

A 2019 Guide to Speech Synthesis with Deep Learning



Derrick Mwiti

Aug 28 · 13 min read

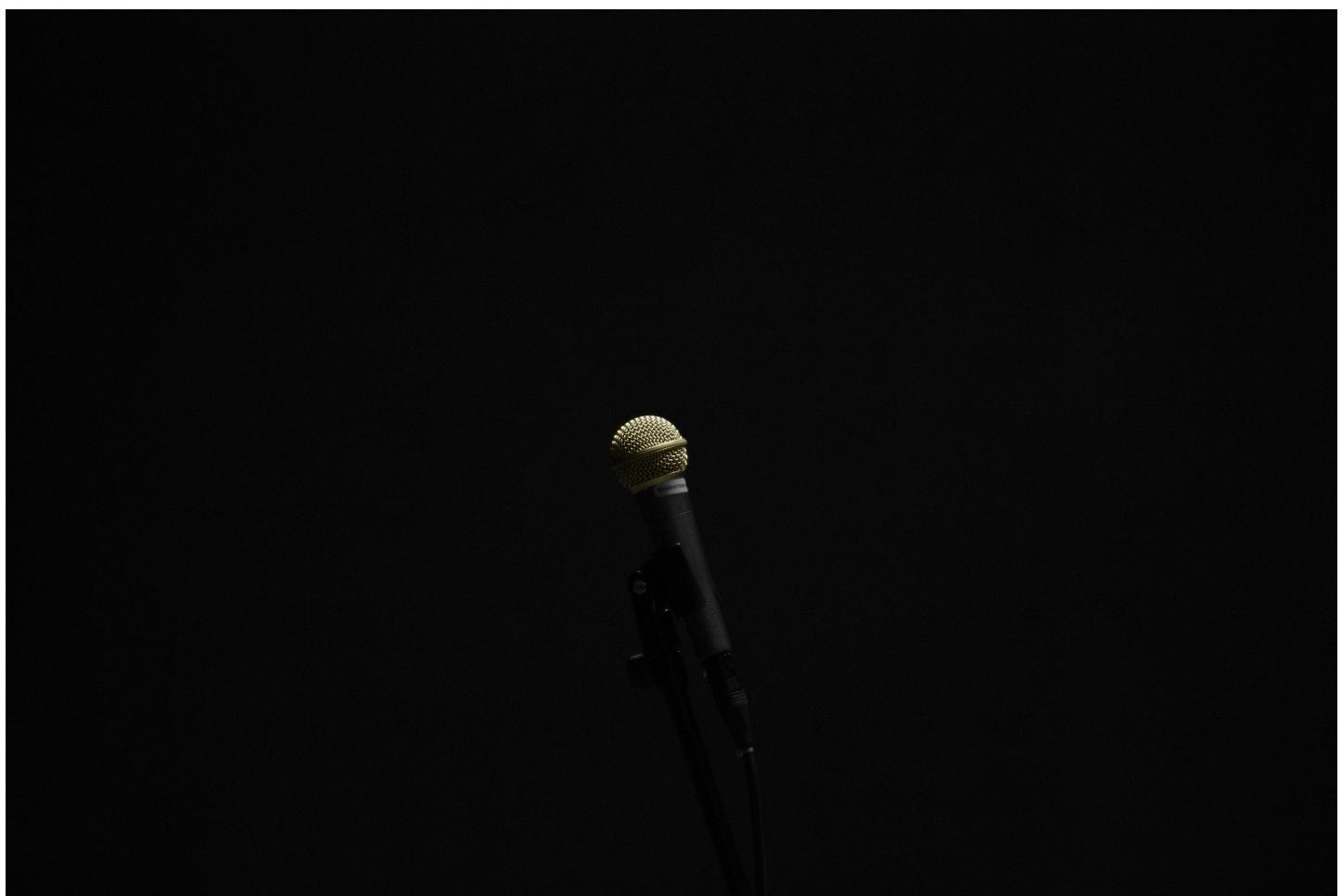


Photo by Daniel Sandvik on Unsplash

Artificial production of human speech is known as speech synthesis. This machine learning-based technique is applicable in text-to-speech, music generation, speech generation, speech-enabled devices, navigation systems, and accessibility for visually-impaired people.

In this article, we'll look at research and model architectures that have been written and developed to do just that using deep learning.

But before we jump in, there are a couple of specific, traditional strategies for speech synthesis that we need to briefly outline: **concatenative** and **parametric**.

In the concatenative approach, speeches from a large database are used to generate new, audible speech. In a case where a different style of speech is needed, a new database of audio voices is used. This limits the scalability of this approach.

The parametric approach uses a recorded human voice and a function with a set of parameters that can be modified to change the voice.

These two approaches represent the old way of doing speech synthesis. Now let's look at the new ways of doing it using deep learning. Here's the research we'll cover in order to examine popular and current approaches to speech synthesis:

- WaveNet: A Generative Model for Raw Audio
 - Tacotron: Towards End-toEnd Speech Synthesis
 - Deep Voice 1: Real-time Neural Text-to-Speech
 - Deep Voice 2: Multi-Speaker Neural Text-to-Speech
 - Deep Voice 3: Scaling Text-to-speech With Convolutional Sequence Learning
 - Parallel WaveNet: Fast High-Fidelity Speech Synthesis
 - Neural Voice Cloning with a Few Samples
 - VoiceLoop: Voice Fitting and Synthesis via A Phonological Loop
 - Natural TTS Synthesis by Conditioning WaveNet on Mel Spectrogram Predictions
- . . .

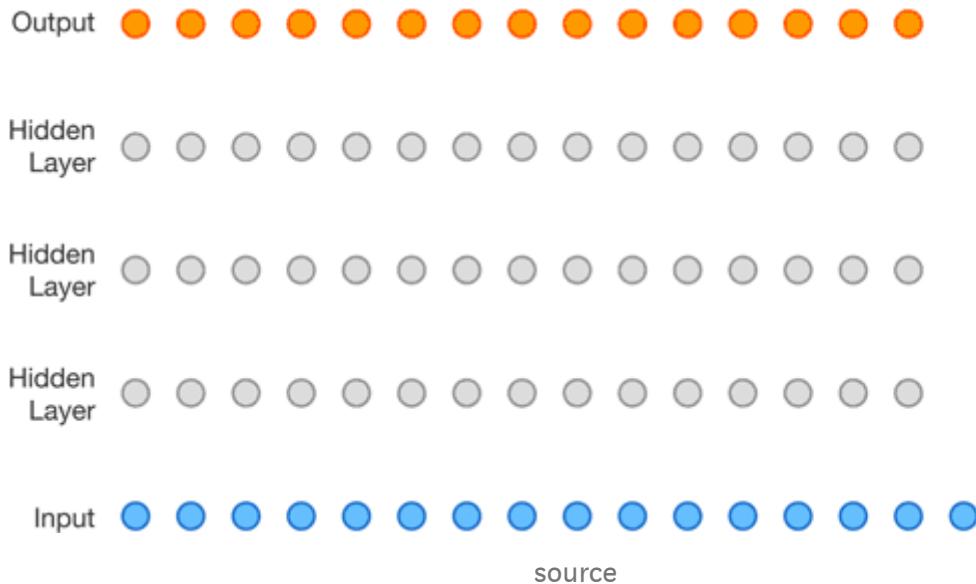
WaveNet: A Generative Model for Raw Audio

The authors of this paper are from Google. They present a neural network for generating raw audio waves. Their model is fully probabilistic and autoregressive, and it generates state-of-the-art text-to-speech results for both English and Mandarin.

WaveNet: A Generative Model for Raw Audio

This paper introduces WaveNet, a deep neural network for generating raw audio waveforms. The model is fully...

arxiv.org



WaveNet is an audio generative model based on the PixelCNN. It's capable of producing audio that's very similar to a human voice.

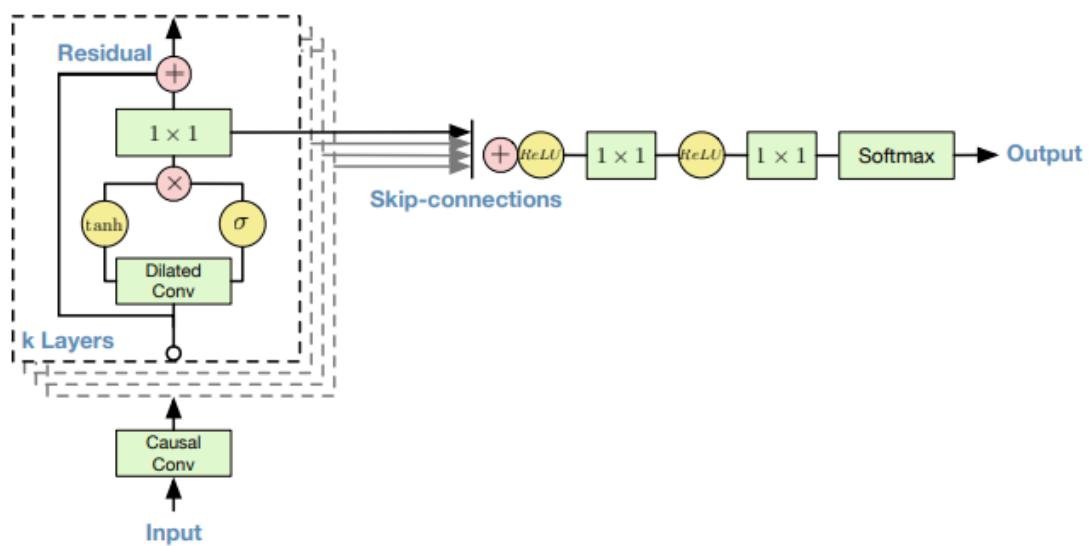


Figure 4: Overview of the residual block and the entire architecture.

source



In this generative model, each audio sample is conditioned on the previous audio sample. The conditional probability is modeled by a stack of convolutional layers. This network doesn't have pooling layers, and the output of the model has the same time dimensionality as the input.

[source](#)

The use of causal convolutions in the architecture ensures that the model doesn't violate the ordering of how the data is modeled. In this model, each predicted voice sample is fed back to the network to aid in predicting the next one. Since causal convolutions don't have a recurrent connection, they're faster to train than RNNs.

One of the major challenges of using causal convolutions is that they require many layers in order to increase the receptive field. To solve this challenge, the authors use

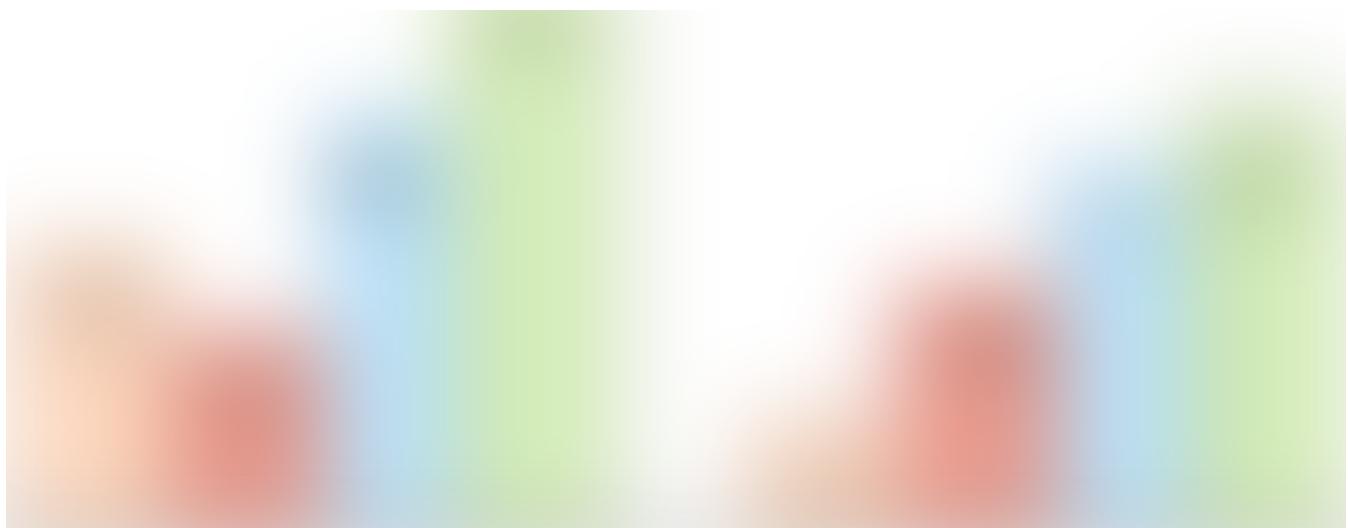
dilated convolutions. Dilated convolutions enable networks to have a large receptive field but with a few layers. Modeling the conditional distributions over the individual audio samples is done using a softmax distribution.

[source](#)

The model is evaluated on multispeaker speech generation, text-to-speech, and music audio modeling. The MOS (Mean Opinion Score) is used for this testing. It measures the quality of voice. It's basically the opinion of a person about the voice quality. It is a number between one and five, with five being the best quality.

[source](#)

The figure below shows the quality of waveNets on a scale of 1–5.



source

· · ·

Machine learning doesn't have to live on servers or in the cloud — it can also live on your smartphone. And Fritz AI has the tools to easily teach mobile apps to see, hear, sense, and think.

· · ·

Tacotron: Towards End-toEnd Speech Synthesis

The authors of this paper are from Google. Tacotron is an end-to-end generative text-to-speech model that synthesizes speech directly from text and audio pairs. Tacotron achieves a 3.82 mean opinion score on US English. Tacotron generates speech at frame-level and is, therefore, faster than sample-level autoregressive methods.

Tacotron: Towards End-to-End Speech Synthesis

A text-to-speech synthesis system typically consists of multiple stages, such as a text analysis frontend, an acoustic...

arxiv.org

The model is trained on audio and text pairs, which makes it very adaptable to new datasets. Tacotron has a seq2seq model that includes an encoder, an attention-based decoder, and a post-processing net. As seen in the architecture diagram below, the model takes characters as input and outputs a raw spectrogram. This spectrogram is then converted to waveforms.

source

The figure below shows what the CBHG module looks like. It consists of 1-D convolution filters, highway networks, and a bidirectional GRU (Gated Recurrent Unit).

[source](#)

A character sequence is fed to the encoder, which extracts sequential representations of text. Each character is represented as a one-hot vector and embedded into a continuous vector. Non-linear transformations are then added, followed by a dropout layer to reduce overfitting. This, in essence, reduces the mispronunciation of words.

The decode used is a *tanh* content-based attention decoder. The waveforms are then generated using the Griffin-Lim algorithm. The hyper-parameters used for this model are shown below.

[source](#)

The figure below shows the performance of Tacotron compared to other alternatives.

[source](#)

• • •

Deep Voice 1: Real-time Neural Text-to-Speech

The authors of this paper are from Baidu's Silicon Valley Artificial Intelligence Lab. Deep Voice is a text-to-speech system developed using deep neural networks.

Deep Voice: Real-time Neural Text-to-Speech

We present Deep Voice, a production-quality text-to-speech system constructed entirely from deep neural networks. Deep...

[arxiv.org](#)

It has five major building blocks:

- A segmentation model for locating phoneme boundaries with deep neural networks using connectionist temporal classification (CTC) loss.
- A grapheme-to-phoneme conversion model (grapheme-to-phoneme is the process of using rules to generate a word's pronunciation).
- A phoneme duration prediction model.
- A fundamental frequency prediction model.
- An audio synthesis model using a variant of WaveNet that uses fewer parameters.

source

The grapheme-to-phoneme model converts English characters to phonemes. The segmentation model identifies where each phoneme begins and ends in an audio file. The phoneme duration model predicts the duration of every phoneme in a phoneme sequence.

The fundamental frequency model predicts whether a phoneme is voiced. The audio synthesis model synthesizes audio by combining the output of the grapheme-to-phoneme, phoneme duration, and fundamental frequency prediction models.

Here's how this model fares compared to other models.

• • •

Deep Voice 2: Multi-Speaker Neural Text-to-Speech

This paper represents the second iteration of Deep Voice by Baidu Silicon Valley Artificial Intelligence Lab. They introduce a method for augmenting neural text-to-speech with low dimensional trainable speaker embeddings to produce various voices from a single model.

The model is based on a similar pipeline as DeepVoice 1. However, it represents a significant improvement in audio quality. The model is able to learn hundreds of unique voices from less than half an hour of data per speaker.

Deep Voice 2: Multi-Speaker Neural Text-to-Speech

We introduce a technique for augmenting neural text-to-speech (TTS) with low dimensional trainable speaker embeddings to...

[arxiv.org](#)

The authors also introduce a WaveNet-based spectrogram-to-audio neural vocoder, which is then used with Tacotron in place of Griffin-Lim audio generation. The main focus of this paper is to handle multiple speakers with fewer data from each speaker. The general architecture is similar to Deep Voice 1. The training process of Deep Voice 2 is depicted in the figure below.

[source](#)

The major difference between Deep Voice 2 and Deep Voice 1 is the separation of the phoneme duration and frequency models. Deep Voice 1 has a single model for jointly predicting the phoneme duration and frequency profile; in Deep Voice 2, the phoneme durations are predicted first and then they are used as inputs to the frequency model.

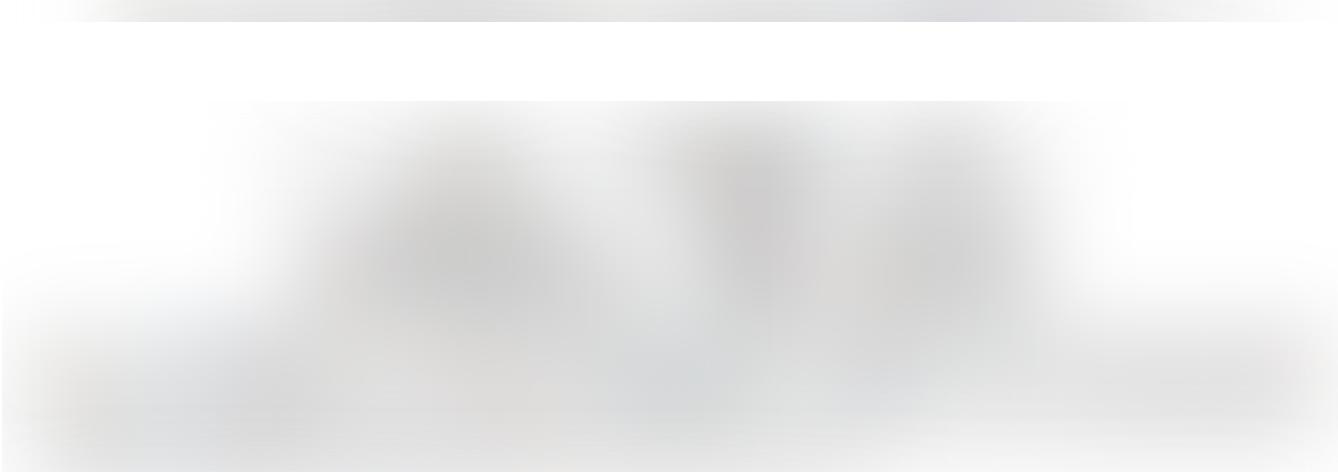
The segmentation model in Deep Voice 2 is a convolutional-recurrent architecture with connectionist temporal classification (CTC) loss applied to classify phoneme pairs. The major modification in Deep Voice 2 is the addition of batch normalization and residual connections in the convolutional layers. Its vocal model is based on a WaveNet architecture.

Synthesizing speech from multiple speakers is done by augmenting each model with a single low-dimensional level speaker embedding vector per speaker. Weight sharing between speakers is achieved by storing speaker-dependent parameters in a very low-dimensional vector.

The initial states of the recurrent neural network (RNN) are produced using speaker embeddings. A uniform distribution is used to randomly initialize the speaker embeddings and trained jointly using backpropagation. Speaker embeddings are incorporated in multiple portions of the model in order to ensure that each speaker's unique voice signature is factored in.

source

Let's now see how this model performs in comparison to other models.



source

· · ·

The latest in deep learning — from a source you can trust. Sign up for a weekly dive into all things deep learning, curated by experts working in the field.

· · ·

Deep Voice 3: Scaling Text-to-speech With Convolutional Sequence Learning

In the third iteration of Deep Voice, the authors introduce is a fully-convolutional attention-based neural text-to-speech (TTS) system.

Deep Voice 3: Scaling Text-to-Speech with Convolutional Sequence Learning

We present Deep Voice 3, a fully-convolutional attention-based neural text-to-speech (TTS) system. Deep Voice 3 matches...

The authors propose a fully-convolutional character-to-spectrogram architecture that enables fully parallel computation. The architecture is an attention-based sequence-to-sequence model. The model was trained on the LibriSpeech ASR dataset.

The proposed architecture is able to convert textual features such as characters, phonemes, and stresses into different vocoder parameters. Some of these include mel-band spectrograms, linear-scale log magnitude spectrograms, fundamental frequency, spectral envelope, and aperiodicity parameters. These vocoder parameters are then used as the input for the audio waveform synthesis model.



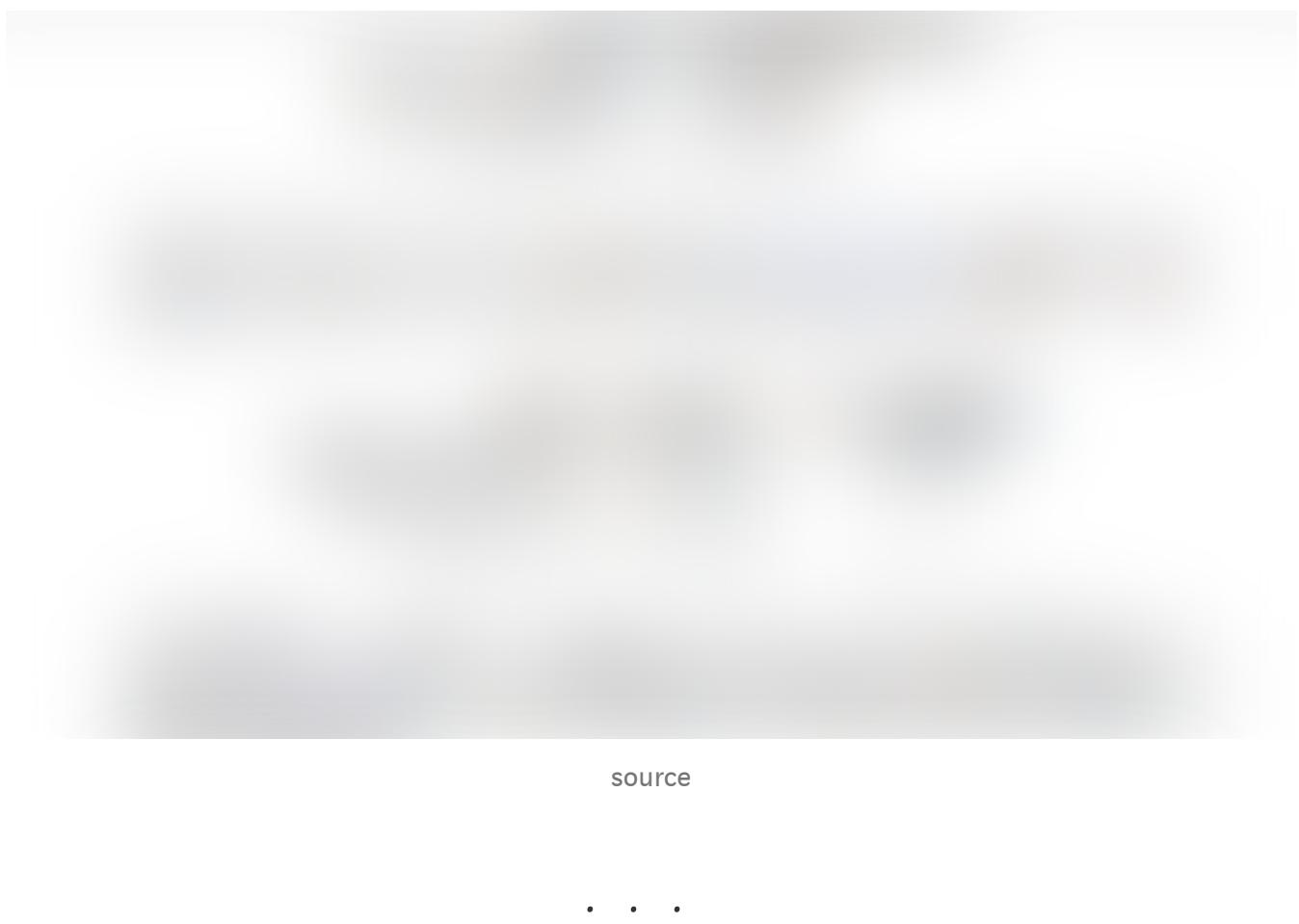
source

The architecture is composed of the following:

- Encoder — a fully-convolutional encoder that converts textual features to an internal learned representation.
- Decoder — a fully-convolutional causal decoder that decodes the learned representations in an autoregressive manner.
- Converter — a fully-convolutional post-processing network that predicts the final vocoder parameters.

For text pre-processing, the authors' uppercase text input characters, remove punctuation marks, end each utterance with a period or question mark, and replace spaces with a special character that indicates the length of a pause.

The figure below is a comparison of the performance of this model with other alternative models.



source

...

Parallel WaveNet: Fast High-Fidelity Speech Synthesis

The authors of this paper are from Google. They introduce a method known as *Probability Density Distillation*, which trains a parallel feed-forward network from a trained WaveNet. The method has been built by marrying the best features of Inverse autoregressive flows (IAFs) and WaveNet. These features represent the efficient training of WaveNet and the efficient sampling of IAF networks.

Parallel WaveNet: Fast High-Fidelity Speech Synthesis

The recently-developed WaveNet architecture is the current state of the art in realistic speech synthesis, consistently...

arxiv.org

For training, the authors use a trained WaveNet as a ‘teacher’, and the parallel WaveNet ‘student’ learns from this. The goal here is to have the student match the probability of its own samples under the distribution learned from the teacher.

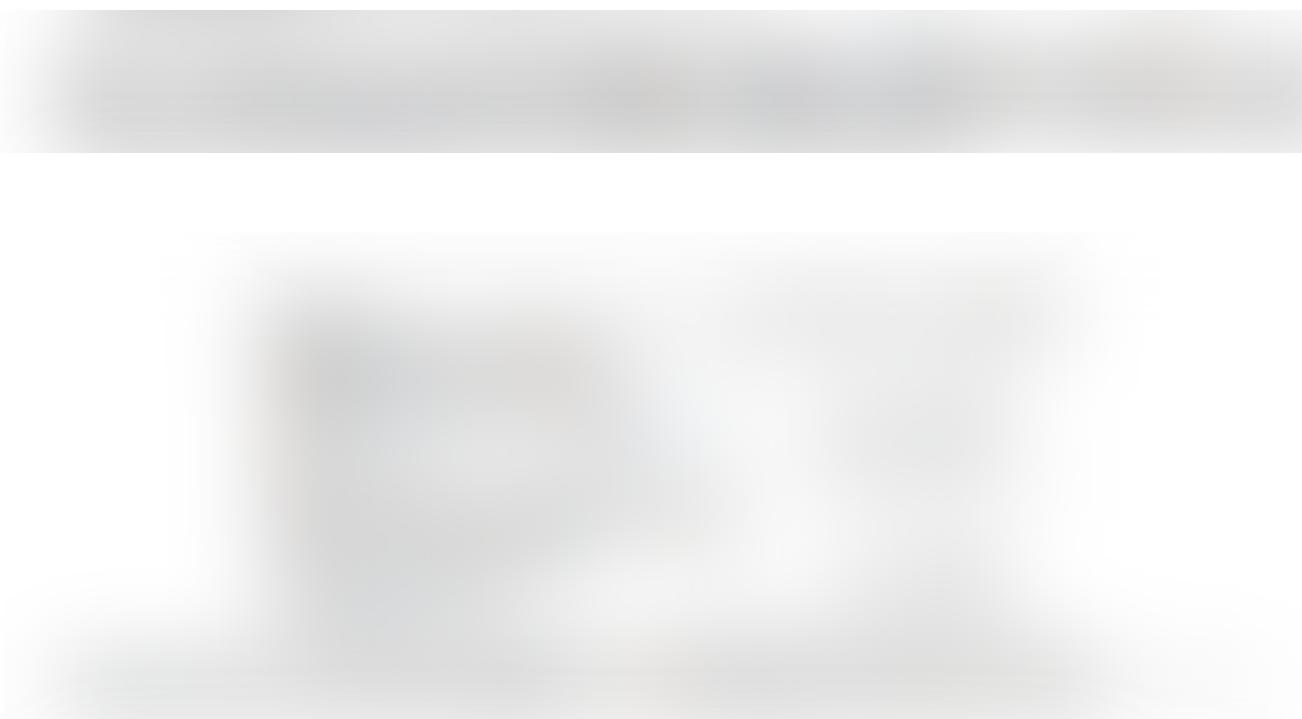


source

The authors also propose additional loss functions for guiding the student in generating high-quality audio streams:

- Power loss — to ensure that the power in different frequency bands of the speeches is used, as in human speech.
- Perceptual loss — for this loss, the authors experimented with feature reconstruction loss (the Euclidean distance between feature maps in the classifier) and style loss (the Euclidean distance between the Gram matrices). They found that style loss produced better results.
- Contrastive loss that penalizes waveforms that have high likelihood regardless of the conditioning vector.

The figure below shows the performance of this model.



source

• • •

Neural Voice Cloning with a Few Samples

The authors of this paper are from Baidu Research. They introduce a neural voice cloning system that learns to synthesize a person's voice from a few audio samples.

The two approaches used are speaker adaptation and speaker encoding. Speaker adaptation works by fine-tuning a multi-speaker generative model, while speaker encoding works by training a separate model to directly infer a new speaker embedding that's applied to the multi-speaker generative model.

Neural Voice Cloning with a Few Samples

Voice cloning is a highly desired feature for personalized speech interfaces.
Neural network based speech synthesis has...

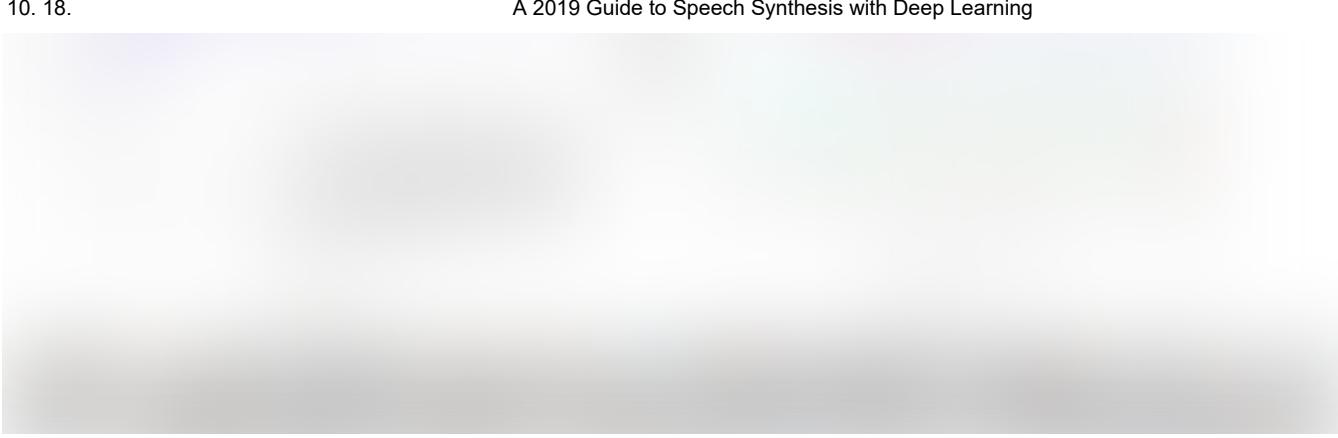
arxiv.org

This paper uses Deep Voice 3 as the baseline for the multi-speaker model. For voice cloning, the authors extract speaker characteristics from a speaker and generate audio provided that text from a given speaker is available.

The performance metrics used for the generated audio are speech naturalness and speaker similarity. They propose a speaker encoding method that directly estimates a speaker's embeddings from the audio samples of an unseen speaker.

source

Below is a look at how voice cloning performs.



source

• • •

VoiceLoop: Voice Fitting and Synthesis via A Phonological Loop

The authors of this paper are from Facebook AI Research. They introduce a neural text-to-speech (TTS) technique that can transform text into speech from voices that have been sampled from the wild.

VoiceLoop: Voice Fitting and Synthesis via a Phonological Loop

We present a new neural text to speech (TTS) method that is able to transform text to speech in voices that are sampled...

arxiv.org

VoiceLoop is inspired by a working memory model known as a phonological loop, which holds verbal information for a short time. It's comprised of a phonological store that's constantly being replaced, and a rehearsal process that maintains longer-term representations in the phonological store.

VoiceLoop constructs a phonological store by implementing a shifting buffer as a matrix. Sentences are represented as a list of phonemes. A short vector is then decoded from each of the phonemes. The current context vector is generated by weighing the encoding of the phonemes and summing them at each time point.

Some of the properties that make VoiceLoop different include the use of a memory buffer instead of the conventional RNNs, memory sharing between all processes, and using shallow, fully-connected networks for all computations.



source

Below is a look at how the model performs in comparison to other alternatives.

source

• • •

Natural TTS Synthesis by Conditioning WaveNet on Mel Spectrogram Predictions

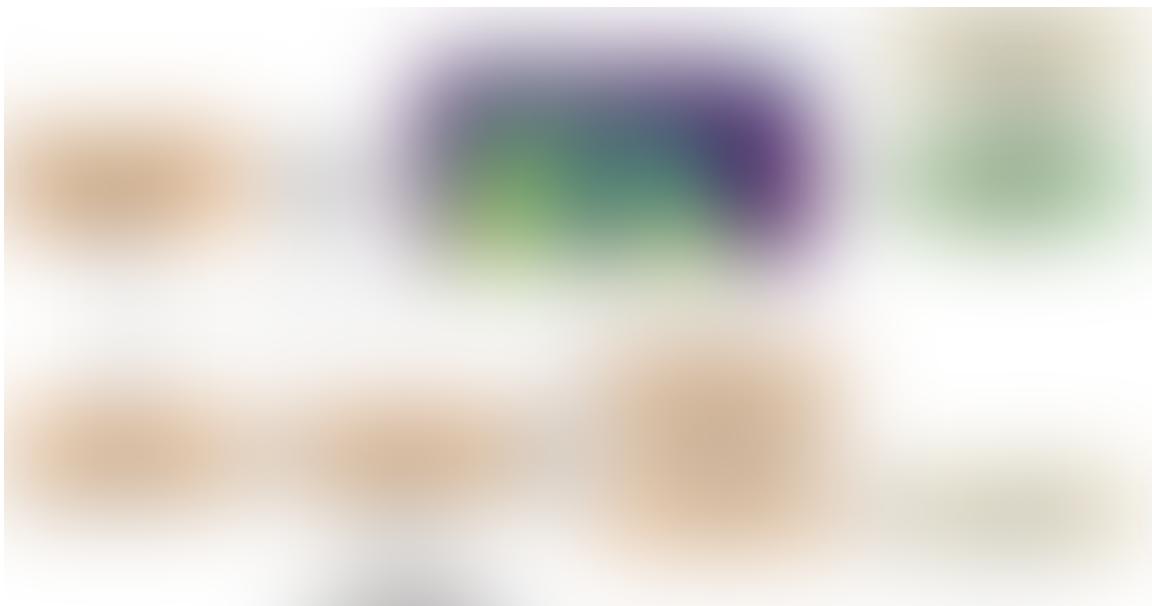
The authors of this paper are from Google and the University of California, Berkeley. They introduce Tacotron 2, a neural network architecture for speech synthesis from text.

Natural TTS Synthesis by Conditioning WaveNet on Mel Spectrogram Predictions

This paper describes Tacotron 2, a neural network architecture for speech synthesis directly from text. The system is...

[arxiv.org](https://arxiv.org/abs/1712.05884)

It's comprised of a recurrent sequence-to-sequence feature prediction network that maps character embeddings to mel-scale spectrograms. This is then followed by a WaveNet model that's been modified. This model acts as a vocoder that synthesizes time-domain waves from the spectrograms. The model achieves a mean opinion score (MOS) of 4.53.





source

This model has been built by combining the best features of Tacotron and WaveNet.
Below is the performance of the model in comparison to alternative models.



source

• • •

Conclusion

We should now be up to speed on some of the most common — and a couple of very recent — techniques for performing speech synthesis in a variety of contexts.

The papers/abstracts mentioned and linked to above also contain links to their code implementations. We'd be happy to see the results you obtain after testing them.

• • •

Editor's Note: Heartbeat is a contributor-driven online publication and community dedicated to exploring the emerging intersection of mobile app development and machine learning. We're committed to supporting and inspiring developers and engineers from all walks of life.

Editorially independent, Heartbeat is sponsored and published by Fritz AI, the machine learning platform that helps developers teach devices to see, hear, sense, and think. We pay our contributors, and we don't sell ads.

If you'd like to contribute, head on over to our call for contributors. You can also sign up to receive our weekly newsletters (Deep Learning Weekly and Heartbeat), join us on Slack, and follow Fritz AI on Twitter for all the latest in mobile machine learning.

Thanks to Austin Kodra.

[Speech Synthesis](#) [Machine Learning](#) [Heartbeat](#) [Guides And Tutorials](#) [Tech](#)

[About](#) [Help](#) [Legal](#)