

ON USING BACKPROPAGATION FOR SPEECH TEXTURE GENERATION AND VOICE CONVERSION

Jan Chorowski, Ron J. Weiss, Rif A. Saurous, Samy Bengio

Google Brain

{chorowski, ronw, rif, bengio}@google.com

ABSTRACT

Inspired by recent work on neural network image generation which rely on backpropagation towards the network inputs, we present a proof-of-concept system for speech texture synthesis and voice conversion based on two mechanisms: approximate inversion of the representation learned by a speech recognition neural network, and on matching statistics of neuron activations between different source and target utterances. Similar to image texture synthesis and neural style transfer, the system works by optimizing a cost function with respect to the input waveform samples. To this end we use a differentiable mel-filterbank feature extraction pipeline and train a convolutional CTC speech recognition network. Our system is able to extract speaker characteristics from very limited amounts of target speaker data, as little as a few seconds, and can be used to generate realistic speech babble or reconstruct an utterance in a different voice.

Index Terms— Texture synthesis, voice conversion, style transfer, deep neural networks, convolutional networks, CTC

1. INTRODUCTION

Deep neural networks are a family of flexible and powerful machine learning models. Trained discriminatively, they have become the technique of choice in many applications, including image recognition [1], speech recognition [2, 3], and machine translation [4, 5, 6]. Additionally, neural networks can be used to generate new data, having been applied to speech synthesis [7, 8], image generation [9], and image inpainting and superresolution [10].

The representation learned by a discriminatively trained deep neural network can be approximately inverted, turning a classification model into a generator. While exact inversion is impossible, the backpropagation algorithm can be used to find inputs which activate the network in the desired manner. This technique has been applied to the computer vision domain in order to gain insights into network operation [11], find adversarial examples which make imperceptible modifications to image inputs in order to change the network’s predictions [12], synthesize textures [13], and regenerate an image according to the style (essentially matching the low-level texture) of another, referred to as *style transfer* [14].

In this work we investigate the possibility of converting a discriminatively trained CTC speech recognition network into a generator. In particular, we investigate: (i) generating waveforms based solely on the activations of selected network layers, giving insights into the nature of the network’s internal representations, (ii) speech texture synthesis by generating waveforms which result in neuron activations in shallow layers whose statistics are similar to those of real speech, and (iii) voice conversion, the speech analog of image style transfer, where the previous two methods are combined to generate waveforms

which match the high level network activations from a *content* utterance while simultaneously matching low level statistics computed from lower level activations from a *style (identity)* utterance.

2. BACKGROUND

2.1. Texture synthesis based on matching statistics

Julesz [15] proposed that visual texture discrimination is a function of an image’s low level statistical properties. McDermott et al. [16, 17] applied the same idea to sound, showing that perception of sound textures relies on matching certain low level signal statistics. Furthermore, following earlier work on image texture synthesis [18], they demonstrated that simple sound textures, such as rain or fire, can be synthesized using a gradient-based optimization procedure to iteratively update a white noise signal to match the statistics of observed texture signals.

Recently, Gatys et al. [13] proposed a similar statistic matching algorithm to synthesize visual textures. However, instead of manually designing the relevant statistics as a function of the image pixels, they utilized a deep convolutional neural network discriminatively trained on an image classification task. Specifically, they proposed to match uncentered correlations between neuron activations in a selected network layer. Formally, let $C^{(n)} \in \mathbb{R}^{W \times H \times D}$ denote the activations of the n -th convolutional layer, where W is the width of the layer, H is its height, and D is the number of filters. The Gram matrix of uncentered correlations $G^{(n)} \in \mathbb{R}^{D \times D}$ is defined as:

$$G_{i,j}^{(n)} = \frac{1}{WH} \sum_{w=1}^W \sum_{h=1}^H C_{whi}^{(n)} C_{whj}^{(n)}. \quad (1)$$

Gatys et al. demonstrated that realistic visual textures can be synthesized by matching the Gram matrices. In other words, the statistics necessary for texture synthesis are the correlations between the values of two convolutional filters taken over all the pixels in a given convolutional filter map. We note that the Gram features in equation (1) are averaged over all image pixels, and therefore are stationary with respect to the pixel location.

2.2. Style transfer

Approximate network inversions and statistic-matching texture synthesis both generate images by minimizing a loss function with backpropagation towards the inputs. These two approaches can be combined to sample images whose content is similar to a seed image, and whose texture is similar to another one [14]. This approach to style transfer is attractive because it leverages a pretrained neural network which has learned the distribution of natural images, and therefore does not require a large dataset at generation time – a single image of

a given style is all that is required, and it need not be related to the images used to train the network.

3. SPEECH RECOGNITION INPUT RECONSTRUCTION

3.1. Network architecture

To apply the texture generation and stylization techniques to speech we train a fully convolutional speech recognition network following [19] on the Wall Street Journal dataset. The network is trained to predict character sequences in an end-to-end fashion using the CTC [20] criterion. We use parameters typical for a speech recognition network: waveforms sampled at 16kHz are segmented into 25ms windows taken every 10ms. From each window we extract 80 log-mel filterbank features augmented with deltas and delta-deltas. The 13 layer network architecture is derived from [19];

- C0** 128-dimensional 5×5 convolution with 2×2 max-pooling,
- C1** 128-dimensional 5×5 convolution with 1×2 max-pooling,
- C2** 128-dimensional 5×3 convolution,
- C3** 256-dimensional 5×3 convolution with 1×2 max-pooling,
- C4-9** six blocks of 256-dimensional 5×3 convolution,
- FC0-1** two 1024-dimensional fully connected layers,

CTC a fully connected layer and CTC cost over characters,

where filter and pooling window sizes are specified in time \times frequency. All layers use batch normalization, ReLU activations, and dropout regularization. Convolutional layers C0-9 use dropout keep probability 0.75, and fully connected layers use keep probability 0.9.

The network is trained using 10 asynchronous workers with the Adam [21] optimizer using $\beta_1 = 0.9$, $\beta_2 = 0.999$, $\epsilon = 10^{-6}$ and learning rate annealing from 10^{-3} to 10^{-6} . We also use L2 weight decay of 10^{-6} . When decoded using the extended trigram language model from the Kaldi WSJ S5 recipe [22], the model reaches an eval WER of 7.8% on *eval92*. While our network does not reach state-of-the-art accuracy on this dataset, it has reasonable performance and is easily amenable to backpropagation towards inputs. Even though the network was trained on the WSJ¹ corpus, we use the VCTK dataset² for all subsequent experiments.

3.2. Waveform sample reconstruction

Our goal is to generate waveforms that will result in a particular neuron activation pattern when processed by a deep network. Ideally, waveform samples would be optimized directly using the backpropagation algorithm. One possibility is to train networks that operate on raw waveforms as in [23]. However, it is also possible to implement the typical speech feature pipeline in a differentiable way. We follow the second approach, which is facilitated by readily available Tensorflow implementation of signal processing routines [24]:

1. Waveform framing and Hamming window application.
2. DFT computation, which multiplies waveform frames by a complex-valued DFT matrix.
3. Smooth approximate modulus computation, implemented as $\text{abs}(x) \approx \sqrt{\epsilon + \text{re}(x)^2 + \text{im}(x)^2}$, with $\epsilon = 10^{-3}$.
4. Filterbank³ feature computation, which can be implemented as a matrix multiplication.

¹<https://catalog.ldc.upenn.edu/ldc93s6a>, <https://catalog.ldc.upenn.edu/ldc94s13a>

²<http://homepages.inf.ed.ac.uk/jyamagis/page3/page58/page58.html>

³We use a mixed linear and mel scale, where frequencies below 1 kHz are copied from the STFT and higher frequencies are compressed using the mel scale. We employ this scaling below 1 kHz because the mel scale allocates too

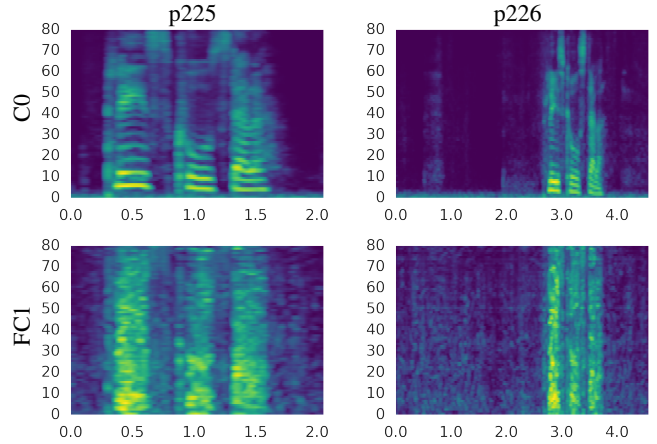


Fig. 1. Mel spectrograms of waveforms reconstructed from layers C0 and FC1 of speakers p225 (female) and p226 (male) from the VCTK dataset. The reconstructions from C0 are nearly exact, while the reconstructions from FC1 are very noisy and are barely intelligible.

5. Taking the elementwise logarithm of the filterbank features.

6. Computing deltas and delta-deltas using convolution over time.

This feature extraction pipeline facilitates two methods for reconstructing waveform samples: (i) gradient-based optimization with backpropagation directly to the waveform, or (ii) gradient-based optimization of the linear spectrogram, followed by Griffin-Lim [25] phase reconstruction. We find that a dual strategy works best, where we first perform spectrogram reconstruction, then invert the spectrogram to yield an initial waveform which is further optimized directly. We use the L-BFGS optimizer [26] for both optimization stages.

3.3. Speech reconstruction from network activations

We implement waveform reconstruction based on network activations following the ReLU non-linearity in a specified layer. Figure 1 shows the spectrograms of waveform reconstructions for speakers p225 and p226 from the VCTK dataset⁴. We have qualitatively established that waveforms reconstructed from shallow network layers are intelligible and the speaker can be clearly identified. Audible phase artifacts are introduced in reconstructions from layer C3 and above, after the final pooling operation over time. While the speech quality degrades, many speaker characteristics are preserved in the reconstructions up to the fully connected layers. Listening to reconstructions from layer C9 it remains possible to recognize the speaker's gender.

In order to reconstruct the waveforms from activations in the fully connected layers FC0 and FC1, we find that the reconstruction cost must be extended with a term penalizing differences between the total energy in each feature frame of the reference and reconstruction. We hypothesize that the network's representation in deeper layers has learned a degree of invariance to the signal magnitude, which hampers reconstruction of realistic signals. For example, the network reliably predicts the CTC blank symbol both for silence and white noise at different amplitudes. The addition of this energy matching penalty enables the network to correctly reconstruct silent segments. However, even with this additional penalty, reconstructions from layers FC0

many bands to low frequencies, some of which are always zero when using 80 mel bands and 256 FFT bins, which was found to be optimal for recognition.

⁴Sound samples are available at https://google.github.io/speech_style_transfer/samples.html

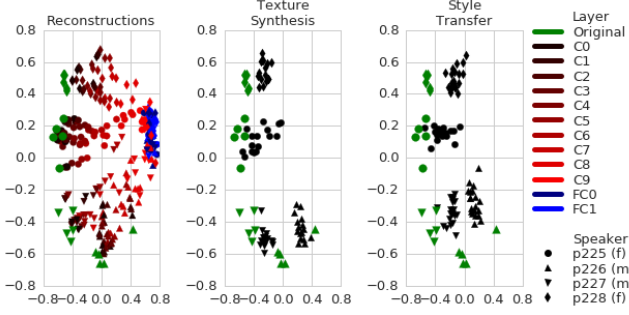


Fig. 2. MDS embeddings of speaker vectors computed on original VCTK recordings, reconstructions from the network, synthesized waveforms, and voice converted waveforms. The synthesized and voice converted utterances are close to the original utterances and reconstructions from early layers. Reconstructions from deep layers converge to a single point, indicating that the speaker identity is lost.

and FC1 are highly distorted. The words are only intelligible with difficulty and the speaker identity is lost.

To evaluate how well reconstructions based on different layers capture characteristics of different speakers, we visualize embedding vectors computed using an internal speaker identification system that uses a Resnet-50 architecture [27] trained on LibriVox⁵ using a triplet-loss [28]. Nearest neighbor classification using these embeddings obtains nearly perfect accuracy on the original VCTK signals. Figure 2 shows a two-dimensional MDS [29] embedding of these vectors. In reconstructions from early layers, signals from each speaker cluster together with no overlap. As the depth increases, the embeddings for all speakers begin to converge on a single point, indicating that the speakers become progressively more difficult to recognize. From this we can conclude that the network’s internal representation becomes progressively more speaker invariant with increasing depth, a desirable property for speaker-independent speech recognition.

3.4. Speech texture synthesis

Unlike image textures whose statistics can be assumed to be stationary across both spatial dimensions, the two dimensions of speech spectrogram features, i.e. time and frequency, have different semantics and should be treated differently. Sound textures are stationary over time but are nonstationary across frequency. This suggests that features extracted from layer activations should involve correlations over time alone. Let $C^{(n)} \in \mathbb{R}^{T \times F \times D}$ be the tensor of activations of the n -th layer of the network which consists of D filters computed for T frames and F frequencies. The temporally stationary Gram tensor, $G^{(n)} \in \mathbb{R}^{F \times F \times D \times D}$, can be written as:

$$G_{ijkl}^{(n)} = \frac{1}{T} \sum_{t=1}^T C_{tik}^{(n)} C_{tjl}^{(n)} \quad (2)$$

We demonstrate that these Gram tensors capture speaker identity by using them as features in a simple nearest neighbor speaker identification system. Figure 3 shows speaker identification accuracy of this system over the first 15 utterances of 30 first speakers of the VCTK dataset. Using the lower network layers (up to C3) yields an accuracy close to 95%, whereas using similar Gram tensors of raw mel-spectrograms extended with deltas and delta-deltas yields only

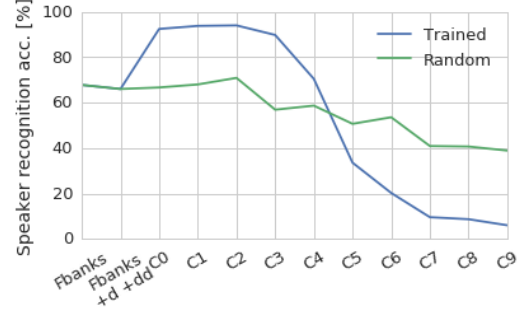


Fig. 3. Accuracy of nearest-neighbor speaker classification using Gram tensors extracted from different network layers.

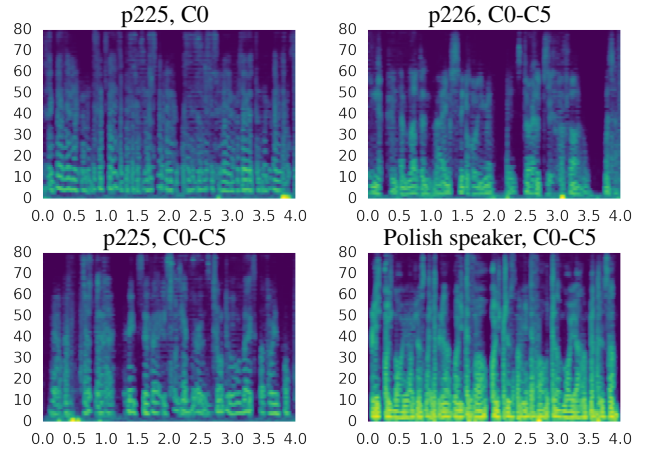


Fig. 4. Mel spectrograms of textures synthesized from Gram matrices computed on 20 utterances from VCTK speakers p225 (female) and p226 (male), as well as a short (1s) utterance in Polish (male). When deeper layers are used, the generated sound captures more temporal structure. Intuitively, listening to “p225, C0” it is hard to discern words, whereas one can hear word boundaries in “p225, C0-C5”. One can also see the characteristic lower pitch in synthesized male voices.

65% accuracy. Deeper layers of the network become progressively less speaker sensitive, mirroring our observations from Figure 2.

We also observe that network training is crucial for Gram features to become speaker-selective and for the texture synthesis to work. After a random initialization the network behaves differently than it does after training: the Gram tensors computed on shallow layers of the untrained network are less sensitive to speaker identity than the corresponding layers in the trained network, while their deeper layers don’t exhibit as dramatic decrease in speaker sensitivity. In contrast, image texture synthesis and style transfer have been reported to work with randomly initialized networks [30].

Figure 4 shows spectrograms of generated speech textures based on speech from the VCTK dataset and a male native Polish speaker. The Gram tensor computed on first layer activations captures the fundamental frequency and harmonics but yields a fairly uniform temporal structure. When features computed on deeper layers are used, longer term phonemic structure can be seen, although the overall speech is not intelligible. This is a consequence of the increased temporal receptive field of filters in deeper layers, where a single activation is a function of structure spanning tens of frames, enabling

⁵<https://librivox.org/>

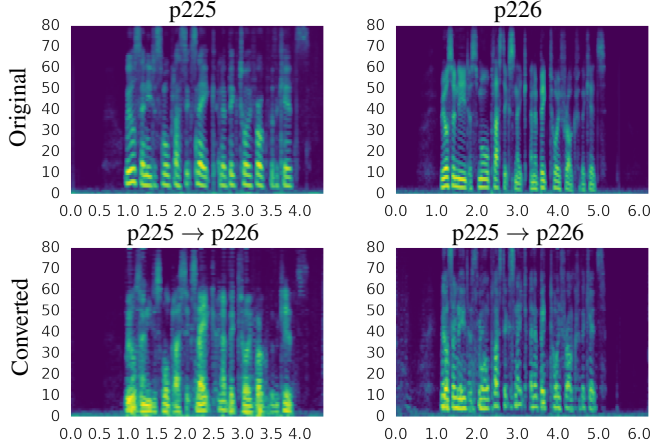


Fig. 5. Mel spectrograms of voice conversion, mapping VCTK utterance 004 between speakers p225 (female) and p226 (male) performed by matching neuron activations to those from a content utterance and Gram features to those computed from 19 speaker identity utterances (about 2 minutes).

the reconstruction of realistic speech babble sounds.

3.5. Voice conversion

The methods described in the previous two sections can be combined to produce the speech analog of image style transfer: a voice conversion system. Specifically we reconstruct the **deep-layer activations of a content utterance**, and the **shallow-layer Gram features of identity or style utterances**.

Listening to the converted samples, we found that **a good tradeoff between matching the target speaker’s voice and sound quality occurs when optimizing a loss that spans all layers**, with the layers C0-C5 matched to style utterances using Gram features, and layers C6-FC1 matched to the content utterance. We normalize the contribution of each layer to the cost by dividing the squared difference between the Gram or activation matrices by their dimensionality. Furthermore, the Gram features of style layers C0-C5 use a weight of 10^5 , activations of content layers C6-C9 use weight 0.2 and activations of layers FC0-FC1 use weight 10, to base the reconstruction on the deepest layers, but provide some signal from those in the middle which are responsible for final voice quality. While the speaker remains identifiable in reconstructions from layers C6-C9 as described in Section 3.3, we find that including these layers in the content loss leads to more natural sounding synthesis. The speaker identity still changes when the Gram feature weight is sufficiently large.

Spectrograms of utterance generated using this procedure are shown in Figure 5. From the spectrograms one can see that the converted utterances contain very different pitch, consistent with the opposite gender. However, because the content loss is applied directly to neuron activations, the exact temporal structure of the content utterance is retained. This highlights a limitation of this approach: the fixed temporal alignment to the content utterance means that it is unable to model temporal variation characteristic to different speakers, such as changes in speaking rate.

4. RELATED WORK

The success of neural image style transfer has prompted a few attempts to apply it to audio. Roberts et al. [31] trained audio clip embeddings using a convolutional network applied directly to raw waveforms and attempted to generate waveforms by maximizing activations of neurons in selected layers. The authors claim noisy results and attribute it to the low quality of the learned filters. Ulyanov et al. [32] used an untrained single-layer network to synthesize simple audio textures such as keyboard and machine gun, and attempted audio style transfer between different musical pieces. The recent work of Wyse [33] is most similar to ours. He examines the application of pretrained convolutional networks for image recognition and for environmental sound classification. **An example of style transfer from human speech to a crowing rooster demonstrates the importance of using a network that has been trained on audio features**, which is in line with our findings. To the best of our knowledge, **our work is the first to demonstrate that style transfer techniques applied to speech recognition networks can be used for voice conversion**.

Speech babble sounds have been previously generated using an unconditioned WaveNet [8] model trained to synthesize speech waveforms. In contrast, we demonstrate that such complex sound textures can be generated from a speech recognition network, using very limited amounts of data from the target speaker.

Typical voice conversion systems rely on advanced speech representations, such as STRAIGHT [34], and use a dedicated conversion function trained on aligned, parallel corpora of different speakers. An overview of the state-of-the-art in this area can be seen in the recent Voice Conversion Challenge [35]. While our system produces samples that have an inferior quality, it operates using a different and novel principle: rather than learning a frame-to-frame conversion, it uses a speech recognition network to define a speaker similarity cost that can be optimized to change the perceived identity of the speaker.

5. LIMITATIONS AND FUTURE WORK

We demonstrate a proof-of-concept speech texture synthesis and voice conversion system that derives a statistical description of the target voice from the activations of a deep convolutional neural network trained to perform speech recognition. The main benefit of the proposed approach is the ability to utilize very limited amounts of data from the target speaker. Leveraging the distribution of natural speech captured by the pretrained network, a few seconds of speech are sufficient to synthesize recognizable characteristics of the target voice. However, the proposed approach is also quite slow, requiring several thousand gradient descent steps. In addition, the synthesized utterances are of relatively low quality.

The proposed approach can be extended in many ways. First, **analogously to the fast image style transfer algorithms [36, 37, 38], the Gram tensor loss can be used as additional supervision for a speech synthesis neural network such as WaveNet [8] or Tacotron [39]. For example, it might be feasible to use the style loss to extend a neural speech synthesis system to a wide set of speakers given only a few seconds of recorded speech from each one**. Second, the method depends on a pretrained speech recognition network. In this work we used a fairly basic network using feature extraction parameters tuned for speech recognition. Synthesis quality could probably be improved by using higher sampling rates, increasing the window overlap and running the network on linear-, rather than mel-filterbank features.

6. ACKNOWLEDGMENTS

Authors thank Yoram Singer, Colin Raffel, Matt Hoffman, Joseph Anagnostini, and Navdeep Jaitly for helpful discussions and inspirations, RJ Skerry-Ryan for signal processing in TF, and Aren Jansen and Sourish Chaudhuri for help with the speaker identification system.

7. REFERENCES

- [1] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” in *Advances in Neural Information Processing Systems*, 2012.
- [2] A. Graves and N. Jaitly, “Towards end-to-end speech recognition with recurrent neural networks,” in *Proc. ICML*, 2014, pp. 1764–1772.
- [3] J. K. Chorowski, D. Bahdanau, D. Serdyuk, K. Cho, and Y. Bengio, “Attention-based models for speech recognition,” in *Advances in Neural Information Processing Systems*, 2015.
- [4] D. Bahdanau, K. Cho, and Y. Bengio, “Neural machine translation by jointly learning to align and translate,” *arXiv:1409.0473*, 2014.
- [5] I. Sutskever, O. Vinyals, and Q. V. Le, “Sequence to sequence learning with neural networks,” in *Advances in Neural Information Processing Systems*, 2014, pp. 3104–3112.
- [6] Y. Wu, M. Schuster, Z. Chen, et al., “Google’s neural machine translation system: Bridging the gap between human and machine translation,” *arXiv:1609.08144*, 2016.
- [7] H. Zen, Y. Agiomyrgiannakis, N. Egberts, F. Henderson, and P. Szczepaniak, “Fast, compact, and high quality lstm-rnn based statistical parametric speech synthesizers for mobile devices,” *arXiv:1606.06061*, 2016.
- [8] A. v. d. Oord, S. Dieleman, H. Zen, et al., “Wavenet: A generative model for raw audio,” *arXiv:1609.03499*, 2016.
- [9] A. Radford, L. Metz, and S. Chintala, “Unsupervised representation learning with deep convolutional generative adversarial networks,” *arXiv:1511.06434*, 2015.
- [10] C. Dong, C. C. Loy, K. He, and X. Tang, “Learning a deep convolutional network for image super-resolution,” in *European Conference on Computer Vision*. Springer, 2014, pp. 184–199.
- [11] K. Simonyan, A. Vedaldi, and A. Zisserman, “Deep inside convolutional networks: Visualising image classification models and saliency maps,” *arXiv:1312.6034*, 2013.
- [12] I. J. Goodfellow, J. Shlens, and C. Szegedy, “Explaining and harnessing adversarial examples,” *arXiv:1412.6572*, 2014.
- [13] L. Gatys, A. S. Ecker, and M. Bethge, “Texture synthesis using convolutional neural networks,” in *Advances in Neural Information Processing Systems*, 2015, pp. 262–270.
- [14] L. A. Gatys, A. S. Ecker, and M. Bethge, “A neural algorithm of artistic style,” *arXiv:1508.06576*, 2015.
- [15] B. Julesz, “Visual pattern discrimination,” *IRE Transactions on Information Theory*, vol. 8, no. 2, pp. 84–92, 1962.
- [16] J. H. McDermott, A. J. Oxenham, and E. P. Simoncelli, “Sound texture synthesis via filter statistics,” in *IEEE Workshop on Applications of Signal Processing to Audio and Acoustics*, 2009.
- [17] J. H. McDermott and E. P. Simoncelli, “Sound texture perception via statistics of the auditory periphery: evidence from sound synthesis,” *Neuron*, vol. 71, no. 5, pp. 926–940, 2011.
- [18] J. Portilla and E. P. Simoncelli, “A parametric texture model based on joint statistics of complex wavelet coefficients,” *International Journal of Computer Vision*, vol. 40, no. 1, 2000.
- [19] Y. Zhang, M. Pezeshki, P. Brakel, et al., “Towards end-to-end speech recognition with deep convolutional neural networks,” *arXiv:1701.02720*, 2017.
- [20] A. Graves, S. Fernández, F. Gomez, and J. Schmidhuber, “Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks,” in *Proc. ICML*. ACM, 2006, pp. 369–376.
- [21] D. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *arXiv:1412.6980*, 2014.
- [22] D. Povey, A. Ghoshal, G. Boulianne, et al., “The kaldi speech recognition toolkit,” in *Proc. ASRU*, 2011.
- [23] T. N. Sainath, R. J. Weiss, A. Senior, K. W. Wilson, and O. Vinyals, “Learning the speech front-end with raw waveform CLDNNs,” in *Proc. Interspeech*, 2015.
- [24] M. Abadi, A. Agarwal, P. Barham, et al., “Tensorflow: Large-scale machine learning on heterogeneous distributed systems,” *arXiv:1603.04467*, 2016.
- [25] D. Griffin and J. Lim, “Signal estimation from modified short-time Fourier transform,” *IEEE Trans. on Acoustics, Speech, and Signal Processing*, vol. 32, no. 2, pp. 236–243, 1984.
- [26] C. Zhu, R. H. Byrd, P. Lu, and J. Nocedal, “Algorithm 778: LBFGS-B: Fortran subroutines for large-scale bound-constrained optimization,” *ACM Transactions on Mathematical Software*, vol. 23, no. 4, pp. 550–560, 1997.
- [27] S. Hershey, S. Chaudhuri, D. P. Ellis, et al., “CNN architectures for large-scale audio classification,” in *Proc. ICASSP*. IEEE, 2017, pp. 131–135.
- [28] H. Bredin, “Tristounet: triplet loss for speaker turn embedding,” in *Proc. ICASSP*, 2017, pp. 5430–5434.
- [29] J. B. Kruskal and M. Wish, *Multidimensional scaling*, vol. 11, Sage, 1978.
- [30] K. He, Y. Wang, and J. Hopcroft, “A powerful generative model using random weights for the deep image representation,” in *Advances in Neural Information Processing Systems*, 2016.
- [31] A. Roberts, C. Resnick, D. Ardila, and D. Eck, “Audio deepdream: Optimizing raw audio with convolutional networks,” in *Proc. ISMIR*, 2016.
- [32] D. Ulyanov and V. Lebedev, “Audio texture synthesis and style transfer,” <https://dmitryulyanov.github.io/audio-texture-synthesis-and-style-transfer/>, 2016.
- [33] L. Wyse, “Audio spectrogram representations for processing with convolutional neural networks,” *arXiv:1706.09559*, 2017.
- [34] H. Kawahara, I. Masuda-Katsuse, and A. De Cheveigne, “Restructuring speech representations using a pitch-adaptive time-frequency smoothing and an instantaneous-frequency-based f0 extraction: Possible role of a repetitive structure in sounds,” *Speech Communication*, vol. 27, no. 3, pp. 187–207, 1999.
- [35] T. Toda, L.-H. Chen, D. Saito, et al., “The voice conversion challenge 2016,” in *Proc. Interspeech*, 2016, pp. 1632–1636.
- [36] D. Ulyanov, V. Lebedev, A. Vedaldi, and V. S. Lempitsky, “Texture networks: Feed-forward synthesis of textures and stylized images,” in *Proc. ICML*, 2016, pp. 1349–1357.

- [37] J. Johnson, A. Alahi, and L. Fei-Fei, “Perceptual losses for real-time style transfer and super-resolution,” in *European Conference on Computer Vision*. Springer, 2016, pp. 694–711.
- [38] V. Dumoulin, J. Shlens, and M. Kudlur, “A learned representation for artistic style,” in *Proc. ICLR*, 2017.
- [39] Y. Wang, R. Skerry-Ryan, D. Stanton, et al., “Tacotron: A fully end-to-end text-to-speech synthesis model,” in *Proc. Interspeech*, 2017.