

AUDIO2FACE: GENERATING SPEECH/FACE ANIMATION FROM SINGLE AUDIO WITH ATTENTION-BASED BIDIRECTIONAL LSTM NETWORKS

Guanzhong Tian¹, Yi Yuan², Yong Liu¹

¹ Institute of Cyber-systems and Control, Zhejiang University, ² Fuxi AI Lab, Netease

ABSTRACT

We propose an end to end deep learning approach for generating real-time facial animation from just audio. Specifically, our deep architecture employs deep bidirectional long short-term memory network and attention mechanism to discover the latent representations of time-varying contextual information within the speech and recognize the significance of different information contributed to certain face status. Therefore, our model is able to drive different levels of facial movements at inference and automatically keep up with the corresponding pitch and latent speaking style in the input audio, with no assumption or further human intervention. Evaluation results show that our method could not only generate accurate lip movements from audio, but also successfully regress the speaker's time-varying facial movements.

Index Terms— Animation, Long short-term memory network, Attention mechanism

1. INTRODUCTION

As an essential part of human machine interaction, facial animation has been applied on many applications, such as computer games, animated movies, man-machine conversation, virtual human agents, etc. Currently, great success has been achieved by vision-based facial capture techniques in synthesizing virtual human face. Vision-based performance capture rigs usually take advantage of motion sensors/markers and drive the animated face with the captured motion of a human face. Although the capture systems are able to provide a pretty acceptable results, active capture rigs tend to be expensive and time-consuming to use because of the elaborate setups and processing steps. Considering the constraints of facial capture methods, other input modalities, such as audio, may be used as an indispensable complementarity to infer facial motions. As a matter of fact, a growing number of animated film now use vision capture systems to produce just key animation, and rely on approaches based on audio or transcript to produce the vast majority of detail.

Although inferring facial motions and expressions from speech is a fundamental human reasoning functionality, it is extremely challenging for computer. On one hand, real-time applications such as inter-person telecommunication requires

strong adaptability of wide variability in user voices and stringent latency requirements. On the other hand, using just audio track to generate facial shapes seems unmanageable due to the inherent uncertainty of the problem— different facial expression may comes from exactly the same audio [1]. However, we try to address this challenge by using deep neural network, which has show the effectiveness and strong recovery capability by its widely spread on super resolution and image generation problems.

In this work, our primary goal is to recreate a plausible 3D virtual talking avatar that can make reasonable lip movements and then further generate natural facial movements. Our approach is designed exclusively on the basis of audio track, without any other auxiliary input such as images, which makes it increasingly challenging as we attempt to regress visual space from vocal sequence. Another challenge is that facial motions involve multifold activations of correlated regions on the geometric surface of face, which makes it difficult to generate life-like and consistent facial deformations for avatars. To recover as much details from the original input face with low latency, we adopt blendshape models to output the activation parameters of a fixed-topology facial mesh model. First, we extract handcrafted, high-level acoustic features, i.e. mel-frequency cepstral coefficients (MFCCs), from raw audio as the input to our network. Second, in order to map the long range context of the audio signal to effective face shape parameters, we employ bidirectional long short-term memory neural networks with attention mechanism to model meaningful acoustic feature representations, taking advantage of the locality and shift invariance in the frequency domain of audio sequence. Moreover, we combine attention embedded LSTM layers with fully-connected layers into an end-to-end framework, which aims at generating temporal transition of lip/facial movements, as well as spontaneous actions.

2. RELATED WROK

Face/Speech Animation Synthesis Quite a few methods have been developed to synthesize a face model from either transcript [2] [3] or speech [4] [5]. Prior to the advent of deep learning, most of the work was based on analyzing explicit knowledge about the phoneme content of languages, and then build a map between phonemes and target output face model

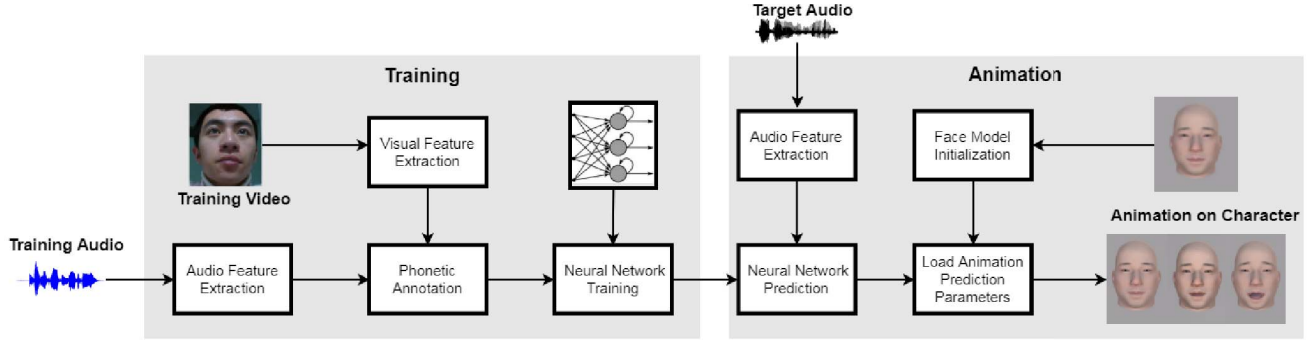


Fig. 1. Our proposed end-to end deep learning approach to learn a map from phoneme labels to facial/speech animation.

coefficients. The animated movements is then generated by the counterpart coefficients of input phonemes. There exists weaknesses of the linguistic-based method: the lack of a automatic way to generate facial animation from different languages and the incapability for responding to non-phoneme audio inputs. Due to these inherent constrains and the arise of deep learning models, neural networks have been successfully applied to speech synthesis [6] [7] and facial animation [8] [9] [10] with superior performance in recent years. The reason for the effectiveness of deep neural networks work lies in the capability of extracting the internal correlation and long-term relationship of high-dimensional input data, as well as the highly non-linear mapping from input features to target output space. Fan et al. [9] suggest to combine acoustic and text features to acquire AAM (active appearance model) coefficients of the mouth area, which then will be used to drive a life-like talking head. Taylor et al. [11] introduce a deep neural network(DNN) model to automatically regress AAM coefficients from input phonemes, which can be retargeted to drive lip movements for different 3D face models.

Recurrent Neural Networks Recurrent neural network (RNN) [12] is a mutant from traditional neural networks and mainly address the issue of sequence modeling. It is designed with loops to possess the ability of allowing information to persist and integrating temporal contextual information. Due to the problem of vanishing gradients, a new building unit for RNN called long short-term memory(LSTM) [13] appears and extend the memory of RNNs. In our work, we design the framework based on assumption that animation frames are connected and a model with longer memory could provide more details for the modeling of facial movements. Therefore, we choose LSTM network since it is born to help memorize long source input sequence in neural machine mapping.

Attention Mechanism Attention is, to some extent, motivated by how we pay attention to different regions of an image or correlate words in one sentence [14]. It has been applied in deep learning, for speech recognition, image caption and sentence summarization modeling [15] [16] [17], etc. A dilemma we are facing when design the network is that not all extracted

features from audio input contribute to the output animation vectors equally. Therefore, the attention mechanism is an ideal and important complement for our architecture, which can be used to train the model to learn to pay selective attention to audio sequence representations.

3. APPROACH

In this section, we introduce the architecture of our network, along with details on bidirectional LSTM and audio feature extraction. In order to fulfill our goal of generating plausible 3D facial/speech animation, we formulate three requirements for our synthesis approach: Lifelike, Real-time and Retargetable. Lifelike (Requirement 1) means the generated animations should be with high fidelity and accurately reflect complex speaking patterns present in visible speech motion. Real-time (Requirement 2) requires the system to generate animations with low latency. Retargetable (Requirement 3) expects the system to be able to retarget the generated animations to other character rig. To satisfy all of the requirements, we first take a vision-based rigs to capture the complex facial movements of actors and build an appropriate training set. Then we build several 3D blendshape-based cartoon face models with counterpart parameters to control different parts of the face: eyebrows, eyes, lip, jaw, etc. Further, we bring attention mechanism into bidirectional LSTM networks to map input audio feature to animation parameters. Fig. 1 depicts the pipeline of our system.

3.1. END-TO-END NETWORK ARCHITECTURE

Input a short span of audio, the goal of our framework is to correspondingly return a series of successive lip/facial movements. Our neural network first take in audio feature as input and feed it into bi-directional LSTM layers for high level feature extraction, then the attention layer can learn to pay attention to only the main relevant features and output abstract representations. Finally, we feed these abstract representations to an output network to generate the final 51 di-

mensional vectors to control the 3D face model. The output network is composed of a pair of dense layers that are used as linear transformer on the data. The first dense layer maps the set of input features to a linear time-varying basis vectors, and the second layer outputs the final blendshape parameter vectors as a weighted sum over the corresponding basis vectors. After the network is trained, we animate the face model by sliding a window over the vocal audio track and evaluate the network at every time step. Due to the bi-LSTM network and the embedded attention mechanism, the network itself has a memory of past input audio features and choose which one can affect the current animation frame.

3.2. Bidirectional LSTM

Our end-to-end deep neural network consists of two bidirectional LSTM layers followed with attention layer. The architecture of bidirectional LSTM model with attention can be seen in Fig 2. W is the weight matrices, where W_{xh}, W_{hh} are the input-hidden, hidden-hidden weight matrices, b_h is the hidden and output bias vectors and σ denotes the nonlinear activation function for hidden nodes. Bi-RNNs compute both forward state sequence \vec{h} and backward state sequence \overleftarrow{h} , as formulated below:

$$\vec{h}_t = \sigma(W_{xh} \vec{x}_i + W_{hh} \vec{h}_{t-1} + b_h) \quad (1)$$

$$\overleftarrow{h}_t = \sigma(W_{xh} \overleftarrow{x}_i + W_{hh} \overleftarrow{h}_{t-1} + b_h) \quad (2)$$

$$h_t = W_h \vec{h}_t + W_h \overleftarrow{h}_t + b \quad (3)$$

Next, Let H be a matrix consisting of output vectors $H = [h_1, h_2, \dots, h_t, \dots, h_T]$ that the LSTM layer produced, where T is the sentence length, then the representation y of the sequence is formed by a weighted sum of these output vectors through the attention layer:

$$\alpha = softmax(W^T \sigma(H)) \quad (4)$$

$$y = H\alpha^T \quad (5)$$

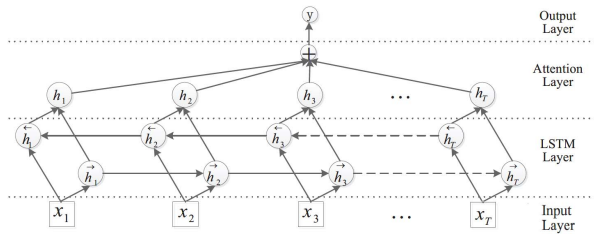


Fig. 2. Bi-RNNs with attention in our system

3.3. Audio Feature Extraction

The analysis step for most of the previous work on audio feature extraction is typically based on a specialized techniques such as Mel-frequency cepstral coefficients (MFCCs), linear predictive coding (LPC), etc. These techniques are widely adopted mainly due to the good linear separation of phonemes, as well as their fitness for hidden Markov models. In our work, we tested two different input data representations (MFCCs and LPC) and observed that the MFCC coefficients tend to be superior for our case.

Every input audio sequence in our system is synchronized to the corresponding video at 30 FPS and the sampling rate for audio track is at 44100Hz. As a result, there are 1470 audio samples for each corresponding video frame. In practice, we include additional samples from the previous video frame, so that for each video frame there is enough audio data to extract two windows of 33.33ms. The input audio window is divided into 64 audio frame with 2x overlap. Therefore, we use 2.13s worth of audio as input, i.e., 1.067s of past and future samples. By choosing this value, the network is able to capture relevant effects like phoneme coarticulation without providing too much data that might lead to high inference time and overfitting. Finally, we calculate $K=39$ Mel-frequency cepstral coefficients to yield a total of 64×39 scalars for the input audio window. In summary, the input feature vector for every video frame has 2496 dimensions.

4. TRAINING DETAILS

In this section, we describe the aspects relevant to training our network: how the training dataset were obtained, loss function design, and implementation details.

4.1. Training Dataset

In this work, we use 51 dimensional blendshape parameters to depict the overall shape of the whole face. Each parameter, range from 0 to 100, controls certain part of the avatar's face and the value of the parameter controls the range of actions. Fig 3 shows several examples of how parameters affect the face shape. We obtained the driven parameters for 3D facial model as training targets using the commercial Faceshift system, which employs depth cameras to directly capture dynamic RGB and depth information of the actor's facial movements. This system allows us to bypass the complex and expensive facial rigging while obtaining a reasonable and acceptable ground truth blendshape parameters, which animated control vertices for each frame. By processing information captured by depth camera, we can finally obtain the driving blendshape parameters to drive the pre-designed face model. These parameters - or more precisely the representations from a specific face shape - are the expecting outputs of our network when feeding a window of audio signal.



Fig. 3. Examples of face shape changes as certain blendshape parameter alters. From left to right: all parameters are set to zero, dimension 12 is set to 100, dimension 13 is set to 100, dimension 51 is set to 100

In general, our training sets consists of two parts: A female actor and a male actor. The female actor was required to read four scripts from different scenes with shorter length (each script last for about 15 minutes). In these four datasets (denoted as F1, F2, F3 and F4), we mainly focus on the task of speech animation and just record the lower part of female actor’s face. On the hand, we captured the whole facial shape changes of the male actor for intact facial animation when he was reading two longer scripts(each script last for about 30 minutes), denoted as M1 and M2. Normally, as the training data expands the inference quality will increase correspondingly. However, due to the cost of capturing high-quality training set, we choose to build a training set with acceptable size and we find that 60 minutes per actor has reach a pretty good result.

4.2. LOSS FUNCTION

Motivated by the unit selection in concatenative speech synthesis and considered the ambiguity nature of the optimization target, we design a specialized lost function that contains two losses: the target loss and the smooth loss:

$$L = \frac{1}{n} \sum_{i=1}^n w_1 L_t(y_i, y_i^p) + \sum_{i=2}^n w_2 L_s(y_{i-1}^p, y_i^p) \quad (6)$$

where y_i and y_i^p are the desired output and the output produced by the network, while L_t and L_s represent the target loss and the smoothness loss, and w_1 , w_2 are target weight and smooth weight, respectively.

Target loss We use the target loss to ensure that each output blendshape parameter is roughly correct. We adopt Huber loss as our target loss because it is less sensitive to outliers in dataset than MSE (mean square error) loss and is able to converge faster than both MSE loss and MAE (mean absolute error) loss to the minimum:

$$L_\delta(a) = \begin{cases} \frac{1}{2}(y_i - y_i^p)^2 & \text{for } |y_i - y_i^p| \leq \delta \\ \delta|y_i - y_i^p| - \frac{1}{2}\delta^2 & \text{otherwise} \end{cases} \quad (7)$$

δ is a hyperparameter can be tuned, which we set to 1 in practice.

Smooth Loss The target term is used to output a roughly correct result for animation task, but itself is insufficient to produce facial animation with high quality: As the coarticulation phenomenon and similarity between frames, training the network with the just target loss could lead to a considerable amount of temporal instability, as well as a weakly response to individual phonemes. Therefore, we bring in the smooth loss, taking the smoothness between adjoining selected frames into consideration. Inspired by the work of [9], we use cosine distance to measure the smooth cost between two adjacent frame sample:

$$L_s(y_i^p, y_j^p) = \cos(y_i^p, y_j^p) = \frac{y_i^p \cdot y_j^p}{\|y_i^p\| \cdot \|y_j^p\|} \quad (8)$$

where i, j represent the i th, j th animation frame.

4.3. Implementation Details

Our framework is implemented in Python, based on the deep learning toolkit pytorch. We train the deep models in 500 epochs, using Adam [18] with the default parameters, where learning rate is chosen as 0.0001 and the weights are initialized with a Gaussian random distribution. Each epoch processes all training samples in minibatches of 100. The dimension of MFCCs from audio feature extraction process is 39. We train and test our neural network on an Nvidia GTX 1080i.

5. EVALUATION

To assess the quality of our output, we conducted an empirical analysis evaluating our approach using both quantitative and subjective measures against several baselines. In the study, we compared the result of our method with video-based capture system from Faceshift and LSTM method [9], as well as classical HMM-based synthesis [4]. We choose audio clips randomly from the validation dataset of the corresponding actor as our test sets, which were not used at the training process.

5.1. Comparison with Other Methods

In our benchmark evaluation experiments, we evaluate all methods on our 6 self-build datasets mentioned above. We evaluate the performance of all the approaches using the same loss function, as mentioned above. HMM based method is denoted as HMM and long short-term memory networks regression is denoted as LSTM.

For the qualitative comparison, Fig 4 (a) shows the speech animation result of one example sequence from the test set. Experiment results reveal that our output speech-driven animation quality from the input sequence is quite plausible compared to ground truth. Fig 4 (b) is the result of our facial animation experiment. We can see that our model is able

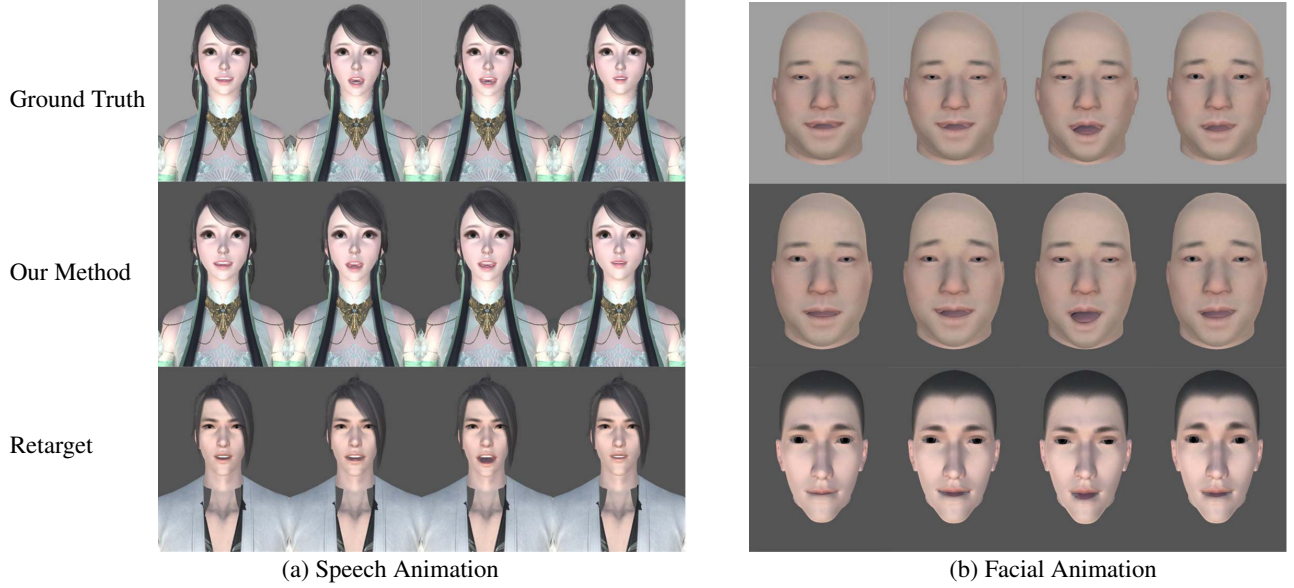


Fig. 4. Comparison of held-out vision-captured ground truth of the reference speaker compared with our reference model rendered predictions. Predicted mouth regions are rendered onto the original face for speech animation (a) and predicted overall facial regions are rendered for facial animation (b).

to effectively reproduce accurate face shapes on the test audio. However, there are certain limitations in our method: our network could not precisely follow blink pattern of the actor. This is an intractable problem because the speaker’s blink hardly has any relationship with their speech. We are considering fixing this by adding a constant blink pattern into the framework in the future work.

For the quantitative evaluation, we tested the performance of network topologies with different number of nodes (128, 256 and 512) and compared the performance of our method with that of HMM and LSTM, as listed in Table 1. The metrics in this table is root mean squared error (RMSE) of parameters over all video frames in the held-out test set:

$$\varepsilon = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - y_i^p)^2} \quad (9)$$

We can figure out from Table 1 that our approach consistently achieves the lowest RMSE. Results show that 256 nodes per layer could get an acceptable result for the speech animation task while 512 nodes perform superior on facial animation task. This is in accordance with our prediction because animation on entire face is more complex.

5.2. Retarget Test

After our model is trained, we retarget it to a different cartoon character. Fig. 4 shows the performance of our method when generalizing the predicted animation to a new charac-

ter’s output face model. We can figure that basically the facial movements and expressions can match between characters.

5.3. Run-Time Performance

To evaluate the run-time performance of our approach (model with 256 nodes) we perform experiments on a sample of random 10 sequences from the test sets with each sequence last for about 5 seconds. Each audio window in the animation last for 33.3ms because the frame rate is 30. The experiments were performed on an Intel Xeon(R) CPU E5-2680 v4 @ 2.40GHz × 16 CPU and an Nvidia GTX 1080i GPU. In general, it requires 0.14ms on average to process one audio window and update the blendshape parameters in the 3D morphable face model on average. Besides, it need an additional 0.52ms to extract audio feature for each audio window. As is illustrated above, the other deployments in the framework do not need extra time in the inference process. Therefore, the average time for the feature extraction and a forward pass of our neural network for one audio window is 0.68ms, which makes it possible for the implementation of real-time animation.

6. CONCLUSION

In this paper, we introduce a deep learning approach for speech/facial animation from exclusively audio input. Our deep framework is based on bi-directional long short-term memory recurrent neural network, combining with attention

Table 1. RMSE when Compared with typical HMM and LSTM methods on Different Test Sets

Algorithm	Mean	F1	F2	F3	F4	M1	M2
HMM	0.3222	0.2345	0.2593	0.3211	0.4123	0.3756	0.3301
LSTM	0.3046	0.2132	0.2410	0.2914	0.4038	0.3644	0.3140
Our Method(128 nodes)	0.2965	0.1875	0.2365	0.2876	0.3714	0.3751	0.3210
Our Method(256 nodes)	0.2891	0.1817	0.2336	0.2819	0.3648	0.3624	0.3104
Our Method(512 nodes)	0.2883	0.1821	0.2328	0.2798	0.3661	0.3589	0.3098

mechanism. Evaluation results demonstrate the effectiveness and generalization ability of our approach in learning the affective state and facial movement of speakers.

7. ACKNOWLEDGMENT

This work is supported in part by the National Key Research and Development Program of China under Grant 2017YFB1302003 and by the Fundamental Research Funds for the Central Universities(2018QNA5012).

8. REFERENCES

- [1] Valery A. Petrushin, "How well can people and computers recognize emotions in speech?," *Proc Aaa Fall Symposium*, 1998.
- [2] E. Cosatto, J. Ostermann, H. P. Graf, and J. Schroeter, "Lifelike talking faces for interactive services," *Proceedings of the IEEE*, vol. 91, no. 9, pp. 1406–1429, 2003.
- [3] Alice Wang, Michael Emmi, and Petros Faloutsos, "Assembling an expressive facial animation system," in *ACM SIGGRAPH Symposium on Video Games*, 2007, pp. 21–26.
- [4] Lijuan Wang and Frank K. Soong, "Hmm trajectory-guided sample selection for photo-realistic talking head," *Multimedia Tools & Applications*, vol. 74, no. 22, pp. 1–21, 2015.
- [5] Lei Xie and Zhi Qiang Liu, "Realistic mouth-synching for speech-driven talking face using articulatory modelling," *IEEE Transactions on Multimedia*, vol. 9, no. 3, pp. 500–510, 2007.
- [6] Yao Qian, Yuchen Fan, Wenping Hu, and Frank K. Soong, "On the training aspects of deep neural network (dnn) for parametric tts synthesis," in *IEEE International Conference on Acoustics, Speech and Signal Processing*, 2014, pp. 3829–3833.
- [7] Heiga Ze, Andrew Senior, and Mike Schuster, "Statistical parametric speech synthesis using deep neural networks," in *IEEE International Conference on Acoustics, Speech and Signal Processing*, 2013, pp. 7962–7966.
- [8] Chuang Ding, Lei Xie, and Pengcheng Zhu, "Head motion synthesis from speech using deep neural networks," *Multimedia Tools & Applications*, vol. 74, no. 22, pp. 9871–9888, 2015.
- [9] Bo Fan, Lei Xie, Shan Yang, Lijuan Wang, and Frank K. Soong, "A deep bidirectional lstm approach for video-realistic talking head," *Multimedia Tools & Applications*, vol. 75, no. 9, pp. 5287–5309, 2016.
- [10] Xinjian Zhang, Lijuan Wang, Gang Li, Frank Seide, and Frank K Soong, "A new language independent, photo-realistic talking head driven by voice only.," in *Inter-speech*, 2013, pp. 2743–2747.
- [11] Sarah Taylor, Taehwan Kim, Yisong Yue, Moshe Mahler, James Krahe, Anastasio Garcia Rodriguez, Jessica Hodgins, and Iain Matthews, "A deep learning approach for generalized speech animation," *Acm Transactions on Graphics*, vol. 36, no. 4, pp. 93, 2017.
- [12] Ronald J Williams and David Zipser, "A learning algorithm for continually running fully recurrent neural networks," *Neural computation*, vol. 1, no. 2, pp. 270–280, 1989.
- [13] Alex Graves, "Long short-term memory," *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [14] Volodymyr Mnih, Nicolas Heess, Alex Graves, and Koray Kavukcuoglu, "Recurrent models of visual attention," vol. 3, pp. 2204–2212, 2014.
- [15] Jan K Chorowski, Dzmitry Bahdanau, Dmitriy Serdyuk, Kyunghyun Cho, and Yoshua Bengio, "Attention-based models for speech recognition," in *Advances in neural information processing systems*, 2015, pp. 577–585.
- [16] Attend Show, "Tell: Neural image caption generation with visual attention," *Kelvin Xu et. al.. arXiv Pre-Print*, vol. 83, pp. 89, 2015.
- [17] Alexander M Rush, Sumit Chopra, and Jason Weston, "A neural attention model for abstractive sentence summarization," *arXiv preprint arXiv:1509.00685*, 2015.
- [18] Diederik Kingma and Jimmy Ba, "Adam: A method for stochastic optimization," *Computer Science*, 2014.