

UNIVERSITÉ DE LIÈGE

FACULTÉ DES SCIENCES APPLIQUÉES

---

# Automatic Multispeaker Voice Cloning

---

Author:

Corentin JEMINE

Supervisor:

Prof. Gilles LOUPPE

Academic year 2018 - 2019



*Graduation studies conducted for obtaining the Master's degree  
in Data Science by Corentin Jemine*

Possibly, start with a TTS lexicon with some definitions. Meanwhile, I'll make a list of TTS-specific words that may be worth explaining:

coarticulation

linguistic context

spectral envelope

fundamental frequency

contour (a way in which something varies especially the pitch of music or the pattern of tones in an utterance.)

supra-segmental

grapheme

# 1 Abstract

We present our implementation of a multispeaker text-to-speech model that performs voice cloning in real time. This work is largely based on (Jia et al., 2018), an offline framework for multispeaker voice cloning. Our contributions are as follows:

- Building the first (to our knowledge) system able to clone in real-time a high quality voice with good fidelity using only a few seconds of reference speech. **Let's hope I won't have to remove this!**
- Reproducing and open-sourcing the first public implementation of (Jia et al., 2018).
- Adapting this framework to run in real time.
- Creating a software that integrates this model to clone multiple speakers in a conversation.

- propose improvements?
- transfer across languages?

Audio samples are available here [link to samples](#) .

## 2 Introduction

### 2.1 Problem definition

We aim to develop a framework that, given an audio segment of speech, is able to infer a representation of the voice of the speaker and to generate speech from arbitrary text in that same voice. We also want to meet the following constraints:

- The framework should be able to operate in a zero-shot setting, that is, for speakers unseen during training.
- It should operate in real-time:
  - The speaker's voice should be incorporated after only a few seconds of speech.
  - The framework should be able to generate speech in a time shorter or equal to the duration of the produced speech.

Try to answer:

- What is the goal of the application? What are its requirements, what is the setting, what kind of data are we going to use it on?
- What is zero-shot voice cloning (SV2TTS 2.4)? How does it fit in here (difference between an online and offline approach)?
- What are the particularities of our implementation (both model and datasets), what are its upsides and downsides (for example: requires huge datasets but fast inference)?
- What did we ultimately achieve? How good are our results?

**Update this section to reflect what was actually achieved and is presented in this document**

### 2.2 Statistical parametric speech synthesis

**Rewrite all of this and the SOTA to be muuuch shorter** Statistical parametric speech synthesis (SPSS) refers to a group of data-driven TTS synthesis methods that emerged in the late 1990s. In SPSS, the relation between features computed on the input text and output acoustic features is modeled by a statistical generative model (called the acoustic model). A complete SPSS framework thus also includes a pipeline to extract features from the text to synthesize as well as a system able to reconstruct an audio waveform from the acoustic features produced by the acoustic model (such a system is called a vocoder). Unlike the acoustic model, these two parts of the framework may be entirely engineered and make use

of no statistical methods. If it is possible to condition parts of the framework in such a way that the characteristics of the generated voice are modified, then the framework is a multispeaker TTS synthesis system.

The processing of text into features can be nearly inexistent as it can be very extensive. Speech is an intricate process that depends on a wide range of linguistic contexts. Providing these contexts greatly reduces the extent of the task to be learned by the acoustic model, but may require complex natural language processing (NLP) techniques or accuracy trade-offs, especially for rare or unknown words. Linguistic contexts are retrieved on different levels: utterance, phoneme, syllable, word and phrase. For each of those elements, their neighbouring elements of the same level are usually considered, as well as the elements lower in the hierarchy it comprises. For example, a given frame will contain a word, the two previous words, the two following words and the syllables contained in all those words. The position of each element with regard to its parent element can be included (e.g. fifth word in a sentence), as well as grammatical information such as part of speech. For syllables, the lexical stress and accent can be predicted by a statistical model such as a decision tree. For prosody, ToBI (Beckman and Elam, 1997) is often used. Linguistic contexts are often represented and concatenated into a single vector in order to be exploitable by statistical models. Categorical features are encoded using one-hot representations. It is common for the resulting vector to contain hundreds of values.

Talk about evaluation metrics (mainly MOS)? 50 features? Arpabet? Remove last segment on linguistic contexts? → rewrite and summarize methods over the years

## 3 State of the art

### 3.1 Text-to-speech systems

Previous state of the art in SPSS includes hidden Markov models (HMM) based speech synthesis (Tokuda, 2013). The speech generation pipeline is laid out in figure 1. In this framework, the acoustic model is a set of HMMs. The input features are rich linguistic contexts. Ideally, one would train an HMM for each possible context; but as the number of contexts increases exponentially with the number of factors considered, it is not practical to do so. Indeed, not every context will be found in a typical dataset and the training set would then have to be partitioned over the different contexts, which is data inefficient. Instead, contexts are clustered using decision trees and an HMM is learned for each cluster (Yoshimura et al., 1999). Note that this does not solve entirely the training set fragmentation problem. The HMMs are trained to produce a distribution over mel-frequency cepstral coefficients (MFCC) with energy (called static features), their delta and delta-delta coefficients (called dynamic features) as well as a binary flag that indicates which parts of the audio should contain voice. This is shown in figure 2. A new sequence of static features is retrieved from these static and dynamic features using the maximum likelihood parameter generation (MLPG) algorithm (Tokuda et al., 2000). These static features are then fed through the MLSA vocoder (Imai, 1983). It is possible to modify the voice generated by conditioning on a speaker or tuning the generated speech parameters with adaptation or interpolation techniques (Yoshimura et al., 1997) *elaborate a bit on these techniques?*, making HMM-based speech synthesis a multispeaker TTS system. *Compare with concatenative see (Zen et al., 2013) and [ieeexplore.ieee.org/document/541110](http://ieeexplore.ieee.org/document/541110).*

Improvements to this framework were later brought by feed-forward deep neural networks (DNN), as a result of progress in both hardware and software. (Zen et al., 2013) proposes to replace entirely the decision tree-clustered HMMs in favor of a DNN. They argue for better data efficiency as the training set is no longer fragmented in different clusters of contexts, *and for a more powerful model?*. They demonstrate improvements over the speech quality with a number of parameters similar to that of the HMM-based approach. Their best model is a DNN with 4 layers of 256 units using a sigmoid activation function. Subjects assessing the quality of the generated audio samples report that the DNN-based models produces speech that sounds less muffled than that of the HMM-based models. Later researches corroborate these findings (Qian et al., 2014). (Hashimoto et al., 2015) additionally studies the effect of replacing MLPG with another DNN. The combinations of HMM/DNN and MLPG/DNN give rise to four possible frameworks, the novel ones being HMM+DNN and DNN+DNN<sup>1</sup>, while HMM+MLPG and DNN+MLPG are the frameworks described respectively in (Tokuda, 2013) and (Zen et al., 2013). Each

<sup>1</sup>Note that since the two networks are consecutive in the framework, they can be considered a single network.

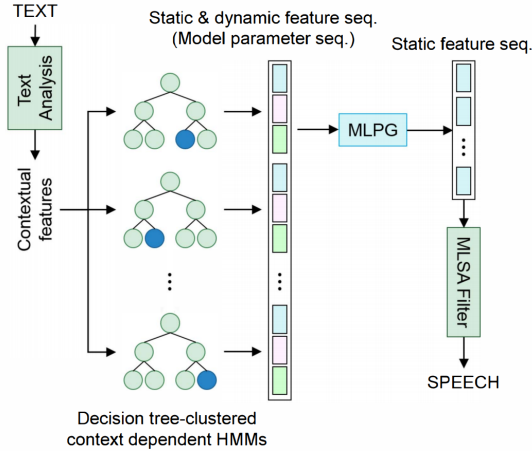


Figure 1: The general HMM-based TTS synthesis approach.

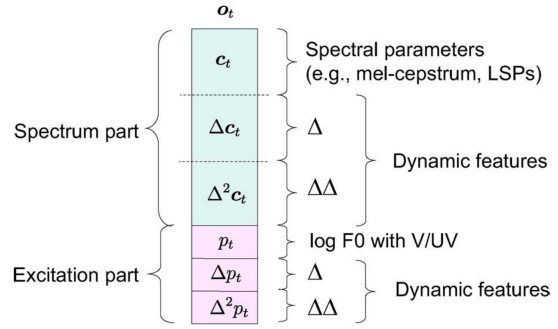


Figure 2: Dynamic and static features produced by the HMMs. F0 is the fundamental frequency and V/UV is the voicing flag.

DNN they use is 3 layers deep with 1024 units using a sigmoid activation function. MOS results confirm that DNN+MLPG is significantly better than HMM+MLPG. The DNN+DNN approach performs as well as HMM+MLPG while HMM+DNN is worse. In another experiment, they introduce a DNN before and after MLPG. While both approaches yield a MOS similar to DNN+MLPG, neither are statistically significantly better. **make this part nicer to read or maybe just remove it? Quid of the MOS? Mention DMDN?**

(Fan et al., 2014) supports that RNNs make natural acoustic models as they are able to learn a compact representation of complex and long-span functions, better so than shallow architectures such as the decision trees used in HMM-based TTS. Furthermore, the internal state of RNNs makes the mapping from an input frame to output acoustic features no longer deterministic, allowing for more variety in the synthesized audio. As RNNs are fit to generate temporally consistent series, the static features can directly be determined by the acoustic model, alleviating the need for dynamic features and MLPG. The authors present two RNNs: Hybrid\_A with three feed-forward layers followed by one bidirectional (BDLSTM) layer and Hybrid\_B with two feed-forward layers followed by two BDLSTM layers. They argue that deeper structures of BDLSTM would worsen the performance due to imprecise gradient computation. They compare these networks against the HMM and DNN based approaches described previously, using objective and subjective measures. The objective measures are compared with the ground truth: log spectral distance (LSD), voiced/unvoiced (V/UV) error rate and fundamental frequency (F0) distortion in root mean squared error (RMSE). The subjective measure is a preference test where participants must choose between two audio samples of different models and have the option to select neither. Results are shown in figures 3 and 4. DNN\_A is 6 layers deep with 512 units per layer while DNN\_B is 3 layers deep with 1024 units per layer. These two networks perform very similarly. Hybrid\_B systematically performs better than the other approaches.

Later, a substantial breakthrough in TTS is achieved with the coming of WaveNet (van den Oord et al., 2016). WaveNet is a deep convolutional neural network that, for a raw audio waveform, models the distribution of a single sample given all previous ones. It is thus possible to generate audio by predicting samples one at a time in an autoregressive fashion. WaveNet leverages stacks of one-dimensional dilated convolutions with a dilation factor increasing exponentially with the layer depth, which allows for a very large receptive field. The output is a categorical distribution (using a softmax layer) over sample values. For tractability reasons, the output signal is restricted to 256 values that correspond to a  $\mu$ -law quantization, which is invertible. The MOS resulting from one of the experiments suggests that there is no statistically significant degradation in the naturalness of speech quantized in this manner. The architecture of WaveNet comes with several non-trivial components, described in later sections of this document [link to the section](#). On its own, a trained WaveNet generates sound alike the training data but without semantics. The network must be locally conditioned on linguistic contexts to achieve TTS synthesis. Furthermore, it allows for conditioning with respect to a vector constant at every timestep (global conditioning), which can be used to designate the speaker identity. The authors suggest that

Model \ Measures	LSD (dB)	V/U Error rate	F0 RMSE (Hz)
HMM (2.89M)	3.74	5.8%	17.7
DNN_A (1.55M)	3.73	5.8%	15.8
DNN_B (2.59M)	3.73	5.9%	15.9
Hybrid_A (2.30M)	3.61	5.7%	16.4
Hybrid_B (3.61M)	3.54	5.6%	15.8

Figure 3: Performance of the different frameworks evaluated on objective measures. In parentheses are the number of parameters of each acoustic model.

44% Hybrid_B	29% Neutral	27% Hybrid_A
59% Hybrid_B	19% Neutral	22% HMM
55% Hybrid_B	25% Neutral	20% DNN_B

Figure 4: Two by two comparisons of some of the frameworks in terms of preferences.

WaveNet is able to encode the embedding of several speakers seen in training with a shared internal representation. Speaker identities are given as a one-hot encoding. The audio generated by WaveNet requires no post-processing other than the inversion of the  $\mu$ -law. TTS with WaveNet does not exactly fit as an SPSS system considering that it does not produce clear intermediate acoustic features and instead serves as both a statistical model and a vocoder in a single pipeline. Performance scores are reported in Figure 5. The parametric approach is an LSTM-based system while the other is an HMM-driven unit selection concatenative system (not detailed in this document). Notice how the results vary between US English and Mandarin Chinese, showing that TTS performance is not language agnostic. **recheck all the info on WaveNet later + F0??**

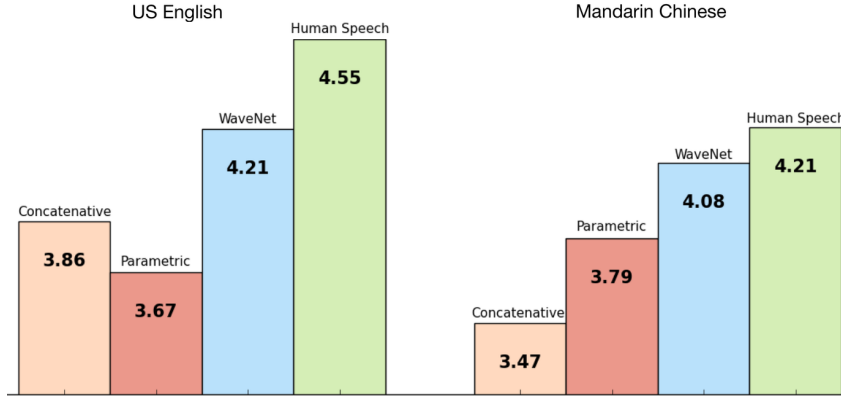


Figure 5: MOS of WaveNet’s performance compared with a parametric and concatenative approach as well as with natural speech.

Deep Voice (Arik et al., 2017) proposes a fully neural TTS framework that exploits WaveNet among other deep architectures. Deep Voice stands out by making use of only a few intermediate features: phonemes with stress annotations, phoneme durations, and F0. This has the advantage of making the framework easily transferable to new domains with little engineering effort as complex linguistic features need not be derived. The neural networks intervening in Deep Voice are: a grapheme-to-phoneme model that converts text to phonemes, a segmentation model that aligns a sequence of phonemes to an audio segment, a phoneme duration model that predicts the duration of each phoneme at generation time, a fundamental frequency model that predicts the V/UV flag with F0 values for voiced parts and finally, WaveNet which acts as an audio synthesis model. In all the works we presented previously, only the audio synthesis model was not manually engineered. While others have researched the use of neural networks for some of these components, Deep Voice is the first to simultaneously employ them all in a single framework. However, the grapheme-to-phoneme model is only used as a fallback for words not present in a phoneme dictionary. The roles and interactions of these components at training and inference time are shown in figure 6. Deep Voice also improves on the inference time of WaveNet, with a speedup of

up to 400. This allows for real-time or near real-time execution, with a tunable speed/quality trade-off. Deep Voice does not yield state of the art results but instead serves as groundwork for future researches. [discuss results](#) .

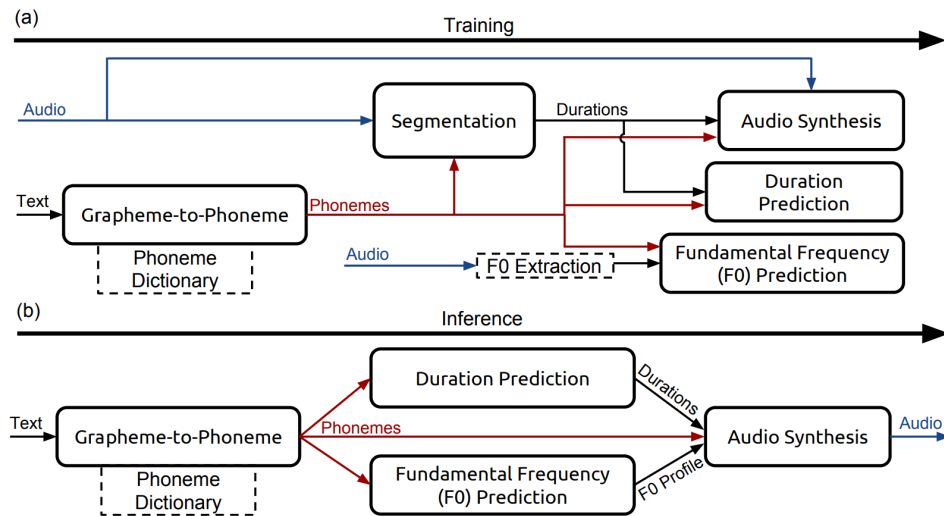


Figure 6: The training (a) and inference (b) procedure of Deep Voice. On the left of the image are the inputs, on the right are the outputs. Note that the segmentation model is only used for training.

[Skipping the Deep Voice 2 & 3 papers. I'll see later whether or not I should include them](#)

Published approximately at the same time, Tacotron (Wang et al., 2017) is a sequence-to-sequence model that produces a spectrogram from a sequence of characters alone, further reducing the need for domain expertise. An audio waveform can be estimated from the spectrogram using the Griffin-Lin algorithm. Tacotron is also built fully with neural networks and is trained in an end-to-end fashion. It uses an encoder-decoder architecture where, at each step, the decoder operates on a weighted sum of the encoder outputs. This mechanism described in (Bahdanau et al., 2014) lets the network decide which steps of the input sequence are important with respect to each steps of the output sequence. Tacotron achieves a MOS of 3.85 on a US English dataset, which is more than the 3.69 score obtained in the parametric approach of (Zen et al., 2016) but less than the 4.09 score obtained by the concatenative approach of (Gonzalvo et al., 2016). The authors mention that Tacotron is merely a step towards a better framework. By the end of 2017, Tacotron 2 is published (Shen et al., 2017). The architecture of Tacotron 2 remains that of an encoder-decoder with attention although several changes to the types of layers are made. The main difference with Tacotron is the addition of a modified WaveNet for vocoder. On the same dataset, Tacotron 2 achieves a MOS of 4.53 compared to 4.58 for human speech (the difference is not statistically significant), achieving the all-time highest MOS. In a preference study, Tacotron 2 was found to be only slightly less preferred on average than ground truth samples. The ratings from that study are shown in figure

[Few samples](#)

[SV2TTS](#)

[Extensions?](#)

[Link audio samples for everything where available](#)

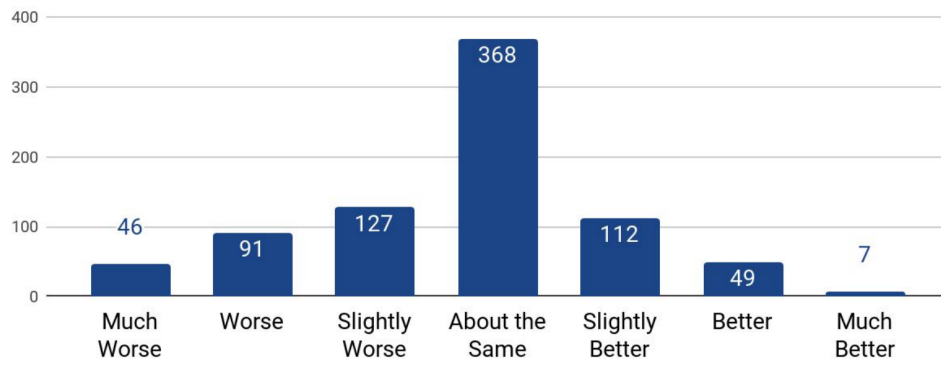


Figure 7: Preference ratings between Tacotron 2 and ground truth samples. 800 ratings on 100 items. The labels are expressed with respect to Tacotron 2.



## References

- Ye Jia, Yu Zhang, Ron J. Weiss, Quan Wang, Jonathan Shen, Fei Ren, Zhifeng Chen, Patrick Nguyen, Ruoming Pang, Ignacio Lopez-Moreno, and Yonghui Wu. Transfer learning from speaker verification to multispeaker text-to-speech synthesis. CoRR, abs/1806.04558, 2018. URL <http://arxiv.org/abs/1806.04558>.
- Mary E. Beckman and Gayle Ayers Elam. Guidelines for tobi labelling, 03 1997.
- Yoshihiko; Toda Tomoki; Zen Heiga; Yamagishi Junichi; Oura Keiichiro Tokuda, Keiichi; Nankaku. Speech synthesis based on hidden markov models. Proceedings of the IEEE, 101, 05 2013. doi: 10.1109/JPROC.2013.2251852.
- Takayoshi Yoshimura, Keiichi Tokuda, Takashi Masuko, Takao Kobayashi, and Tadashi Kitamura. Simultaneous modeling of spectrum, pitch and duration in hmm-based speech synthesis. In EUROSPEECH, 1999.
- K. Tokuda, T. Yoshimura, T. Masuko, T. Kobayashi, and T. Kitamura. Speech parameter generation algorithms for hmm-based speech synthesis. In 2000 IEEE International Conference on Acoustics, Speech, and Signal Processing. Proceedings (Cat. No.00CH37100), volume 3, pages 1315–1318 vol.3, June 2000. doi: 10.1109/ICASSP.2000.861820.
- S. Imai. Cepstral analysis synthesis on the mel frequency scale. In ICASSP ’83. IEEE International Conference on Acoustics, Speech, and Signal Processing, volume 8, pages 93–96, April 1983. doi: 10.1109/ICASSP.1983.1172250.
- Takayoshi Yoshimura, Takashi Masuko, Keiichi Tokuda, Takao Kobayashi, and Tadashi Kitamura. Speaker interpolation in hmm-based speech synthesis system. In EUROSPEECH, 1997.
- H. Zen, A. Senior, and M. Schuster. Statistical parametric speech synthesis using deep neural networks. In 2013 IEEE International Conference on Acoustics, Speech and Signal Processing, pages 7962–7966, May 2013. doi: 10.1109/ICASSP.2013.6639215.
- Y. Qian, Y. Fan, W. Hu, and F. K. Soong. On the training aspects of deep neural network (dnn) for parametric tts synthesis. In 2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pages 3829–3833, May 2014. doi: 10.1109/ICASSP.2014.6854318.
- K. Hashimoto, K. Oura, Y. Nankaku, and K. Tokuda. The effect of neural networks in statistical parametric speech synthesis. In 2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pages 4455–4459, April 2015. doi: 10.1109/ICASSP.2015.7178813.
- Y Fan, Yuqian Qian, Feng-Long Xie, and Frank Soong. Tts synthesis with bidirectional lstm based recurrent neural networks. pages 1964–1968, 01 2014.
- Aäron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves, Nal Kalchbrenner, Andrew W. Senior, and Koray Kavukcuoglu. Wavenet: A generative model for raw audio. CoRR, abs/1609.03499, 2016. URL <http://arxiv.org/abs/1609.03499>.
- Sercan Ömer Arik, Mike Chrzanowski, Adam Coates, Greg Diamos, Andrew Gibiansky, Yongguo Kang, Xian Li, John Miller, Jonathan Raiman, Shubho Sengupta, and Mohammad Shoeybi. Deep voice: Real-time neural text-to-speech. CoRR, abs/1702.07825, 2017. URL <http://arxiv.org/abs/1702.07825>.
- Yuxuan Wang, R. J. Skerry-Ryan, Daisy Stanton, Yonghui Wu, Ron J. Weiss, Navdeep Jaitly, Zongheng Yang, Ying Xiao, Zhifeng Chen, Samy Bengio, Quoc V. Le, Yannis Agiomyrgiannakis, Rob Clark, and Rif A. Saurous. Tacotron: A fully end-to-end text-to-speech synthesis model. CoRR, abs/1703.10135, 2017. URL <http://arxiv.org/abs/1703.10135>.
- Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. CoRR, abs/1409.0473, 2014. URL <http://arxiv.org/abs/1409.0473>.

- Heiga Zen, Yannis Agiomyrgiannakis, Niels Egberts, Fergus Henderson, and Przemyslaw Szczepaniak. Fast, compact, and high quality LSTM-RNN based statistical parametric speech synthesizers for mobile devices. CoRR, abs/1606.06061, 2016. URL <http://arxiv.org/abs/1606.06061>.
- Xavi Gonzalvo, Siamak Tazari, Chun-an Chan, Markus Becker, Alexander Gutkin, and Hanna Silen. Recent advances in google real-time hmm-driven unit selection synthesizer. In Interspeech, pages 2238–2242, 2016.
- Jonathan Shen, Ruoming Pang, Ron J. Weiss, Mike Schuster, Navdeep Jaitly, Zongheng Yang, Zhifeng Chen, Yu Zhang, Yuxuan Wang, R. J. Skerry-Ryan, Rif A. Saurous, Yannis Agiomyrgiannakis, and Yonghui Wu. Natural TTS synthesis by conditioning wavenet on mel spectrogram predictions. CoRR, abs/1712.05884, 2017. URL <http://arxiv.org/abs/1712.05884>.