

Week 8 Tuesday Paper 1: [WaveNet: A Generative Model for Raw Audio](#)

Explained by

- ▶ **WaveNet by Google DeepMind | Two Minute Papers #93**
- ▶ **WaveNet (Continued) | Lecture 75 (Part 1) | Applied Deep L...**
- ▶ **WaveNet (Theory and Implementation)**

Supplemental: [Fast Wavenet Generation Algorithm](#)

1. What are the input and output in “WaveNet”? How is the output of the model evaluated? Be precise.

Ans: The input to a WaveNet model is a raw audio waveform of the speech signal. This waveform is typically preprocessed through μ -law companding transformation and quantization to 256 possible values to make the audio representation more tractable for the model.

The output of the WaveNet model is a sequence of audio samples. Specifically, it generates a categorical distribution over possible quantized audio sample values for the next timestep, conditioned on all previous timesteps.

The output of the WaveNet model is evaluated subjectively through listening tests where human raters assess the naturalness of the generated speech. The model is also evaluated on its ability to generate speech with different voices by conditioning on the speaker identity and on its performance in music generation and speech recognition tasks.

2. How are conditional distributions over the individual audio samples modeled in the paper?

Ans: In the paper, conditional distributions over individual audio samples are modeled as categorical distributions using a softmax function. Despite the continuous nature of raw audio, this approach is preferred as it does not make assumptions about the shape of the audio distribution, allowing the model to capture arbitrary distributions effectively.

3. Explain convolutions, causal convolutions and dilated causal convolutions.

Ans: Convolutions are a mathematical operation used for filtering signals or computing features from data. They combine information from input data with a filter to produce a transformed version of the data.

Causal Convolutions ensure that the prediction for the current audio sample does not depend on future samples, maintaining the temporal ordering of the audio data. This is achieved by padding the input in such a way that the convolution outputs for a given time step depend only on that timestep and previous timesteps.

Dilated Causal Convolutions extend causal convolutions by introducing gaps in the input data they process, which allows the network to increase its receptive field exponentially with depth without a significant increase in computational cost. This is crucial for modeling the long-range dependencies in audio data.

4. Are residual and skip connections used in the paper? What is their role?

Ans: Yes, residual and skip connections are used in the WaveNet paper to facilitate training of deeper networks by addressing issues like vanishing gradients. Residual connections help in propagating gradients back through the network by adding the input of a convolutional block to its output, while skip connections aggregate signals from different layers and pass them to a final softmax layer to produce the output.

5. What do you think is the computational bottleneck while implementing WaveNet? How would you address it?

Ans: The computational bottleneck in implementing WaveNet is primarily due to the autoregressive nature of the model, which requires sequential generation of audio samples. This makes parallelization difficult and increases generation time.

One approach to address this issue could be the use of parallel WaveNet, a student-teacher framework where a faster, parallelizable student network is trained to imitate the slower, autoregressive teacher network. Another approach could be reducing the model complexity or employing more efficient sampling methods.

6. What is meant by conditional WaveNet?

Ans: Conditional WaveNet refers to versions of the model that are conditioned on additional inputs such as text for speech synthesis, speaker identity, or musical tags. This conditioning allows the network to generate audio that adheres to specific attributes dictated by these inputs.

7. Can WaveNets be conditioned in a local way? Explain with an example.

Ans: WaveNets can be locally conditioned by providing a secondary time series (e.g., linguistic features in text-to-speech synthesis) that influences the audio generation process. This local conditioning allows for variations in the generated audio that correspond to the input features, like generating speech with the correct phonemes and intonations corresponding to the input text.

8. Can WaveNets be conditioned in a global way? Explain with an example.

Ans: Global conditioning in WaveNets involves conditioning the model on global attributes such as speaker identity. This means that the generated audio will possess characteristics consistent with the global condition, such as a particular speaker's voice, while the specifics of the audio content can still vary. An example would be generating speech in the same voice across different sentences or contexts.

9. What questions do you have about the paper?

Ans: What are the other ways in which we can address the computational bottleneck while implementing WaveNet?

Week 8 Tuesday Paper 2: [Conformer: Convolution-augmented Transformer for Speech Recognition](#)

Explained by

📺 **Conformer: Convolution-augmented Transformer for Speech...**

Supplemental: [GitHub - sooftware/conformer: \[Unofficial\] PyTorch implementation of "Conformer: Convolution-augmented Transformer for Speech Recognition" \(INTERSPEECH 2020\)](#)

1. What is the central idea of a “Conformer”?

Ans: The central idea of the Conformer for speech recognition, is to combine the strengths of both Transformer and CNN models to model both local and global dependencies of an audio sequence in a parameter-efficient way. Transformers capture content-based global interactions, while CNNs exploit local features effectively. The Conformer integrates these approaches to provide improved speech recognition performance.

2. What is a half-step residual connection? What is its role in the paper?

Ans: A half-step residual connection refers to the novel architecture within each Conformer block, where two feed-forward modules sandwich the self-attention and convolution modules. Each FFN contributes only half of its output to the next module in the sequence, creating a half-step residual connection. This is inspired by the Macaron-Net architecture and is used instead of a single feed-forward layer. This structure has been shown to improve the model's performance by providing a significant improvement over having a single feed-forward module.

3. What is meant by layer normalization? Why is it used in “Conformer”?

Ans: Layer normalization is a technique used within neural networks to normalize the inputs across the features for each layer. In the Conformer model, layer normalization is applied within each residual unit, as well as before the first linear layer in the feed-forward module. It is used to help stabilize the training and improve the convergence of deeper models. This normalization technique is essential in

Conformer's architecture for regularizing and improving the efficiency and effectiveness of the network.

4. Is the combination of convolution and self-attention beneficial? Explain in detail.

Ans: The combination of convolution and self-attention in the Conformer is beneficial as it allows the model to learn both position-wise local features and content-based global interactions. This combination leads to improved modeling of audio sequences by capturing the detailed, fine-grained local features through convolutions while maintaining the ability to understand broader, context-based features through self-attention. The paper demonstrates that this combination significantly outperforms models that rely on either approach individually, achieving state-of-the-art results.

5. What baselines are used in the paper? What datasets are used in the paper?

Ans: The baselines used in the Conformer paper include Hybrid Transformer, CTC, QuartzNet, LAS, Transformer, and various versions of ContextNet. The dataset used for evaluating the Conformer model is the widely used LibriSpeech benchmark, which consists of 970 hours of labeled speech.

6. Explain in detail the ablation studies carried out.

Ans: The ablation studies in the Conformer paper explore various configurations and components of the Conformer model, including the impact of convolution blocks, the Macaron-style feed-forward layers, the number of attention heads, and the convolution kernel sizes. The studies show the significance of each component and configuration on the model's performance, demonstrating the effectiveness of the Conformer's architectural choices.

7. What are the limitations of this work?

Ans: The only limitations are the model's complexity and the computational requirements.

8. What is meant by macaron-like feed-forward layers?

Ans: Macaron-like feed-forward layers in the Conformer model refer to the architectural choice where each Conformer block contains two feed-forward modules, one before and one after the self-attention and convolution modules. These are used with half-step residuals, splitting the traditional feed-forward layer of the Transformer block into two parts and providing improvements in performance and efficiency.

