

WAV2PIX: SPEECH-CONDITIONED FACE GENERATION USING GENERATIVE ADVERSARIAL NETWORKS

Amanda Duarte^{†*}, Francisco Roldan^{*}, Miquel Tubau^{*}, Janna Escur^{*}, Santiago Pascual^{*}
Amaia Salvador^{*}, Eva Mohedano[‡], Kevin McGuinness[‡], Jordi Torres^{†*}, Xavier Giro-i-Nieto^{†*}

^{*} Universitat Politècnica de Catalunya, Barcelona, Catalunya/Spain

[†] Barcelona Supercomputing Center, Spain

[‡] Insight Centre for Data Analytics - DCU, Ireland

ABSTRACT

Speech is a rich biometric signal that contains information about the identity, gender and emotional state of the speaker. In this work, we explore its potential to generate face images of a speaker by conditioning a Generative Adversarial Network (GAN) with raw speech input. We propose a deep neural network that is trained from scratch in an end-to-end fashion, generating a face directly from the raw speech waveform without any additional identity information (e.g. reference image or one-hot encoding). Our model is trained in a self-supervised approach by exploiting the audio and visual signals naturally aligned in videos. With the purpose of training from video data, we present a novel dataset collected for this work, with high-quality videos of youtubers with notable expressiveness in both the speech and visual signals.

Index Terms— deep learning, adversarial learning, face synthesis, computer vision

1. INTRODUCTION

Audio and visual signals are the most common modalities used by humans to identify other humans and sense their emotional state. Features extracted from these two signals are often highly correlated, allowing us to imagine the visual appearance of a person just by listening to their voice, or build some expectations about the tone or pitch of their voice just by looking at a picture of the speaker. When it comes to image generation, however, this multimodal correlation is still under-explored.

In this paper, we focus on cross-modal visual generation, more specifically, the generation of facial images given a speech signal. Two recent approaches have recently popularized this research venue [1, 2]. Chung *et al.* [1] present a method for generating a video of a talking face starting from audio features and an image of him/her (identity). Suwajanakorn *et al.* focus on animating a point-based lip model to later synthesize high quality videos of President Barack Obama [2]. Unlike the aforementioned works, however, we aim to generate the whole face image at pixel level, conditioning only on the raw speech signal (*i.e.* without the use of any hand-crafted features) and without requiring any previous knowledge (e.g. speaker image or face model).

To this end, we propose a conditional generative adversarial model (shown in Figure 1) that is trained using the aligned audio and video channels in a self-supervised way. For learning such a model, high quality, aligned samples are required. This makes the most commonly used datasets such as *Lip Reading in the wild* [3], or *VoxCeleb* [4] unsuitable for our approach, as the position of

the speaker, the background, and the quality of the videos and the acoustic signal can vary significantly across different samples. We therefore built a new video dataset from YouTube, composed of videos uploaded to the platform by well-established users (commonly known as *youtubers*), who recorded themselves speaking in front of the camera in their personal home studios. Such videos are usually of high quality, with the faces of the subject featured in a prominent way and with notable expressiveness in both the speech and face. Hence, our main contributions can be summarized as follows:

1) We present a conditional GAN that is able to generate face images directly from the raw speech signal, which we call *Wav2Pix*.

2) We present a manually curated dataset of videos from youtubers, that contains high-quality data with notable *expressiveness* in both the speech and face signals.

3) We show that our approach is able to generate realistic and diverse faces.

The developed model, software and dataset are publicly released¹.

2. RELATED WORKS

Generative Adversarial Networks. (GANs) [5] are a state of the art deep generative model that consist of two networks, a Generator G and a Discriminator D , playing a min-max game against each other. This means both networks are optimized to fulfill their own objective: G has to generate realistic samples and D has to be good at rejecting G samples and accepting real ones. This joint learning adversarial process lasts for as long as G begins generating samples which are as good enough as to fool D into making as many mistakes as possible. The way Generator can create novel data mimicking real one is by mapping samples $z \in \mathbb{R}^n$ of arbitrary dimensions coming from some simple prior distribution \mathcal{Z} to samples x from the real data distribution \mathcal{X} (in this case we work with images, so $\mathbf{x} \in \mathbb{R}^{w \times h \times c}$ where $w \times h$ are spatial dimensions width and height and c is the amount of channels). This means each \mathbf{z} forward is like sampling from \mathcal{X} . On the other hand the discriminator is typically a binary classifier as it distinguishes *real* samples from *fake* ones generated by G . One can further condition G and D on a variable $e \in \mathbb{R}^k$ of arbitrary dimensions to derive the conditional GANs [6] formulation, with the conditioning variable being of any type, *e.g.* a class label or text captions [7]. In our work, we generate images conditioned on raw speech waveforms.

¹<https://imatge-upc.github.io/wav2pix/>

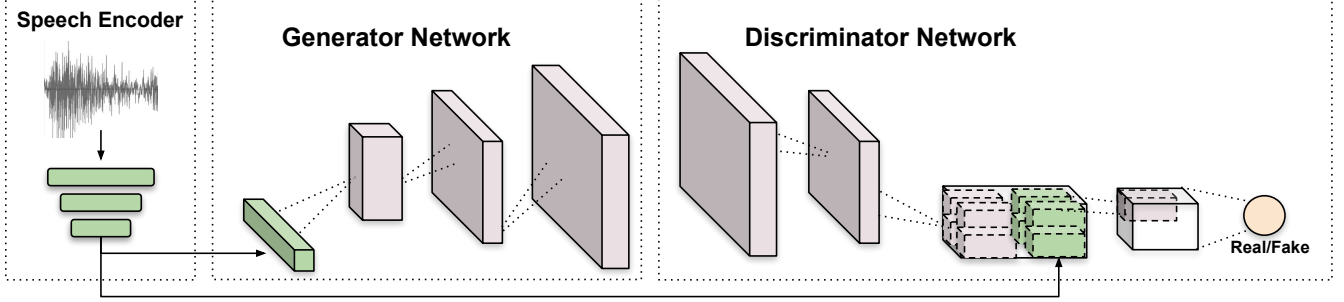


Fig. 1. The overall diagram of our speech-conditioned face generation GAN architecture. The network consists of a speech encoder, a generator and a discriminator network. An audio embedding (green) is used by both the generator and discriminator, but its error is just back-propagated at the generator. It is encoded and projected to a lower dimension (vector of size 128). Pink blocks represent convolutional/deconvolutional stages.

Numerous improvements to the GANs methodology have been presented lately. Many focusing on stabilizing the training process and enhance the quality of the generated samples [8, 9]. Others aim to tackle the vanishing gradients problem due to the sigmoid activation and the log-loss in the end of the classifier [10, 11, 12]. To solve this, the least-squares GAN (LSGAN) approach [12] proposed to use a least-squares function with binary coding (1 for real, 0 for fake). We thus use this conditional GAN variant with the objective function is given by:

$$\min_D V_{\text{LSGAN}}(D) = \frac{1}{2} \mathbb{E}_{\mathbf{x}, \mathbf{e} \sim p_{\text{data}}(\mathbf{x}, \mathbf{e})} [(D(\mathbf{x}, \mathbf{e}) - 1)^2] + \frac{1}{2} \mathbb{E}_{\mathbf{z} \sim p_z(\mathbf{z}), \mathbf{e} \sim p_{\text{data}}(\mathbf{e})} [D(G(\mathbf{z}, \mathbf{e}), \mathbf{e})^2]. \quad (1)$$

$$\min_G V_{\text{LSGAN}}(G) = \frac{1}{2} \mathbb{E}_{\mathbf{z} \sim p_z(\mathbf{z}), \mathbf{y} \sim p_{\text{data}}(\mathbf{y})} [(D(G(\mathbf{z}, \mathbf{e}), \mathbf{e}) - 1)^2], \quad (2)$$

Multi-modal generation. Data generation across modalities is becoming increasingly popular [13, 7, 14, 15]. Several works [7, 15] present different approaches for synthesizing realistic images given a text description. Recently, a number of approaches combining audio and vision have appeared, with tasks such as generating speech from a video [16] or generating images from audio/speech [17]. In this paper we will focus on the latter.

Most works on audio conditioned image generation adopt non end-to-end approaches and exploit previous knowledge about the data. Typically, speech has been encoded with handcrafted features which have been very well engineered to represent human speech. At the visual part, point-based models of the face [18] or the lips [2] have been adopted. In contrast to that, our network is trained entirely end-to-end solely from raw speech to generate image pixels.

A direct synthesis of facial pixels was obtained in [1] with a discriminative model whose input were a pair of audio features and a visual example of the face to predict. In that case, the model had the help of an additional identity information (image of the speaker) to help in the prediction, so the network learned how to modify this input to match with the speech utterance. Following a similar architecture, a generative model trained with adversarial training was proposed in [19]. In this case, they introduced a temporal regularization to improve the smoothness of the output video sequence. Our work differs from theirs in that we use raw speech instead of hand-crafted features, and we do not need any image of the speaker as all identity information is extracted from the speech only.

3. YOUTUBERS DATASET

In this section we describe the multi-stage pipeline adopted to collect the new audio-visual dataset of human speech used in this work. We collected videos uploaded to YouTube by well-established users (so-called *youtubers*), who tend to record themselves speaking in front of the camera in a well controlled environment. Such videos are usually of high quality, with the faces of the subject featured in a prominent way and with notable expressiveness in both the speech and face. The Youtubers dataset is composed of two sets: the complete noisy dataset automatically generated, and a clean subset which was manually curated to obtain high quality data.

In total we collected 168,796 seconds of speech with the corresponding video frames, and cropped faces from a list of 62 youtubers active during the past few years. The dataset is gender balanced and manually cleaned keeping 42,199 faces, each with an associated 1-second speech chunk. The pipeline used for downloading and pre-processing the full dataset is summarized in the next paragraphs. In Figure 2, and the key stages are discussed in the following paragraphs:

Youtubers collection and downloading: A list of 62 different Spanish speaker *youtubers* was built, consisting on 29 males and 33 females from different ethnicity and accents.

Audio preprocessing: The audio was originally downloaded in Advance Audio Coding (AAC) format at 44100 Hz and stereo and converted to WAV, as well as re-sampled to 16 kHz with 16 bits per sample and converted to mono.

Face Detection: The faces were detected using a Haar Feature-based Classifier [20] trained with frontal face features. We prevent the method from having false positives by taking only the most confident detection for each frame.

Audio/faces cropping: From each detection it is saved the bounding box coordinates, an image of the cropped face in BGR format, the full frame and a 4 seconds length speech frame, which encompasses 2 seconds ahead and behind the given frame. Moreover, we keep an identity (name) for each sample. We apply a pre-emphasis step to each speech frame and normalize it between $[-1, 1]$.

As stated in section 5 our model demonstrate a loss of performance when trained with noisy data. Thus, a part of the dataset was manually filtered to obtain the high-quality data required to improve the performance of our network. We took a subset of 10 identities, five female and five male, from our dataset and manually filtered them making sure that all faces were visually clear and all audios

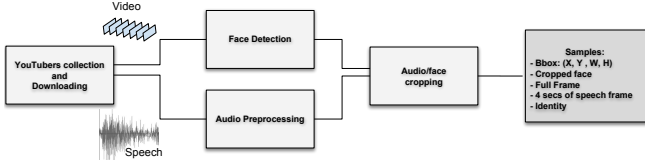


Fig. 2. High level representation of the data collection pipeline. Each detected face is associated with a 4 seconds length audio and the corresponding identity. Besides that we also kept the bounding box coordinates and the original image frame.

contain speech, so that all the silence and music parts were removed. As a result, the cleaned dataset contains a total of 4,860 images and audios (4 seconds length).

4. METHOD

Since our goal is to train a GAN conditioned on raw speech waveforms, our model is divided in three modules trained altogether end-to-end: a speech encoder, a generator network and a discriminator network described in the following paragraphs respectively. The speech encoder was adopted from the discriminator in [21], while both the image generator and discriminator architectures were inspired by [7]. The whole system was trained following a Least Squares GAN [12] scheme. Figure 1 depicts the overall architecture.

Speech Encoder: As mentioned in Section 2, many existing methods [1, 19, 2] require the extraction of handcrafted audio features before feeding the data into the neural network. This could limit the representation learning, as the audio information is extracted manually and not optimized for our generative task. In contrast, SEGAN [21] proposed a method for speech enhancement in which they do not work on the spectral domain, but at the waveform level. We coupled a modified version of the SEGAN discriminator Φ as input to an image generator G . Our speech encoder was modified to have 6 strided one-dimensional convolutional layers of kernel size 15, each one with stride 4 followed by LeakyReLU activations. Moreover we only require one input channel, so our input signal is $\mathbf{s} \in \mathbb{R}^{T \times 1}$, being $T = 16,384$ the amount of waveform samples we inject into the model (roughly one second of speech at 16 kHz). The aforementioned convolutional stack decimates this signal by a factor $4^6 = 4096$ while increasing the feature channels up to 1024. Thus, obtaining a tensor $f(\mathbf{s}) \in \mathbb{R}^{4 \times 1024}$ in the output of the convolutional stack f . This is flattened and injected into three fully connected layers that reduce the final speech embedding dimensions from $1024 \times 4 = 4096$ to 128, obtaining the vector $\mathbf{e} = \Phi(\mathbf{s}) \in \mathbb{R}^{128}$.

Image Generator Network: We take the speech embedding \mathbf{e} as input to generate images such that $\hat{\mathbf{x}} = G(\mathbf{e}) = G(\Phi(\mathbf{s}))$. The inference proceeds with two-dimensional transposed convolutions, where the input is a tensor $\mathbf{e} \in \mathbb{R}^{1 \times 1 \times 128}$ (an image of size 1×1 and 128 channels), based on [13]. The final interpolation can either be $64 \times 64 \times 3$ or $128 \times 128 \times 3$ just by playing with the amount of transposed convolutions (4 or 5). It is important to mention that we have no latent variable \mathbf{z} in G inference as it did not give much variance in predictions in preliminary experiments. To enforce the generative capacity of G we followed a dropout strategy at inference time inspired by [22].

In preliminary experiments, we found it convenient to add a sec-

ondary component to the loss of G : a *softmax* classifier trained over the given speech embedding. This classifier helped the whole network into preserving the identity of the speaker. The magnitude of the classification component is controlled by a new hyper-parameter λ . Therefore, the G loss, follows the LSGAN loss presented in Equation 2 with the addition of this weighted auxiliary loss for identity classification.

Image Discriminator Network: The Discriminator D is designed to process several layers of stride 2 convolution with a kernel size of 4 followed by a spectral normalization [23] and leakyReLU (apart from the last layer). When the spatial dimension of the discriminator is 4×4 , we replicate the speech embedding \mathbf{e} spatially and perform a depth concatenation. The last convolution is performed with stride 1 to obtain a D score as the output.

5. EXPERIMENTS

Model training: The *Wav2Pix* model was trained on the cleaned dataset described in Section 3 combined with a data augmentation strategy. In particular, we copied each image five times, pairing it with 5 different audio chunks of 1 second randomly sampled from the 4 seconds segment. Thus, we obtained $\approx 24\text{k}$ images and paired audio chunks of 1 second used for training our model. Our implementation is based on the PyTorch library [24] and trained on a GeForce Titan X GPU with 12GB memory. We kept the hyper-parameters as suggested in [7], changing the learning rate to 0.0001 in G and 0.0004 in D as suggested in [25]. We use ADAM solver [26] with momentum 0.1.

Evaluation: Figure 5 shows examples of generated images given a raw speech chunk, compared to the original image of the person who the voice belongs to. Different speech waveform produced by the same speaker were fed into the network to produce such images. Although the generated images are blurry, it is possible to observe that the model learns the person’s physical characteristics, preserving the identity, and present different face expressions depending on the input speech². Other examples from six different identities are presented in Figure 3.

To quantify the model’s accuracy regarding the identity preservation, we fine-tuned a pre-trained VGG-Face Descriptor network [27, 28] with our dataset. We predicted the speaker identity from the generated images of both the speech train and test partitions, obtaining an identification accuracy of 76.81% and 50.08%, respectively.

We also assessed the ability of the model to generate realistic faces, regardless of the true speaker identity. To have a more rigorous test than a simple Viola & Jones face detector [20], we measured the ability of an automatic algorithm [29] to correctly identify facial landmarks on images generated by our model. We define detection accuracy as the percentage of images where the algorithm is able to identify *all* 68 key-points. For the proposed model and all images generated for our test set, the detection accuracy is 90.25%, showing that in most cases the generated images retain the basic visual characteristics of a face. This detection rate is much higher than the identification accuracy of 50.08%, as in some cases the model confuses identities, or mixes some of them in a single face. Examples of detected faces together with their numbered facial landmarks can be seen in Figure 4.

²Some examples of images and it correspondent speech as well as more generated images are available at: <https://imatge-upc.github.io/wav2pix/>

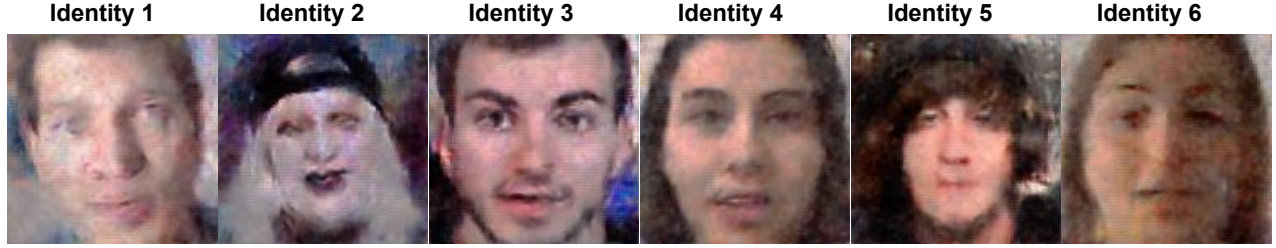


Fig. 3. Generated samples conditioned to raw speech produced by our model.

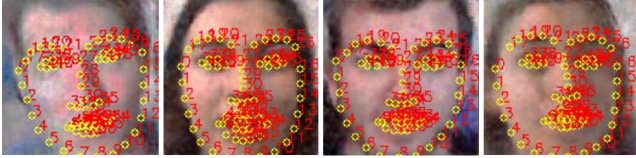


Fig. 4. Examples of the 68 key-points detected on images generated by our model. Yellow circles indicate facial landmarks fitted to the generated faces, numbered in red fonts.



Fig. 5. Examples of generated faces compared to the original image of the person who the voice belongs to. In the generated images, we can observe that our model is able to preserve the physical characteristics and produce different face expressions. In the first row we can see examples of the *youtuber Javier Muñiz*. In the second row we can see examples of the *youtuber Jaime Altozano*.

Additional experiments: We also tried to generate faces with noisy speech, experiments that resulted in failures. Firstly, we used audio snippets from the same youtubers videos that presented background noise, silences or other people’s voice. Secondly, using the well known VoxCeleb 1 dataset [4], which contains a larger amount of images and identities but present a lower audio quality. In both cases, the quality of results was very poor, making it almost impossible to recognize faces. These results show the importance of having clean speech samples to train the proposed model.

We also observed a drop in performance when working with smaller speech chunks and lower image definitions. We observed a visual degradation when using audio chunks of 300 and 700 milliseconds, which was reflected in a decrease of the face detection rate. Detection accuracy when using 300 and 700 ms chunks was 81.16% and 89.12%, respectively, in both cases worse than the 90.25% accuracy achieved when using 1000 ms chunks. Figure 6 (left) shows examples of generated images for the three speech chunk lengths. Figure 6 (right) shows how using a lower definition of 64x64 pixels

instead of 128x128 results into blurrier images.

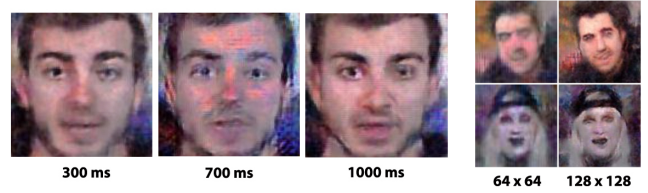


Fig. 6. (Left) Images generated for three speech chunk lengths. (Right) Images generated at two spatial resolutions.

6. CONCLUSIONS

In this work we introduced a simple yet effective cross-modal approach for generating images of faces given only a short segment of speech, and proposed a novel generative adversarial network variant that is conditioned on the raw speech signal.

As high-quality training data are required for this task, we further collected and curated a new dataset, the Youtubers dataset, that contains high quality visual and speech signals. Our experimental validation demonstrates that the proposed approach is able to synthesize plausible facial images with an accuracy of 90.25%, while also being able to preserve the identity of the speaker about 50% of the times. Our ablation experiments further showed the sensitivity of the model to the spatial dimensions of the images, the duration of the speech chunks and, more importantly, on the quality of the training data.

Further steps may address the generation of a sequence of video frames aligned with the conditioning speech, as well exploring the behaviour of the *Wav2Pix* when conditioned on unseen identities.

Acknowledgements

This research has received funding from “la Caixa” Foundation funded by the European Union’s Horizon 2020 research and innovation programme under the Marie Skłodowska-Curie grant agreement No. 713673. This research was partially supported by the Spanish Ministry of Economy and Competitiveness and the European Regional Development fund under contracts TEC 2015-69266-P and TEC 2016-75976-R (MINECO/FEDER, UE). This research was also partially supported by the Science Foundation Ireland (SFI) under grant number SFI/15/SIRG/3283. We gratefully acknowledge the support of NVIDIA Corporation with the donation of GPUs. We would like to acknowledge the youtubers *Javier Muñiz* and *Jaime Altozano* for providing us the rights of publishing their images for this research.

7. REFERENCES

- [1] Joon Son Chung, Amir Jamaludin, and Andrew Zisserman, “You said that?,” in *BMVC*, 2017.
- [2] Supasorn Suwajanakorn, Steven M Seitz, and Ira Kemelmacher-Shlizerman, “Synthesizing obama: learning lip sync from audio,” *ACM Transactions on Graphics (TOG)*, vol. 36, no. 4, pp. 95, 2017.
- [3] Joon Son Chung, Andrew W Senior, Oriol Vinyals, and Andrew Zisserman, “Lip reading sentences in the wild,” in *CVPR*, 2017, pp. 3444–3453.
- [4] Arsha Nagrani, Joon Son Chung, and Andrew Zisserman, “Voxceleb: a large-scale speaker identification dataset,” in *Interspeech*, 2017.
- [5] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio, “Generative adversarial nets,” in *NIPS*, 2014, pp. 2672–2680.
- [6] Mehdi Mirza and Simon Osindero, “Conditional generative adversarial nets,” *arXiv preprint arXiv:1411.1784*, 2014.
- [7] Scott Reed, Zeynep Akata, Xinchun Yan, Lajanugen Logeswaran, Bernt Schiele, and Honglak Lee, “Generative adversarial text to image synthesis,” in *ICML*, 2016.
- [8] Zizhao Zhang, Yuanpu Xie, and Lin Yang, “Photographic text-to-image synthesis with a hierarchically-nested adversarial network,” in *CVPR*, 2018.
- [9] Andrew Brock, Jeff Donahue, and Karen Simonyan, “Large scale gan training for high fidelity natural image synthesis,” in *ICLR*, 2019.
- [10] Martin Arjovsky, Soumith Chintala, and Léon Bottou, “Wasserstein generative adversarial networks,” in *ICML*, 2017.
- [11] David Berthelot, Thomas Schumm, and Luke Metz, “Began: boundary equilibrium generative adversarial networks,” *arXiv preprint arXiv:1703.10717*, 2017.
- [12] Xudong Mao, Qing Li, Haoran Xie, Raymond YK Lau, Zhen Wang, and Stephen Paul Smolley, “Least squares generative adversarial networks,” in *ICCV*. IEEE, 2017, pp. 2813–2821.
- [13] Alec Radford, Luke Metz, and Soumith Chintala, “Unsupervised representation learning with deep convolutional generative adversarial networks,” in *ICLR*, 2016.
- [14] Augustus Odena, Christopher Olah, and Jonathon Shlens, “Conditional image synthesis with auxiliary classifier gans,” in *ICML*, 2017.
- [15] Han Zhang, Tao Xu, Hongsheng Li, Shaoting Zhang, Xiaoqi Huang, Xiaogang Wang, and Dimitris Metaxas, “Stackgan: Text to photo-realistic image synthesis with stacked generative adversarial networks,” in *ICCV*, 2017.
- [16] Ariel Ephrat, Tavi Halperin, and Shmuel Peleg, “Improved speech reconstruction from silent video,” in *ICCV*, 2017.
- [17] Lele Chen, Sudhanshu Srivastava, Zhiyao Duan, and Chenliang Xu, “Deep cross-modal audio-visual generation,” in *Proceedings of the on Thematic Workshops of ACM Multimedia 2017*. ACM, 2017, pp. 349–357.
- [18] Tero Karras, Timo Aila, Samuli Laine, Antti Herva, and Jaakko Lehtinen, “Audio-driven facial animation by joint end-to-end learning of pose and emotion,” *ACM Transactions on Graphics (TOG)*, vol. 36, no. 4, pp. 94, 2017.
- [19] Konstantinos Vougioukas, Stavros Petridis, and Maja Pantic, “End-to-end speech-driven facial animation with temporal gans,” in *BMVC*, 2018.
- [20] Paul Viola and Michael Jones, “Rapid object detection using a boosted cascade of simple features,” in *Computer Vision and Pattern Recognition, 2001. CVPR 2001. Proceedings of the 2001 IEEE Computer Society Conference on*. IEEE, 2001, vol. 1, pp. I–I.
- [21] Santiago Pascual, Antonio Bonafonte, and Joan Serrà, “Segan: Speech enhancement generative adversarial network,” *Interspeech*, pp. 3642–3646, 2017.
- [22] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros, “Image-to-image translation with conditional adversarial networks,” in *CVPR*. IEEE, 2017, pp. 5967–5976.
- [23] Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida, “Spectral normalization for generative adversarial networks,” in *ICLR*, 2018.
- [24] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer, “Automatic differentiation in pytorch,” 2017.
- [25] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter, “Gans trained by a two time-scale update rule converge to a local nash equilibrium,” in *NIPS*, 2017, pp. 6626–6637.
- [26] Diederik P Kingma and Jimmy Ba, “Adam: A method for stochastic optimization,” in *ICLR*, 2015.
- [27] O. M. Parkhi, A. Vedaldi, and A. Zisserman, “Deep face recognition,” in *BMVC*, 2015.
- [28] Qiong Cao, Li Shen, Weidi Xie, Omkar M. Parkhi, and Andrew Zisserman, “Vggface2: A dataset for recognizing faces across pose and age,” 2017, vol. abs/1710.08092.
- [29] Vahid Kazemi and Josephine Sullivan, “One millisecond face alignment with an ensemble of regression trees,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2014, pp. 1867–1874.