

# WaveNet: A Generative Model for Raw Audio

Aäron van den Oord et al., 2016

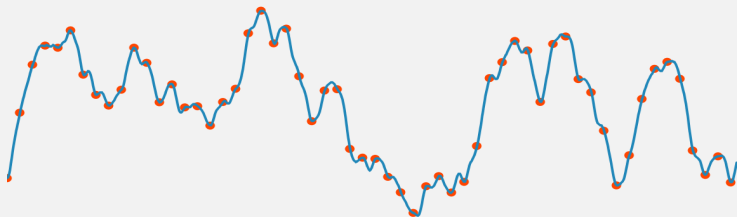
Rui Yoshinaga

July 21, 2021

## WaveNet: A Generative Model for Raw Audio

- Topic: Audio Generation
- Contributions:
  - Text-to-Speech (TTS) : Generating natural speech
  - Develop new architectures
  - Various audio generation by conditioning
  - Application to music generation, etc.
- Generated samples:
  - <https://deepmind.com/blog/article/wavenet-generative-model-raw-audio>

# Composition



waveform  $\mathbf{x} = \{x_1, \dots, x_T\} \in [-1, 1]^T$

## Modeling of $p(\mathbf{x})$

- The joint probability of waveform  $\mathbf{x}$  is factorised:

$$p(\mathbf{x}) = \prod_{t=1}^T p(x_t \mid x_1, \dots, x_{t-1}) \quad (1)$$

- Sample  $x_t$  is conditioned on the previous samples  $x_1, \dots, x_{t-1}$

## Architecture

- NN that takes inputs  $x_{t-L}, \dots, x_{t-1}$  and outputs  $x_t$
- Modeling by CNN (Dilated Causal Convolution)
  - c.f. image generation [1, 2], w/o pooling layer

## Training and Generation

- Training: maximize the log-likelihood
- Generation: autoregression

# Softmax Distribution

## Input/Output

- Input: previous samples  $x_{t-L}, \dots, x_{t-1}$
- Output: the probability distribution of next sample  $x_t$

## $\mu$ -law companding transformation

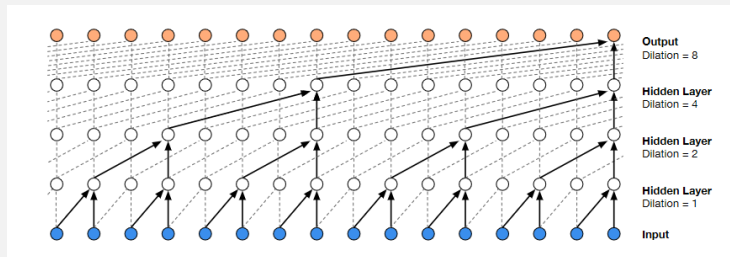
- raw audio: quantized at 16bits, amplitude:  $2^{16}$  levels
- Apply  $\mu$ -law companding transformation  $\rightarrow$  256 levels

$$f(x_t) = \text{sign}(x_t) \frac{\ln(1 + \mu|x_t|)}{\ln(1 + \mu)}$$

where  $-1 < x_t < 1$  and  $\mu = 255$

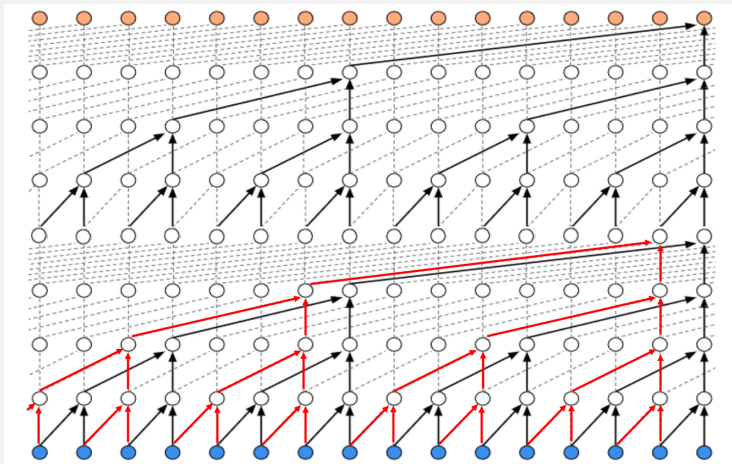


# Dilated Convolution



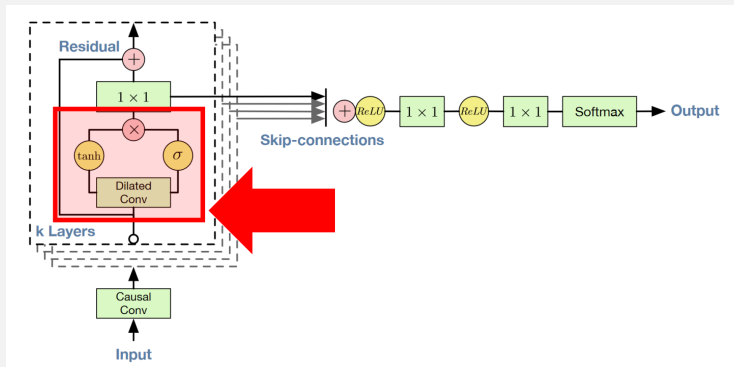
- Skip inputs at fixed intervals
- In WaveNet, dilation is doubled ( $1, 2, 4, \dots, 512, 1, 2, \dots$ )
  - Exponential increase in dilation  $\rightarrow$  larger receptive field
  - Stacking blocks further increases the size of receptive field

# Stacked Dilated Convolution





# Gated Activation Unit

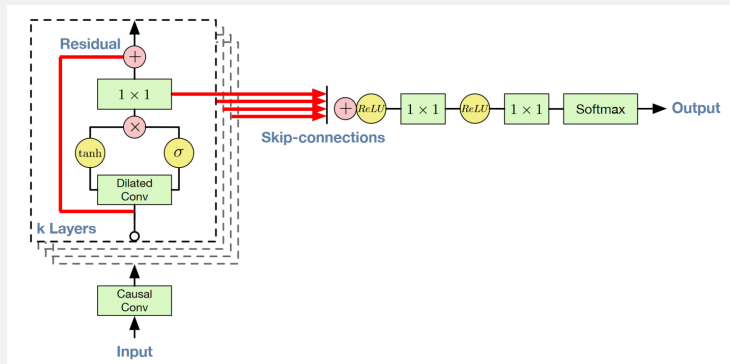


- Gated activation unit [2]

$$z = \tanh(W_{f,k} * \mathbf{x}) \odot \sigma(W_{g,k} * \mathbf{x}) \quad (2)$$

- $W_{*,k}$  : learnable convolution filter in  $k$ -th layer
- $*$  : convolution operator
- Worked significantly better than ReLU

# Residual and Skip Connection



- Residual [3] and Skip Connections
  - Speed up convergence
  - Enable training of much deeper model

- Modeling conditional distribution  $p(\mathbf{x}|\mathbf{h})$  by additional input  $\mathbf{h}$
- Eq.(1) now becomes

$$p(\mathbf{x} \mid \mathbf{h}) = \prod_{t=1}^T p(x_t \mid x_1, \dots, x_{t-1}, \mathbf{h}) \quad (3)$$

- Generate audio with the required characteristics by conditioning
  - e.g.1 speaker ID (who speaks)
  - e.g.2 text for TTS (what to speak)
- Conditioning in 2 different ways:
  - {Global, Local} conditioning

- Conditioning on single  $\mathbf{h}$  that influences output across all timesteps  
e.g. TTS: speaker embedding (who speaks)
- The activation function from Eq.(2) now becomes:

$$\mathbf{z} = \tanh(W_{f,k} * \mathbf{x} + V_{f,k}^T \mathbf{h}) \odot \sigma(W_{g,k} * \mathbf{x} + V_{g,k}^T \mathbf{h})$$

- $W_{*,k}$  : learnable convolution filter in  $k$ -th layer
- $*$  : convolution operator
- $V_{*,k}$  : learnable linear projection in  $k$ -th layer
- $V_{*,k}^T \mathbf{h}$  : broadcast over the time dimension

- Condition on the second time series  $\mathbf{h} = \{h_t\}$ 
  - $\mathbf{h}$  has a lower sampling frequency than audio  $\mathbf{x}$   
e.g. TTS: linguistic features
- Map  $\mathbf{h}$  to a new time series  $\mathbf{y} = f(\mathbf{h})$ 
  - $f$  : transposed convolutional network (learned upsampling)
  - (also possible to use  $V_{f,k} * \mathbf{h}$  repeatedly across time)
- The activation function from Eq.(2) now becomes:
$$\mathbf{z} = \tanh(W_{f,k} * \mathbf{x} + V_{f,k} * \mathbf{y}) \odot \sigma(W_{g,k} * \mathbf{x} + V_{g,k} * \mathbf{y})$$
  - $W_{*,k}, V_{*,k}$  : learnable convolution filter in  $k$ -th layer
  - $*$  : convolution operator

# Experiments (Music)

## Music Datasets

1. MagnaTagATune dataset
  - each 29-second clip is annotated with tags
2. YouTube piano dataset

## Receptive field

- several seconds: long-range consistency ×
- enlarging the receptive field → musical samples

## Generating results

- Harmonic and aesthetically pleasing even no conditioning
- Conditioning on tags
  - Train with one-hot vector corresponds to the tags
  - Control generation by conditioning on the one-hot vector
- <https://www.deepmind.com/blog/wavenet-generative-model-raw-audio/>

## Multi-Speaker Speech Generation

- Conditioned on the speaker: feed the speaker ID as one-hot vector
- Generate non-existent but human language-like words realistically
- Mimicked the breathing and mouth movements of the speakers

## Text-To-Speech

- Conditioned on the linguistic features and the value of  $\log f_0$
- Achieved 5-scale MOSs in naturalness above 4.0
- Significantly better than the baseline

## Speech Recognition

- Partially change architecture and loss
- Obtained the best score as model trained directly on raw audio

- WaveNet: A generative model for raw audio
  - Directly generate audio as waveform (raw audio domain)
- Dilated Causal Convolutions
  - Enables large receptive field with exponentially increasing dilation
- Can be applied in various applications about audio signals
  - TTS, Music generation, Speech recognition



- [1] Aaron van den Oord, Nal Kalchbrenner, and Koray Kavukcuoglu.  
Pixel Recurrent Neural Networks, 2016.
- [2] Aaron van den Oord, Nal Kalchbrenner, Oriol Vinyals, Lasse Espeholt,  
Alex Graves, and Koray Kavukcuoglu.  
Conditional Image Generation with PixelCNN Decoders, 2016.
- [3] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun.  
Deep Residual Learning for Image Recognition, 2015.