipreading approaches relied on HMMs (Hidden Markov Models) [17] or SVMs (Support Vector Machines) [57] to transcribe videos from manually extracted features such as

DCT [17] or mouth geometry [24]. Recently, end-to-end models have attracted attention due to their superior performance over traditional approaches. One of the first end-to-end architectures for lipreading was [4]. This model featured a convolutional encoder as the visual feature extractor and a two-layer BGRU-RNN (Bidirectional GRU recurrent neural network) followed by a linear layer as the classifier, and it achieved state of the art performance for the GRID corpus. This work was followed by [7], whose model relied entirely on CNNs (Convolutional Neural Networks) and RNNs, and was successfully applied to spoken utterances recorded in the wild.

Various works have followed which apply end-to-end deep learning models to achieve competitive lipreading performance. [39, 42] propose an encoder composed of fully connected layers and performs classification using LSTMs (Long-short Term Memory RNNs). Other works choose to use convolutional encoders [47], often featuring residual connections [48], and then apply RNNs to perform classification. Furthermore, these end-to-end architectures have been extended for multi-view lipreading [41] and audiovisual [40] speech recognition.

B. Speech Synthesis

One of the most popular speech synthesis models in recent years has been WaveNet [35], which proposed dilated causal convolutions to compose waveform audio sample by sample, taking advantage of the large receptive field achieved by stacking these layers. This model achieved far more realistic results than any artificial synthesizer proposed before then. Another work [55] introduced a vastly different sequence-to-sequence model that predicted linear-scale spectrograms from text, which were then converted into waveform using the Griffin-Lim Algorithm (GLA) [60]. This process produced very clear and intelligible audio. In the following years, [46] combined these two methodologies to push the state-of-the-art once more, and [36] accelerated and improved the original WaveNet.

The first model to apply GANs for end-to-end speech synthesis was [11], which used simple convolutional networks with large kernels as the generator and discriminator and applied the improved Wasserstein loss [16]. In a later work [56], the original WaveNet vocoder [35] has been combined with the adversarial methodology introduced in [11]. This results in a network which has far less parameters than the original WaveNet, but remains on par with the latest WaveNet-based models. Recently, the first end-to-end adversarial Text-To-Speech model [12] was also proposed, whose performance is comparable to the state-of-the-art.

C. Reconstructing audio from visual speech

To the best of the authors’ knowledge, the first work to attempt the task of video-to-speech synthesis directly was [9]. The proposed model aims to predict the spectral envelope (LPC or mel-filterbanks) from manually extracted visual features (DCT or AAM) using Gaussian Mixture Models (GMMs) or deep neural networks. These acoustic parameters are then fed into an HMM-based vocoder, together with an estimate of the voicing parameters. Through multiple user studies, the speech reconstructed by this model is shown to have fairly low

intelligibility (WER $\approx 50\%$), but shows that this task is indeed achievable. This work was extended in [10], which introduced additional temporal information in the visual features and in the model itself. These improvements yielded an impressive 15% WER for GRID (single speaker), based on user studies.

The next development in this field comes with [14], which uses a deep CNN architecture to predict acoustic features – LPC analysis followed by LSP (Line Spectral Pairs) decomposition, frame by frame – from gray-scale video frames. These are combined with white noise (excitation signal) and fed into a source-filter speech synthesizer which produces unvoiced speech. This model produces intelligible results (WER $< 20\%$) when trained and tested on a single speaker from GRID, and constitutes a step forward given that it no longer relies on handcrafted visual features as input. An improved version of this model was presented in [13], which predicts spectrograms that are then translated into waveform using the Griffin-Lim algorithm. This extension also proposes a new encoder composed of two ResNet-18s followed by a post-processing network which increases temporal resolution. This work is the first to experiment with multiple speakers and achieves much more realistic speech than any previous work for this task.

Lip2Audspec [1] proposes a similar CNN+RNN encoder to predict spectrograms directly from the gray-scale frames of the video. As in [13], the spectrograms are converted to waveform using a phase estimation method. The resulting spectrograms are very close to the original samples, but the reconstructed waveforms sound noticeably robotic. Another recent work [33] uses CNNs+RNNs to predict vocoder parameters (aperiodicity and spectral envelope), rather than spectrograms. Additionally, the model is trained to predict the transcription of the speech, in other words performing speech reconstruction and recognition simultaneously in a multi-task fashion. This approach achieves results which are very impressive when measured with objective speech quality metrics (PESQ, STOI), but yields samples which still sound noticeably robotic.

Finally, a recent work [43] proposes an approach based on the Tacotron 2 architecture [46], predicting mel-frequency spectrograms from video rather than text. To perform this task, it applies a stack of residual 3D convolutional layers as a spatio-temporal encoder for the video, and combines it with an attention-based decoder adapted from [46], which generates the spectrograms. Unlike Tacotron, these spectrograms are decoded into waveform audio using the Griffin-Lim algorithm [60] rather than WaveNet, as the authors claim the generated spectrograms are not as accurate as modern TTS works, and therefore do not perform well with neural vocoders. This work is able to generate remarkably intelligible audio from visual speech and achieves state-of-the-art performance in all presented metrics. However, it focuses on speaker specific speech reconstruction, i. e., it is trained and tested on the same speaker.

An aspect which is worth highlighting is that none of these models attempt to generate the waveform end-to-end from video, instead predicting spectrograms or other features which can be translated into waveform. This is likely due to the notoriously arduous task of generating realistic waveforms, which can be attributed to the lack of suitable loss functions. The only model to perform video-to-waveform speech recon-

struction without the use of intermediate representations is [53]. This work proposes a generative adversarial network based on a convolutional encoder-decoder model (combined with a GRU) which encodes video into visual features and decodes them directly into waveform audio. The generator is trained with an adversarial loss based on a convolutional waveform critic, as well as three other comparative losses. This procedure achieves competitive results for speech reconstruction on both seen and unseen speaker datasets (GRID).

D. Reconstructing audio from multi-view visual speech

The majority of works in video-driven speech reconstruction use frontal views of the face. In this section, we briefly describe a set of works which use multiple-views in order to improve the quality of reconstructed speech.

The first work to use multi-view video for this task was [26]. This model is very similar to [14] in the sense that it applies a CNN to extract visual features directly from video, which then predict vocoder parameters (LPC followed by LSP). This work, however, uses video taken from two different angles for every speaker (Oulu VS2 dataset [2]). The results presented in this paper show that the use of multiple views can substantially improve speech reconstruction performance.

This model has been improved in [28] by replacing the LSTM with a BGRU and using more than two views as input. It is shown that the use of three views can yield improvements of 20% in the quality of reconstructed outputs (measured with PESQ). This has been extended in [27] by including a view classifier to attribute view labels to the input videos and by also generating text transcriptions. The latest work in this field [52] follows the trend seen in single-view speech reconstruction research [13, 14] and speech synthesis in general [51, 55] by switching from LPC coefficients to spectrograms as the predicted audio representation.

E. Audio reconstruction from video in other applications

Finally, a set of past works has approached the application of Video-to-Audio models to domains outside speech [5, 37, 59]. Namely, these papers have focused on a diverse range of datasets which feature a set of generic sounds such as fireworks and drums [59]; different instruments being played [5]; or even objects composed of different materials being hit with a drumstick [37]. The methodology applied to reconstruct audio from video is similar to what is seen for video-to-speech systems. CNNs are applied to encode the video frames, followed by RNNs or fully connected layers to produce acoustic features which are decoded into audio using vocoders. While some of these works struggle to reproduce the corresponding audio, [59] is able to produce remarkably realistic audio (as proven by its user studies) by combining the extraction of optical flow with a neural network-based vocoder.

III. VIDEO-DRIVEN SPEECH RECONSTRUCTION

Our model is composed of a video encoder based on a ResNet-18 combined with a Bidirectional GRU, as well as a convolutional decoder which transforms the visual features into

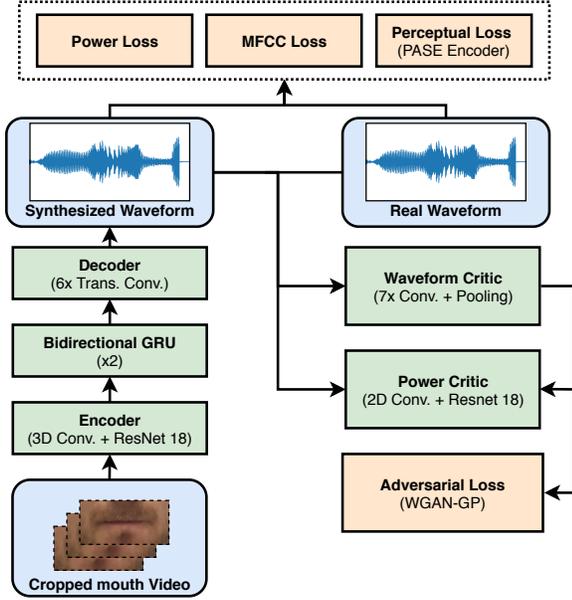


Fig. 1: Architecture of the generator (encoder, bidirectional GRU, decoder) and critics (waveform critic, power critic) used in this work, as well as the losses that are used for training.

waveform audio. This generator is trained using two separate critics, to ensure the realism of the outputs, as well as three L1 losses to minimize the difference between real and synthesized audio for each video.

A. Generator

Given that we aim to synthesize speech directly from video, our generator accomplishes two sequential tasks: encode temporal visual features and decode them into an audio waveform. Firstly, we encode the frames of the video using a Resnet-18 preceded by a spatio-temporal 3D convolutional layer (combined with a max pooling layer). This initial layer has a receptive field of 5 frames centered on the frame it will encode meaning that the encoding for each frame will depend on the previous two frames and on the following two frames. We experimented with different numbers of frames as input to this layer (3 and 7), but found that this did not considerably affect results. The ResNet-18 is composed of 4 blocks of 4 convolutional layers, each followed by batch normalization and ReLU (Rectified Linear Unit) activation, and an adaptive average pooling layer. The features extracted from the ResNet encoder are then fed into a 2-layer bidirectional GRU which temporally correlates the features produced from each set of frames. This architecture is described in detail in Figure 2.

After this, the decoder upsamples the features from each video frame into a waveform segment of N audio samples. The length of each segment is given by:

$$N = \frac{\text{audio sampling rate}}{\text{video frame rate}}. \quad (1)$$

Since we use a sampling rate of 16 kHz and a frame rate of 25 frames per second, N is equal to 640 (corresponding to 40 ms of audio). The decoder is composed of six stacked transposed

convolutional layers, each followed by batch normalization and ReLU activation except for the last layer which uses a hyperbolic tangent activation function. In an attempt to alleviate the issue of abrupt frame transitions, we use an overlap of 50% between the generated waveform frames, as proposed in [14]. The overlapped segments are linearly averaged sample by sample in order to maintain the original waveform scale. The detailed architecture of the decoder is shown in Figure 3.

B. Critics

As demonstrated in recent works [11, 25, 56], the use of a waveform critic can dramatically increase the realism and clarity of synthesized speech. To discriminate the real from the synthesized waveforms, we adapt the critic from [25]. After experimenting extensively with and without weight normalization for this module, as well as for the generator, we find that weight normalization increases the stability of adversarial training but overall leads to worse results. Therefore, we remove weight normalization from this critic but otherwise keep the original architecture: 7 convolutional layers, each followed by Leaky ReLU activation, as shown in Figure 4.

We did not attempt batch normalization, which worked well for the generator, since this interferes with the gradient penalty for our adversarial loss [16]. We compared this architecture to other convolutional critics similar to the one proposed in [11] as well as a one-dimensional ResNet 18, and found that this critic produced the best results. Remarkably, this critic has a far smaller receptive field than any of the critics we experimented with. This may indicate that waveform critics work best when focusing on the small scale.

Inspired by the SpecGAN model [11], we propose to combine the waveform critic, which judges the audio in the temporal domain, with a power critic, which judges the audio in the spectral domain. This module discriminates the spectrograms computed from real and generated audio. We first compute the spectrogram from both the real and generated samples using the short-time Fourier transform (STFT) with a window size of 25 ms, a hop size of 10 ms and frequency bins of size 512. We then compute the natural logarithm of the spectrogram magnitudes, normalize these values to mean 0 and variance 1, clip values outside $[-3, 3]$ and normalize them to $[-1, 1]$, similarly to [11]. In this case, we use a ResNet18 identical to the one presented in our generator, except with a two-dimensional front end convolutional layer in the beginning, since our input is a single image. As with the waveform critic, we cannot use batch normalization in this module due to the gradient penalty, and found that weight normalization did not improve results. The architecture for the power critic is shown in Figure 5.

C. Losses

To train our network, we apply the Wasserstein GAN loss [3], which aims to minimize the Wasserstein Distance between the distributions of real and synthesized data. We also add the gradient penalty [16] in order to satisfy the Lipschitz constraint

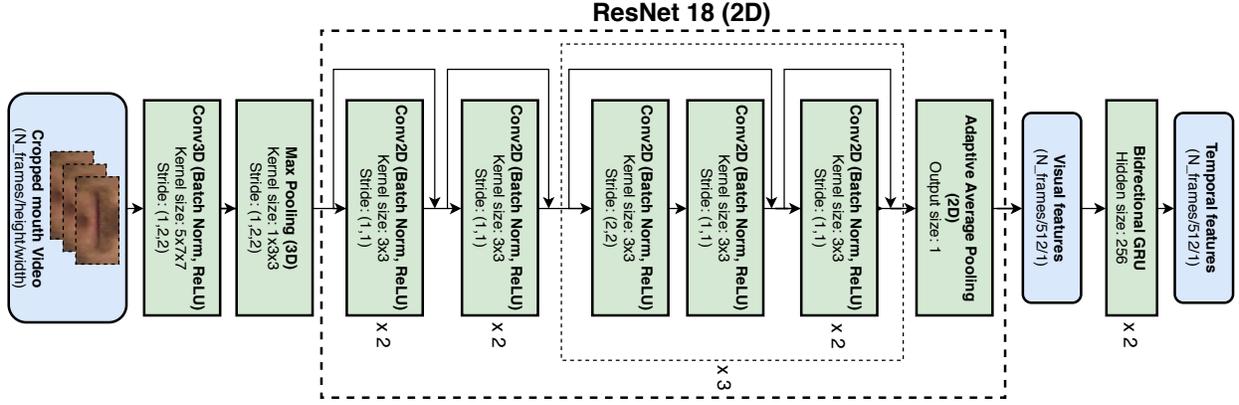


Fig. 2: Description of the layers in the encoder (generator).

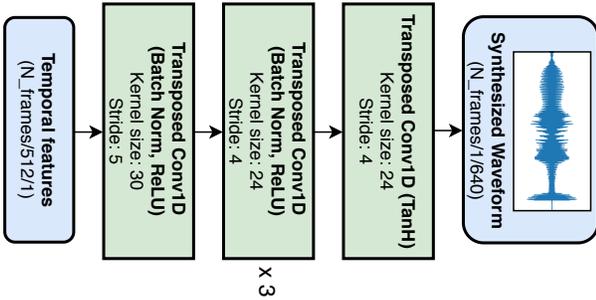


Fig. 3: Description of the layers in the decoder (generator).

in the Wasserstein GAN objective. The losses for the generator and respective critic(s) are defined as:

$$L_G = -\mathbb{E}_{\tilde{x} \sim \mathbb{P}_G} [D(\tilde{x})] + \lambda \mathbb{E}_{\hat{x} \sim \mathbb{P}_{\hat{x}}} [(\|\nabla_{\hat{x}} D(\hat{x})\| - 1)^2] \quad (2)$$

$$L_D = \mathbb{E}_{\tilde{x} \sim \mathbb{P}_G} [D(\tilde{x})] - \mathbb{E}_{x \sim \mathbb{P}_R} [D(x)], \quad (3)$$

where G is the generator, D is the critic, $x \sim \mathbb{P}_R$ are samples from the real distribution, $\tilde{x} \sim \mathbb{P}_G$ are samples from the estimated distribution (produced by the generator) and $\hat{x} \sim \mathbb{P}_{\hat{x}}$ are sampled uniformly between two points from \mathbb{P}_G and \mathbb{P}_R respectively. In this work, we apply two critics: the waveform critic and the power critic. Each critic is trained with their own losses $L_{D_{wave}}$ and $L_{D_{power}}$, whereas the generator combines the losses from the two critics such that:

$$L_{G_{adv}} = L_{G_{wave}} + L_{G_{power}}, \quad (4)$$

where $L_{G_{wave}}$ and $L_{G_{power}}$ are calculated as mentioned in Eq. 2. The coefficient for the gradient penalty λ is kept at the value of 10 for both critics, as proposed in [16].

In addition to this adversarial loss, we also apply three other losses to train the generator. The first is a perceptual loss:

$$L_{PASE} = \|\delta(x) - \delta(\tilde{x})\|, \quad (5)$$

where x is the real waveform, \tilde{x} is the synthesized waveform from the same video and δ is our perceptual feature extractor. In this work, we use the pre-trained PASE model [38] to extract perceptual features $\delta(x)$. PASE has been trained in a self-supervised manner to produce meaningful speech representations. We have also tried using PASE+ [44], which is

an improved version of PASE, however, no improvement in the speech reconstruction quality was observed. Furthermore, we experimented with multiple ASR models as feature extractors, but we found that they also did not improve results.

The second loss we apply is the Power Loss. This function aims to improve the accuracy of the reconstructed audio by attempting to match it with the real audio in the frequency domain. For this purpose, we use the L1 loss between the STFT magnitudes of the real and synthesized audio as follows:

$$L_{power} = \|\log\|STFT(x)\|^2 - \log\|STFT(\tilde{x})\|^2\|, \quad (6)$$

where x is the real waveform, \tilde{x} is the synthesized waveform from the same video and $STFT$ is the Short Time Fourier Transform with a window size of 25 ms, a hop size of 10 ms and frequency bins of size 512 (same parameters used for the power critic). We found that scaling the magnitudes using the natural logarithm and using an L1 Loss rather than the L2 Loss chosen in [36] greatly improve training stability and performance.

The third loss we apply is the MFCC Loss:

$$L_{MFCC} = \|MFCC(x) - MFCC(\tilde{x})\|, \quad (7)$$

where x is the real waveform, \tilde{x} is the synthesized waveform from the same video and $MFCC$ is the MFCC function which extracts 25 mel-frequency cepstral coefficients from the corresponding waveform. The objective of this loss lies in increasing the accuracy and intelligibility of the synthesized speech, given that MFCCs are known to be effective in ASR [18] and emotion recognition [22].

We adapt the function provided on an open-source repository¹.

Finally, the loss for the generator is described based on the losses mentioned above as:

$$L_G = \alpha_1 L_{G_{adv}} + \alpha_2 L_{PASE} + \alpha_3 L_{power} + \alpha_4 L_{MFCC}. \quad (8)$$

We tune the coefficients $\alpha_{1,2,3,4}$ by sequentially training multiple models on GRID (4 speakers, seen speaker split) and incrementally finding the coefficients that yield the best WER on the validation set. Through our search, we find that $\alpha_1 = 1$, $\alpha_2 = 140$, $\alpha_3 = 50$, $\alpha_4 = 0.4$ yield the best results.

¹<https://github.com/skaws2003/pytorch-mfcc>

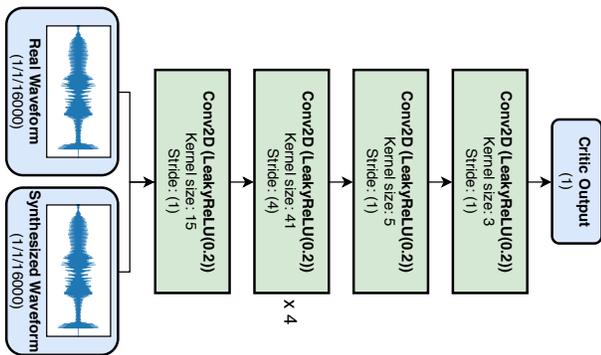


Fig. 4: Description of the layers in the waveform critic used to train our model.

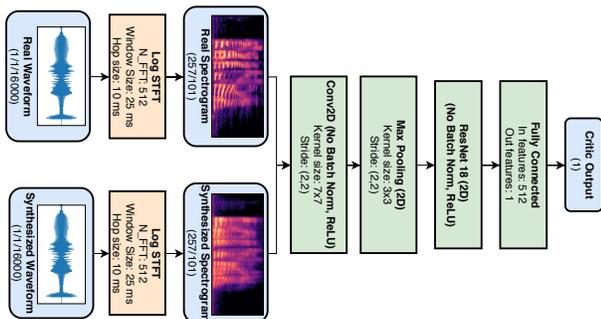


Fig. 5: Description of the layers in the power critic used to train our model, and the process used to extract spectrograms from waveform samples.

D. Training details

We use the Adam optimizer with a learning rate of 0.0001 and $\beta_1 = 0.5$, $\beta_2 = 0.99$ to train our generator and critics end-to-end. Given that the critics should be trained to completion before every generator training step, we perform 6 training steps on the critics before every training step of the generator. It should also be noted that we feed a one second clip randomly sampled from the real and synthesized audio to each of the critics, rather than the entire utterance. The other losses are computed using the entire real and synthesized utterances.

Additionally, we employ two data augmentation methods during training. Firstly, we apply random cropping on the input frame, producing a frame with roughly 90% of the original size. Furthermore, we apply horizontal flipping to each frame with a probability of 50%. These procedures help make our model more robust and provide regularization. During test time, the same cropping is performed on the center of the frame and no horizontal flipping is performed.

Training our model for each of the experiments generally takes approximately one week on an Nvidia RTX 2080 Ti GPU. Synthesizing a 3 second audio clip sampled at 16 kHz from 75 frames of video takes approximately 32 ms on the same high-end GPU, excluding pre-processing.

IV. DATASETS

For the purpose of this work, we use three separate audiovisual datasets to train and evaluate our model: GRID,

Corpus	Training set (clips / hours)	Validation set (clips / hours)	Test set (clips / hours)
GRID (4 speakers, seen speakers)	3576 / 2.98	210 / 0.18	210 / 0.18
GRID (33 speakers, seen speakers)	29584 / 24.65	1642 / 1.37	1641 / 1.37
GRID (33 speakers, unseen speakers)	15888 / 13.24	7000 / 5.83	9982 / 8.32
TCD-TIMIT (3 lipspeakers)	1014 / 1.64	57 / 0.09	60 / 0.09
LRW (full)	488763 / 157.49	25000 / 8.06	25000 / 8.06
FLRW 500 Words	112811 / 36.35	5878 / 1.89	5987 / 1.93
FLRW 100 Words	22055 / 7.11	1151 / 0.37	1144 / 0.37
FLRW 20 Words	4347 / 1.40	266 / 0.09	248 / 0.08

TABLE I: Number of speech clips and total number of hours of speech for each dataset used in our study.

TCD-TIMIT and LRW. GRID contains 33 speakers, each uttering 1000 short sentences composed of 6 simple words from a constrained vocabulary of 51 words. It is the most commonly used dataset for video-driven speech reconstruction [1, 14, 33, 43] due to the clean recording conditions and the limited vocabulary.

TCD-TIMIT is another audiovisual dataset composed of 62 speakers, three of which are trained lipspeakers. In order to compare with previous works [43], we only use the audiovisual data uttered by the three lipspeakers. Each lipspeaker utters 375 unique phonetically rich sentences, as well as two additional sentences which are uttered by all three speakers. This results in a total of 1 131 clips. The video/audio for this data is recorded in studio conditions with exceptional clarity given the particular speaking ability of the professional lipspeakers.

Finally, LRW contains roughly 500 000 speech samples (500 words, up to 1 000 clips per word) uttered by hundreds of different speakers, taken from television broadcasts. Due to the fact that these utterances are recorded ‘in the wild’ from a large variety of speakers, LRW presents a far more substantial challenge for speech reconstruction than the datasets mentioned above. Additionally, we use a subset of this corpus which keeps only the videos that are approximately frontal, i. e., videos with yaw, pitch and roll below 10 degrees. This leads to a corpus containing 124 676 samples in total and will be referred to as *F(rontal)LRW*. We also randomly select 20/100 words from this subset to experiment with different ranges of vocabulary during training/testing. These smaller sets will be referred to as *FLRW20* and *FLRW100*, respectively. Further statistics for each dataset are presented in Table I.

Rather than using the full face as input to our network, as is standard in other speech reconstruction works [1, 13, 43], we crop the mouth of the speaker, and use it as the input for every frame. We do this by performing face detection and alignment using dlib’s 68 landmark model [21], aligning each face to a reference mean face shape and extracting a mouth ROI (Region of Interest) from each frame. The mouth ROI is of size 128x74 for GRID and 96x96 for TCD-TIMIT and LRW.

V. EVALUATION METRICS

Although many metrics have been proposed for evaluating the quality of speech [30], it is widely acknowledged that none of the existing metrics are highly correlated with human perception. For this reason, we evaluate our speech reconstruction model using 4 objective metrics which capture different properties of the audio: PESQ, STOI, MCD and WER.

PESQ (Perceptual Evaluation of Speech Quality) [45] is an objective speech quality metric originally proposed for telephony quality assessment. It consists of a complex series of filters and transforms which result in a speech quality score. For the purposes of our work, we use this metric to measure how clean a speech signal is.

STOI (Short-Time Objective Intelligibility measure) [50] aims to measure how intelligible a speech signal is through a comparative DFT-based (Discrete Fourier Transform) approach. It has been found that it achieves close correlation to human intelligibility scores. In our experiments, we use this metric to measure the intelligibility of the reconstructed samples.

MCD (Mel-Cepstral Distance) [23] is designed to evaluate speech quality based on the cepstrum distance on the mel-scale. In practice, this is calculated as the distance between the MFCCs extracted from two signals. We find that it works quite reliably in measuring perceptual quality in our synthesized outputs, when compared to the original signal.

WER (Word Error Rate) measures the accuracy of a speech recognition system. It is calculated as:

$$WER = \frac{S + D + I}{N}, \quad (9)$$

where S is the number of substitutions, D is the number of deletions, I is the number of insertions and N is the total number of words in an utterance. For our work, we apply pre-trained ASR models to measure WER, which serves as an objective intelligibility metric for the reconstructed speech.

VI. RESULTS ON SEEN SPEAKERS

In this section, we present our experiments for seen speakers. For direct comparison with other works we use the same 4 speakers from GRID (1, 2, 4 and 29) as in [1, 33, 43, 53] and the 3 lipspeakers from TCD-TIMIT as in [43]. In order to investigate the impact of the number of speakers and the amount of training data, we also present results for all 33 speakers from the GRID dataset. We split the utterances in each of these datasets using a 90-5-5% ratio for training, validation and testing respectively similarly to [1, 33, 43, 53], such that the speakers in the validation and test sets are identical to the speakers seen in the training set (but the utterances are different). To measure the Word Error Rate (WER) for our GRID samples, we use a pre-trained ASR model (based on [31]) which was trained and tested on the full GRID dataset (using the split mentioned in Section VII), achieving a baseline of 4.23% WER on the test set. Audio samples, as well as spectrogram and waveform figures are presented on our website² for the experiments presented in sections VI, VII and VIII. Additionally, we present a publicly

available repository³ which can be used to reproduce each of the evaluation metrics presented in this work. We are also available to provide generated test samples for researchers hoping to reproduce or compare with our work.

A. Ablation Study

Results for the ablation study are shown in Table II. For this study, we only consider the 4 subjects from GRID presented above (1,2,4 and 29).

Firstly, we observe that each of the three comparative losses L_{PASE} , L_{power} and L_{MFCC} yield considerable improvements in the verbal accuracy of samples (as shown by the WER), even when only one is removed. We can also observe that L_{MFCC} and L_{power} are particularly impactful on the MCD of the reported samples, which is unsurprising since this is an MFCC-based metric. On the other hand, it is clear that L_{PASE} is essential towards achieving high intelligibility, given its particular impact on STOI. Finally, all three losses also seem to positively impact the PESQ score, indicating an increase in overall audio clarity.

We can see that the simultaneous removal of L_{PASE} and L_{power} greatly decreases PESQ and STOI, indicating that these losses are particularly important towards the clarity of generated samples. We also show that the absence of L_{MFCC} and L_{power} sharply increases MCD, indicating that these two losses greatly increase the similarity between real and synthesized audio. On the other hand, this model maintains a WER below 10%, which means that L_{PASE} alone (together with the adversarial losses) can achieve intelligible audio. Finally, the removal of all three L1 losses results in realistic yet unintelligible audio. This is because the adversarial losses are the only objective used for training, and therefore there is no incentive for the network to learn the exact words corresponding to the input video.

We observe that the use of the waveform critic yields noticeable improvements through our metrics, particularly in WER and STOI, suggesting that its inclusion substantially increases intelligibility. Additionally, the power critic also yields moderate improvements in PESQ, STOI and WER. Finally, we observe that the removal of both critics results in substantially lower MCD and WER, but maintains PESQ and STOI at a similar value. This again indicates that our model can generate intelligible and accurate words without the adversarial losses. However, these synthesized samples lack realism, which drastically improves when the critics are used. To demonstrate this effect, readers are encouraged to listen to examples on our website².

We also experiment with using the full face as input, as this is commonly used in previous studies. Through this ablation, we show that using a cropped mouth region instead of the full face improves our results substantially regarding WER, effectively improving intelligibility. We also prove that the use of overlap improves all metrics slightly, suggesting that its purpose of minimizing the issue of frame transitions is beneficial towards output quality.

²<https://sites.google.com/view/video-to-speech/home>

³<https://github.com/miraodasilva/evalaudio>

Model	PESQ	STOI	MCD	WER
w/o L_{PASE}	2.06	0.597	26.44	8.97 %
w/o L_{power}	2.05	0.575	28.64	9.54 %
w/o L_{MFCC}	2.08	0.591	28.09	9.09 %
w/o L_{PASE} , w/o L_{power}	1.86	0.545	27.47	13.44 %
w/o L_{PASE} , w/o L_{MFCC}	2.02	0.589	28.82	13.33 %
w/o L_{MFCC} , w/o L_{power}	2.00	0.569	31.43	9.71 %
w/o L_{PASE} , w/o L_{power} , w/o L_{MFCC}	1.14	0.311	53.63	89.12 %
w/o waveform critic	2.07	0.583	26.66	8.47 %
w/o power critic	2.08	0.594	26.73	7.30 %
w/o waveform critic, w/o power critic	2.07	0.584	27.45	9.01 %
w/o overlap	2.06	0.590	26.73	7.40 %
w/ full face	2.07	0.596	26.46	9.94 %
full model	2.10	0.595	26.78	7.03 %

TABLE II: Ablation study performed on GRID for seen speaker speech reconstruction.

A qualitative comparison with other works can be seen in Figure 6. Compared to the real audio, our spectrogram is similar overall, but is slightly blurrier and fails to model some of the fine details in the frequency bins, especially in the higher frequencies. The model trained without adversarial critics features a much blurrier spectrogram than the full model, failing to reproduce even the lower frequency bands during voiced speech, highlighting the importance of adversarial training.

B. Comparison with Other Works

We compare our proposed model with previous works on the commonly used 4 GRID speakers as shown in Table III. We note that the metrics reported on Lip2Wav [43] are taken directly from their paper due to test samples not being publicly available, and that their WER was calculated using the Google Speech-to-Text (STT) API rather than our ASR model.

Regarding PESQ, it is clear that our model is superior to the previous approaches by a sizeable margin. This suggests that the quality of our synthesized speech is somewhat higher than past models. Our model also outperforms previous works on STOI, excluding Lip2Wav. This shows that our samples are more intelligible than most other approaches, but are outperformed by the robustness and consistency of the speech produced by Lip2Wav. Furthermore, our generated samples achieve a better MCD than previous works, indicating that our reconstructed audio is more accurate than previous approaches on the frequency domain. Finally, our work achieves the best WER out of all methods, which shows that our model is more accurate than any of the previous approaches by a large factor, outperforming our previous model by more than 10 %.

A qualitative comparison is shown in Figure 7, which displays waveforms, mel-frequency spectrogram, and mel-frequency spectrogram differences, i. e., the element-wise absolute difference between the real and synthesized spectrograms.

Method	PESQ	STOI	MCD	WER
Lip2Audspect [1]	1.81	0.425	63.88	46.36 %
GAN-based [53]	1.70	0.539	45.37	21.11 %
Vocoder-based [33]	1.90	0.553	46.64	22.14 %
Lip2Wav [43]	1.77	0.731	-	14.08 ^a %
Ours	2.10	0.595	26.78	7.03 %

^aReported using Google STT API.

TABLE III: Comparison between our model and previous works, using the GRID subset (4 speakers) with a seen speaker split.

Method	PESQ	STOI	MCD
Lip2Wav [43]	1.35	0.558	-
Ours	1.61	0.295	32.12

TABLE IV: Comparison between our model and Lip2Wav, using TCD-Timit (3 lipspeakers) with a seen speaker split.

This difference is calculated as:

$$\|MelSpec(x) - MelSpec(\tilde{x})\|, \quad (10)$$

where x is the real waveform and \tilde{x} is the synthesized waveform. Through the spectrograms, it is clear that Lip2Audspect is the least accurate in the frequency domain, failing to model many frequencies, particularly in the higher bands. The other three approaches are clearly more accurate, but all feature some inaccuracies during voiced speech and also noise in unvoiced segments. While [53] and [33] feature an excessive amount of low frequency noise, our model seems to accurately emulate the low amount of noise in the real audio and therefore achieves the least substantial spectrogram difference.

We also compare our model to Lip2Wav on TCD-TIMIT (3 lipspeakers) in Table IV. Once more, it is clear that our model outperforms Lip2Wav [43] on PESQ, but achieves lower performance on STOI, which indicates that our model produces clearer, yet somewhat less intelligible audio. Additionally, our samples achieve a reasonably low MCD, indicating moderate similarity in the frequency domain.

C. Performance as a Function of Training Set Size

For the purposes of this study, we use all 33 subjects from GRID and we report results as we vary the size of the training set from 20 % to 100 % in steps of 20 %. Results are shown in Table V. When compared to the results reported for GRID (4 speakers, seen split), we observe comparable performance for 33 speakers when using the full training set. This shows that our network adapts well to larger datasets and is able to model a large amount of speakers with no substantial drop in performance.

Regarding the models which are trained using a smaller subset of the training set, it is clear that the performance drops as the amount of training data is gradually reduced. However, it is worth highlighting that the overall performance remains moderately consistent, even when we use only 20 %

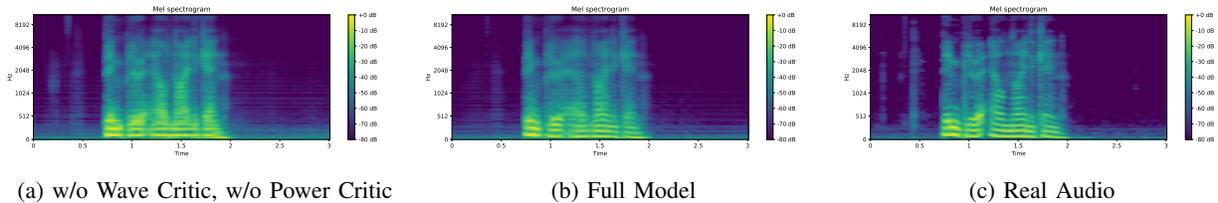


Fig. 6: Mel-frequency spectrograms taken from the audio reconstructed with our seen speaker ablation models. The clip we present is from GRID, speaker 1, utterance ‘Bin blue at L 9 again’.

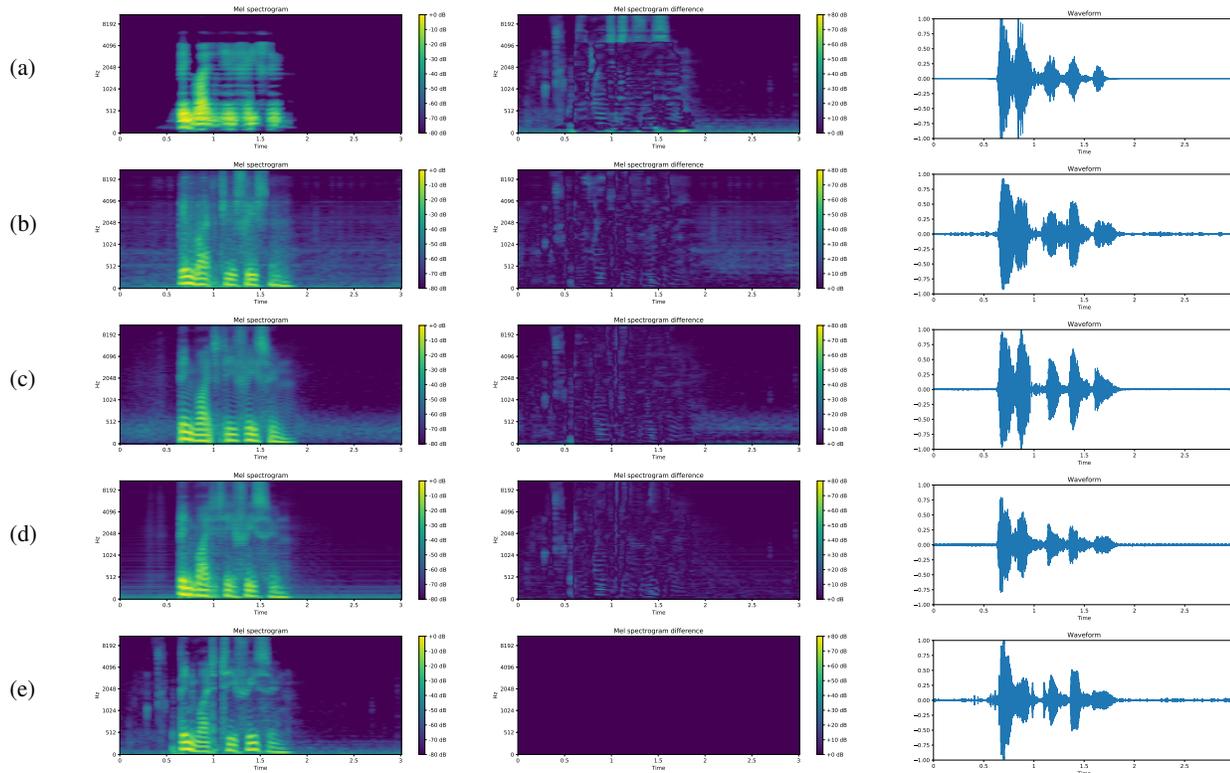


Fig. 7: Mel-frequency spectrograms (left), Mel-frequency spectrogram differences (middle) and waveforms (right) taken from the audio reconstructed with Lip2AudSpec [1] (a), our previous work [53] (b), a previous vocoder-based model [33] (c) and our model (d), as well as the real audio (e) – GRID, Speaker 1, utterance ‘Bin white at T 3 soon’. All models were trained on the same split of GRID (4 speakers, seen speaker split), as presented in our comparison.

% of Training Set	PESQ	STOI	MCD	WER
20 %	1.96	0.583	29.22	11.78 %
40 %	2.00	0.594	28.49	10.10 %
60 %	2.02	0.595	27.94	9.06 %
80 %	2.02	0.596	27.68	8.36 %
100 %	2.02	0.601	27.78	8.03 %

TABLE V: Study on the performance of our speech reconstruction model using varying training set sizes, using the full GRID seen speaker split mentioned in Section IV.

of the training data. This shows that our model adapts well to smaller datasets. We note that all 5 models were trained for the same amount of total training steps to avoid any bias in our comparative results.

VII. RESULTS ON UNSEEN SPEAKERS

In this section, we investigate the performance of the proposed approach on unseen speakers. For the purposes of this study, we use all speakers from the GRID dataset, using a 50-20-30 % split ratio similarly to [33, 53], such that there is no overlap between the speakers featured in the training, validation and test sets. To measure WER, we use the GRID pre-trained model mentioned in the previous section.

A. Ablation Study

In this study, we use all 33 GRID speakers. The results for the ablation study are shown in Table VI.

For this, task, we find that L_{power} provides the greatest impact on the quality of results, providing a substantial improvement in all metrics. On the other hand, L_{PASE} and L_{MFCC} show noticeable improvements in PESQ and

Model	PESQ	STOI	MCD	WER
w/o L_{PASE}	1.44	0.520	38.19	22.66 %
w/o L_{power}	1.37	0.503	39.59	24.32 %
w/o L_{MFCC}	1.44	0.518	39.03	21.70 %
w/o Waveform Critic, w/o Power Critic	1.43	0.516	38.48	22.82 %
Full Model	1.47	0.523	37.91	23.13 %

TABLE VI: Ablation study performed on GRID for unseen speaker speech reconstruction.

Method	PESQ	STOI	MCD	WER
GAN-based [53]	1.24	0.470	51.28	37.10 %
Vocoder-based [33]	1.23	0.477	55.02	55.23 %
Ours	1.47	0.523	37.91	23.13 %

TABLE VII: Comparison between our current and previous model, using full GRID (33 speakers) with an unseen speaker split.

STOI, indicating that these losses contribute to the clarity and intelligibility of the generated samples. Furthermore, we once more find that L_{MFCC} and L_{power} are particularly important towards achieving a low MCD, meaning that these losses are essential towards achieving accurate MFCCs in our synthesized samples.

Regarding the adversarial loss, we can see that, as reported in the seen speaker ablation, PESQ, STOI and MCD improve with the addition of the waveform and power critics. This suggests that these critics have a positive effect on the clarity and intelligibility of samples, and that the accuracy on the frequency domain is improved as well. However, we observe that the WER remains at a similar value with the removal of both critics, indicating that the network is generally capable of reproducing the correct words from the corresponding video samples while relying only on the three proposed L1 losses.

B. Comparison with Other Works

We present our comparison with other works [33, 53] on the subject-independent split of GRID in Table VII. It is clear that our model outperforms previous works in all performance measures. Although, the improvement in PESQ and STOI compared to these works is not as emphatic as the gains reported for seen speakers, WER sees a very substantial reduction. This improvement in WER can easily be observed in our synthesized speech, and clearly shows that our model is far more consistent for this task than previous approaches. Furthermore, the observed MCD is substantially lower in our work, indicating that our synthesized speech yields more accurate spectrograms, which suggests a greater similarity between the content of real and synthesized samples.

C. Additional Experiments

Additionally, we present a study on silent speakers. For this experiment, we artificially produce a video of a speaker from the GRID corpus being silent for five seconds by feeding

Method	PESQ	STOI	MCD	WER
Lip2Wav [43]	1.20	0.543	-	34.20^a %
Ours	1.45	0.556	39.32	42.51 %

^aReported using Google STT API.

TABLE VIII: Comparison between our model and Lip2Wav, using the full LRW dataset.

Brownian noise into the facial animation model proposed in [54]. We then use this video as input for our model trained on the full GRID dataset (33 speakers, unseen speaker split). This aims to measure two distinct factors: firstly, our model’s ability to recognize a silent speaker and not produce any voiced speech; and secondly, the baseline noise that is present in the audio we synthesize with our network, which is clear to observe when the speaker is silent. As discussed in Figure 8, our model performs well in this scenario and produces minimal noise for this silent example.

VIII. RESULTS IN THE WILD

In this section, we investigate the performance of the proposed approach on utterances recorded ‘in the wild’. For this purpose, we use the full LRW dataset, and its subsets FLRW 500 Words, FLRW 100 Words and FLRW 20 Words, which are introduced in Section IV. We split the utterances using the default split for LRW (90-5-5 % ratio), such that there is no overlap between the utterances in the training, validation and test sets. To measure the Word Error Rate (WER) for our samples, we use a pre-trained model (based on [40]) which was trained and tested on full LRW using the same split, and achieve a baseline WER of 1.68 % on the test set.

A. Comparison with Other Works

Our comparison with Lip2Wav [43] on LRW (500 Words) is presented in Table VIII. We compare our model to Lip2Wav on LRW (500 Words), in order to compare our model’s performance “in the wild” to this recent work. Our work shows a great improvement in PESQ compared to Lip2Wav, which suggests that our samples are able to achieve a superior clarity in this regard. On the other hand, our STOI is very similar to the one reported in Lip2Wav, achieving a slight edge which could indicate a minor improvement in intelligibility.

B. Performance for Different Subsets

In order to demonstrate our model’s ability to reconstruct speech in less constrained conditions, we experiment with the LRW dataset, as well as some of its subsets. These subsets present increasing degrees of challenge, culminating with the full LRW dataset which presents the greatest challenge given its large vocabulary and large variance in video perspective.

Regarding the experiments with frontal LRW, we observe that our model maintains a similar overall quality of outputs for larger vocabularies, as demonstrated by the consistency in PESQ, STOI and MCD. However, it is clear that the more difficult task presented by larger vocabularies yields a decrease

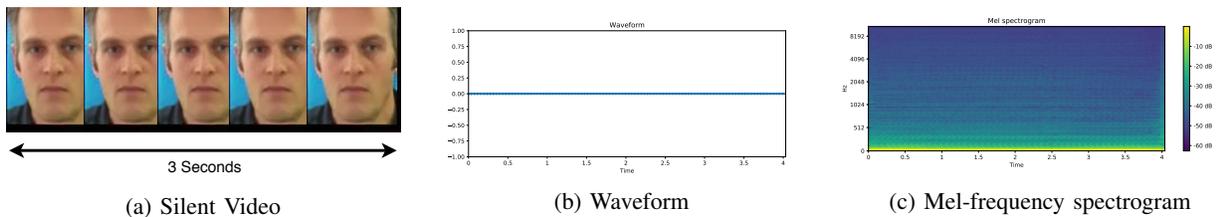


Fig. 8: The spectrogram and waveform for the audio produced by our model for a video of a silent speaker (Speaker 2 from GRID) are portrayed in (a). As displayed in the waveform (b), the audio is almost completely silent, disregarding some low frequency noise which is highlighted in the spectrogram (c). This shows that our model is robust to the scenario of silent speakers and produces minimal baseline noise under these circumstances. This audio sample is also available on our website².

Corpus	PESQ	STOI	MCD	WER
FLRW 20 Words	1.43	0.523	43.87	25.00 %
FLRW 100 Words	1.40	0.528	41.56	36.54 %
FLRW 500 Words	1.44	0.555	39.72	44.28 %
LRW 500 Words	1.45	0.556	39.32	42.51 %

TABLE IX: Study on the performance of our speech reconstruction model for the three subsets of LRW mentioned in Section IV, as well as the full LRW dataset.

in the average accuracy of samples, shown by the increasing WER. This implies that our model scales well with larger datasets, but has difficulties in adapting to larger vocabularies in very unconstrained and inconsistent environments. Even still, the word error rate reported for FLRW 20 Words is noticeably low, implying that our model can realistically reconstruct speech for hundreds of different speakers, even under such ‘wild’ conditions. Finally, we found that the full LRW dataset yields a better performance than our full frontal subset (FLRW 500 Words). Although we expected the frontal data to provide an easier task for the network during training and testing, this result shows that the network benefits strongly from a larger training set, even if the visual data is less consistent.

IX. CONCLUSION

In this work, we have presented our end-to-end video-to-waveform synthesis model using a generative adversarial network with two critics on waveform and spectrogram. First, we showed through an ablation study on GRID that the use of our losses, adversarial critics and other choices in training methodology provide a positive impact on the quality of our results for both seen and unseen speaker video-to-speech. Furthermore, we demonstrated through our experiments on LRW that our model is able to generate intelligible speech for videos recorded entirely in the wild by hundreds of different speakers. Finally, we compared our model to previous video-to-speech models and found that it produces the best results on most metrics for GRID and LRW and achieves state-of-the-art performance on PESQ for TCD-TIMIT.

We observed that the choice of good critics as well as adequate comparative losses is fundamental towards obtaining realistic results. Therefore, we believe that the pursuit of alternative loss functions (including different adversarial losses)

is a promising option for future work. Additionally, we believe that there would be substantial benefit in experimenting with a speaker embedding as input to the generator, in addition to the video, in order to generalize to unseen speakers with a more accurate voice profile, as proposed in [43, 46]. Finally, extending our model towards other practical applications such as speech inpainting i. e., reconstructing missing audio segments in an audiovisual stream, would be a promising research pursuit in order to show the empirical value of video-to-speech synthesis.

ACKNOWLEDGMENTS

All datasets used in the experiments and all training, testing and ablation studies have been conducted at Imperial College. Rodrigo Mira would like to thank Samsung for their continued support of his work on this project. Additionally, the authors would like to thank AWS for providing cloud computation resources for the experiments discussed in this paper.

REFERENCES

- [1] H. Akbari, H. Arora, L. Cao, and N. Mesgarani, “Lip2audspec: Speech reconstruction from silent lip movements video,” in *Proc. of ICASSP*, IEEE, 2018, pp. 2516–2520.
- [2] I. Anina, Z. Zhou, G. Zhao, and M. Pietikäinen, “Ouluvs2: A multi-view audiovisual database for non-rigid mouth motion analysis,” in *Proc. of FG*, IEEE, 2015, pp. 1–5.
- [3] M. Arjovsky, S. Chintala, and L. Bottou, “Wasserstein GAN,” *CoRR*, vol. abs/1701.07875, 2017.
- [4] Y. M. Assael, B. Shillingford, S. Whiteson, and N. de Freitas, “Lipnet: Sentence-level lipreading,” *CoRR*, vol. abs/1611.01599, 2016.
- [5] L. Chen, S. Srivastava, Z. Duan, and C. Xu, “Deep cross-modal audio-visual generation,” in *Proc. of MM*, ACM, 2017, pp. 349–357.
- [6] J. S. Chung and A. Zisserman, “Lip reading in the wild,” in *Asian Conference on Computer Vision*, 2016.
- [7] J. S. Chung and A. Zisserman, “Lip reading in the wild,” in *Proc. of ACCV*, ser. Lecture Notes in Computer Science, vol. 10112, 2016, pp. 87–103.
- [8] M. Cooke, J. Barker, S. Cunningham, and X. Shao, “An audio-visual corpus for speech perception and automatic speech recognition (I),” *The Journal of the Acoustical Society of America*, vol. 120, pp. 2421–4, 2006.
- [9] T. L. Cornu and B. Milner, “Reconstructing intelligible audio speech from visual speech features,” in *Proc. of Interspeech*, ISCA, 2015, pp. 3355–3359.
- [10] T. L. Cornu and B. Milner, “Generating intelligible audio speech from visual speech,” *IEEE ACM Trans. Audio Speech Lang. Process.*, vol. 25, no. 9, pp. 1751–1761, 2017.
- [11] C. Donahue, J. J. McAuley, and M. S. Puckette, “Adversarial audio synthesis,” in *ICLR*, 2019.
- [12] J. Donahue, S. Dieleman, M. Binkowski, E. Elsen, and K. Simonyan, “End-to-end adversarial text-to-speech,” *CoRR*, vol. abs/2006.03575, 2020.

- [13] A. Ephrat, T. Halperin, and S. Peleg, "Improved speech reconstruction from silent video," in *Proc. of ICCV*, IEEE, 2017, pp. 455–462.
- [14] A. Ephrat and S. Peleg, "Vid2speech: Speech reconstruction from silent video," in *Proc. of ICASSP*, IEEE, 2017.
- [15] I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. C. Courville, and Y. Bengio, "Generative adversarial networks," *CoRR*, vol. abs/1406.2661, 2014.
- [16] I. Gulrajani, F. Ahmed, M. Arjovsky, V. Dumoulin, and A. C. Courville, "Improved training of wasserstein gans," in *Proc. of NeurIPS*, 2017, pp. 5767–5777.
- [17] M. Gurban and J. Thiran, "Information theoretic feature extraction for audio-visual speech recognition," *IEEE Trans. Signal Process.*, vol. 57, no. 12, pp. 4765–4776, 2009.
- [18] W. Han, C. Chan, O. C. Choy, and K. Pun, "An efficient MFCC extraction method in speech recognition," in *International Symposium on Circuits and Systems*, 2006.
- [19] N. Harte and E. Gillen, "Tcd-timit: An audio-visual corpus of continuous speech," *IEEE Transactions on Multimedia*, vol. 17, no. 5, pp. 603–615, 2015.
- [20] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. of CVPR*, IEEE, 2016, pp. 770–778.
- [21] D. E. King, "Dlib-ml: A machine learning toolkit," *Journal of Machine Learning Research*, vol. 10, pp. 1755–1758, 2009.
- [22] K. V. Krishna Kishore and P. Krishna Satish, "Emotion recognition in speech using mfcc and wavelet features," in *(Proc. of IACC)*, 2013, pp. 842–847.
- [23] R. Kubichek, "Mel-cepstral distance measure for objective speech quality assessment," in *Proceedings of IEEE Pacific Rim Conference on Communications Computers and Signal Processing*, vol. 1, 1993, 125–128 vol.1.
- [24] K. Kumar, T. Chen, and R. M. Stern, "Profile view lip reading," in *Proc. of ICASSP*, IEEE, 2007, pp. 429–432.
- [25] K. Kumar, R. Kumar, T. de Boissiere, L. Gestin, W. Z. Teoh, J. Sotelo, A. de Brébisson, Y. Bengio, and A. C. Courville, "Melgan: Generative adversarial networks for conditional waveform synthesis," in *Proc. of NeurIPS*, 2019, pp. 14881–14892.
- [26] Y. Kumar, M. Aggarwal, P. Nawal, S. Satoh, R. R. Shah, and R. Zimmermann, "Harnessing AI for speech reconstruction using multi-view silent video feed," in *Proc. of MM*, ACM, 2018.
- [27] Y. Kumar, R. Jain, K. M. Salik, R. R. Shah, Y. Yin, and R. Zimmermann, "Lipper: Synthesizing thy speech using multi-view lipreading," in *Proc. of AAAI*, 2019, pp. 2588–2595.
- [28] Y. Kumar, R. Jain, K. M. Salik, R. R. Shah, R. Zimmermann, and Y. Yin, "Mylipper: A personalized system for speech reconstruction using multi-view visual feeds," in *Proc. of ISM*, IEEE, 2018, pp. 159–166.
- [29] J. Li, V. Lavrukhin, B. Ginsburg, R. Leary, O. Kuchaiev, J. M. Cohen, H. Nguyen, and R. T. Gadde, "Jasper: An end-to-end convolutional neural acoustic model," in *Proc. of Interspeech*, G. Kubin and Z. Kacic, Eds., ISCA, 2019, pp. 71–75.
- [30] P. C. Loizou, "Speech quality assessment," in *Multimedia Analysis, Processing and Communications*, vol. 346, 2011, pp. 623–654.
- [31] P. Ma, S. Petridis, and M. Pantic, "Investigating the lombard effect influence on end-to-end audio-visual speech recognition," in *INTER-SPEECH*, 2019.
- [32] A. L. Maas, Q. V. Le, T. M. O'Neil, O. Vinyals, P. Nguyen, and A. Y. Ng, "Recurrent neural networks for noise reduction in robust ASR," in *Proc. of Interspeech*, ISCA, 2012, pp. 22–25.
- [33] D. Michelsanti, O. Slizovskaia, G. Haro, E. Gómez, Z. Tan, and J. Jensen, "Vocoder-based speech synthesis from silent videos," in *Proc. of Interspeech*, H. Meng, B. Xu, and T. F. Zheng, Eds., ISCA, 2020, pp. 3530–3534.
- [34] M. Morise, F. Yokomori, and K. Ozawa, "WORLD: A vocoder-based high-quality speech synthesis system for real-time applications," *IEICE Trans. Inf. Syst.*, vol. 99-D, no. 7, pp. 1877–1884, 2016.
- [35] A. van den Oord, S. Dieleman, H. Zen, K. Simonyan, O. Vinyals, A. Graves, N. Kalchbrenner, A. W. Senior, and K. Kavukcuoglu, "Wavenet: A generative model for raw audio," in *Proc. of the ISCA Speech Synthesis Workshop*, ISCA, 2016, p. 125.
- [36] A. van den Oord, Y. Li, I. Babuschkin, K. Simonyan, O. Vinyals, K. Kavukcuoglu, G. van den Driessche, E. Lockhart, L. C. Cobo, F. Stimberg, N. Casagrande, D. Grewe, S. Noury, S. Dieleman, E. Elsen, N. Kalchbrenner, H. Zen, A. Graves, H. King, T. Walters, D. Belov, and D. Hassabis, "Parallel wavenet: Fast high-fidelity speech synthesis," in *Proc. of ICML*, PMLR, 2018, pp. 3915–3923.
- [37] A. Owens, P. Isola, J. H. McDermott, A. Torralba, E. H. Adelson, and W. T. Freeman, "Visually indicated sounds," in *Proc. of CVPR*, IEEE, 2016, pp. 2405–2413.
- [38] S. Pascual, M. Ravanelli, J. Serrà, A. Bonafonte, and Y. Bengio, "Learning problem-agnostic speech representations from multiple self-supervised tasks," in *Proc. of Interspeech*, G. Kubin and Z. Kacic, Eds., ISCA, 2019, pp. 161–165.
- [39] S. Petridis, Z. Li, and M. Pantic, "End-to-end visual speech recognition with LSTM," in *Proc. of ICASSP*, IEEE, 2017, pp. 2592–2596.
- [40] S. Petridis, T. Stafylakis, P. Ma, F. Cai, G. Tzimiropoulos, and M. Pantic, "End-to-end audiovisual speech recognition," in *Proc. of ICASSP*, IEEE, 2018, pp. 6548–6552.
- [41] S. Petridis, Y. Wang, Z. Li, and M. Pantic, "End-to-end multi-view lipreading," in *BMVC*, BMVA, 2017.
- [42] S. Petridis, Y. Wang, P. Ma, Z. Li, and M. Pantic, "End-to-end visual speech recognition for small-scale datasets," *Pattern Recognit. Lett.*, vol. 131, pp. 421–427, 2020.
- [43] K. R. Prajwal, R. Mukhopadhyay, V. P. Namboodiri, and C. V. Jawahar, "Learning individual speaking styles for accurate lip to speech synthesis," in *Proc. of CVPR*, IEEE, 2020, pp. 13793–13802.
- [44] M. Ravanelli, J. Zhong, S. Pascual, P. Swietojanski, J. Trmal, and Y. Bengio, "Multi-task self-supervised learning for robust speech recognition," in *Proc. of ICASSP*, IEEE, 2020, pp. 6989–6993.
- [45] A. W. Rix, J. G. Beerends, M. P. Hollier, and A. P. Hekstra, "Perceptual evaluation of speech quality (pesq)-a new method for speech quality assessment of telephone networks and codecs," in *Proc. of ICASSP*, 2001.
- [46] J. Shen, R. Pang, R. J. Weiss, M. Schuster, N. Jaitly, Z. Yang, Z. Chen, Y. Zhang, Y. Wang, R. Ryan, R. A. Saurous, Y. Ajiomyrgiannakis, and Y. Wu, "Natural TTS synthesis by conditioning wavenet on MEL spectrogram predictions," in *Proc. of ICASSP*, IEEE, 2018, pp. 4779–4783.
- [47] B. Shillingford, Y. M. Assael, M. W. Hoffman, T. Paine, C. Hughes, U. Prabhu, H. Liao, H. Sak, K. Rao, L. Bennett, M. Mulville, M. Denil, B. Coppin, B. Laurie, A. W. Senior, and N. de Freitas, "Large-scale visual speech recognition," in *Proc. of Interspeech*, ISCA, 2019, pp. 4135–4139.
- [48] T. Stafylakis and G. Tzimiropoulos, "Combining residual networks with lstms for lipreading," in *Proc. of Interspeech*, ISCA, 2017, pp. 3652–3656.
- [49] D. Stewart, R. Seymour, A. Pass, and J. Ming, "Robust audio-visual speech recognition under noisy audio-video conditions," *IEEE Trans. Cybern.*, vol. 44, no. 2, pp. 175–184, 2014.
- [50] C. H. Taal, R. C. Hendriks, R. Heusdens, and J. Jensen, "An algorithm for intelligibility prediction of time-frequency weighted noisy speech," *IEEE Trans. Speech Audio Process.*, vol. 19, no. 7, pp. 2125–2136, 2011.
- [51] Y. Tabet and M. Boughazi, "Speech synthesis techniques. a survey," in *Proc. of the International Workshop on Systems, Signal Processing and their Applications*, WOSSPA, 2011.
- [52] S. Uttam, Y. Kumar, D. Sahrawat, M. Aggarwal, R. R. Shah, D. Mahata, and A. Stent, "Hush-hush speak: Speech reconstruction using silent videos," in *Proc. of Interspeech*, ISCA, 2019, pp. 136–140.
- [53] K. Vougioukas, P. Ma, S. Petridis, and M. Pantic, "Video-driven speech reconstruction using generative adversarial networks," in *Proc. of Interspeech*, G. Kubin and Z. Kacic, Eds., ISCA, 2019, pp. 4125–4129.
- [54] K. Vougioukas, S. Petridis, and M. Pantic, "Realistic speech-driven facial animation with gans," *Int. J. Comput. Vis.*, vol. 128, no. 5, pp. 1398–1413, 2020.
- [55] Y. Wang, R. J. Skerry-Ryan, D. Stanton, Y. Wu, R. J. Weiss, N. Jaitly, Z. Yang, Y. Xiao, Z. Chen, S. Bengio, Q. V. Le, Y. Ajiomyrgiannakis, R. Clark, and R. A. Saurous, "Tacotron: Towards end-to-end speech synthesis," in *Proc. of Interspeech*, ISCA, 2017, pp. 4006–4010.
- [56] R. Yamamoto, E. Song, and J. Kim, "Parallel wavenet: A fast waveform generation model based on generative adversarial networks with multi-resolution spectrogram," in *Proc. of ICASSP*, IEEE, 2020, pp. 6199–6203.
- [57] G. Zhao, M. Barnard, and M. Pietikäinen, "Lipreading with local spatiotemporal descriptors," *IEEE Trans. Multimedia*, vol. 11, no. 7, pp. 1254–1265, 2009.
- [58] H. Zhou, Z. Liu, X. Xu, P. Luo, and X. Wang, "Vision-infused deep audio inpainting," in *Proc. of ICCV*, CVF / IEEE, 2019, pp. 283–292.
- [59] Y. Zhou, Z. Wang, C. Fang, T. Bui, and T. L. Berg, "Visual to sound: Generating natural sound for videos in the wild," in *Proc. of CVPR*, IEEE, 2018, pp. 3550–3558.
- [60] X. Zhu, G. Beauregard, and L. L. Wyse, "Real-time signal estimation from modified short-time fourier transform magnitude spectra," *IEEE Trans. Audio, Speech & Language Processing*, vol. 15, no. 5, pp. 1645–1653, 2007.