

# Parallel WaveNet: Fast High-Fidelity Speech Synthesis

Kamlyk Erik 171

# WaveNet vs. Parallel WaveNet



Figure 1: A second of generated speech.

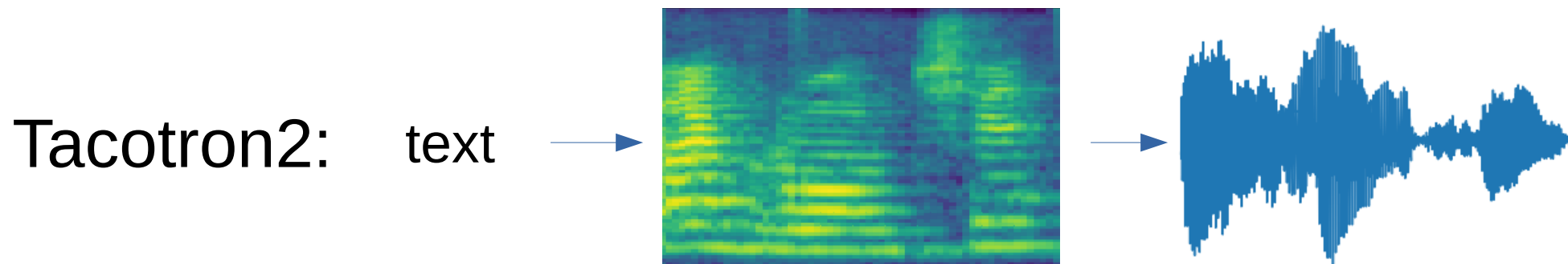
## WaveNet (2016)

- high fidelity
- autoregressive
- slow inference

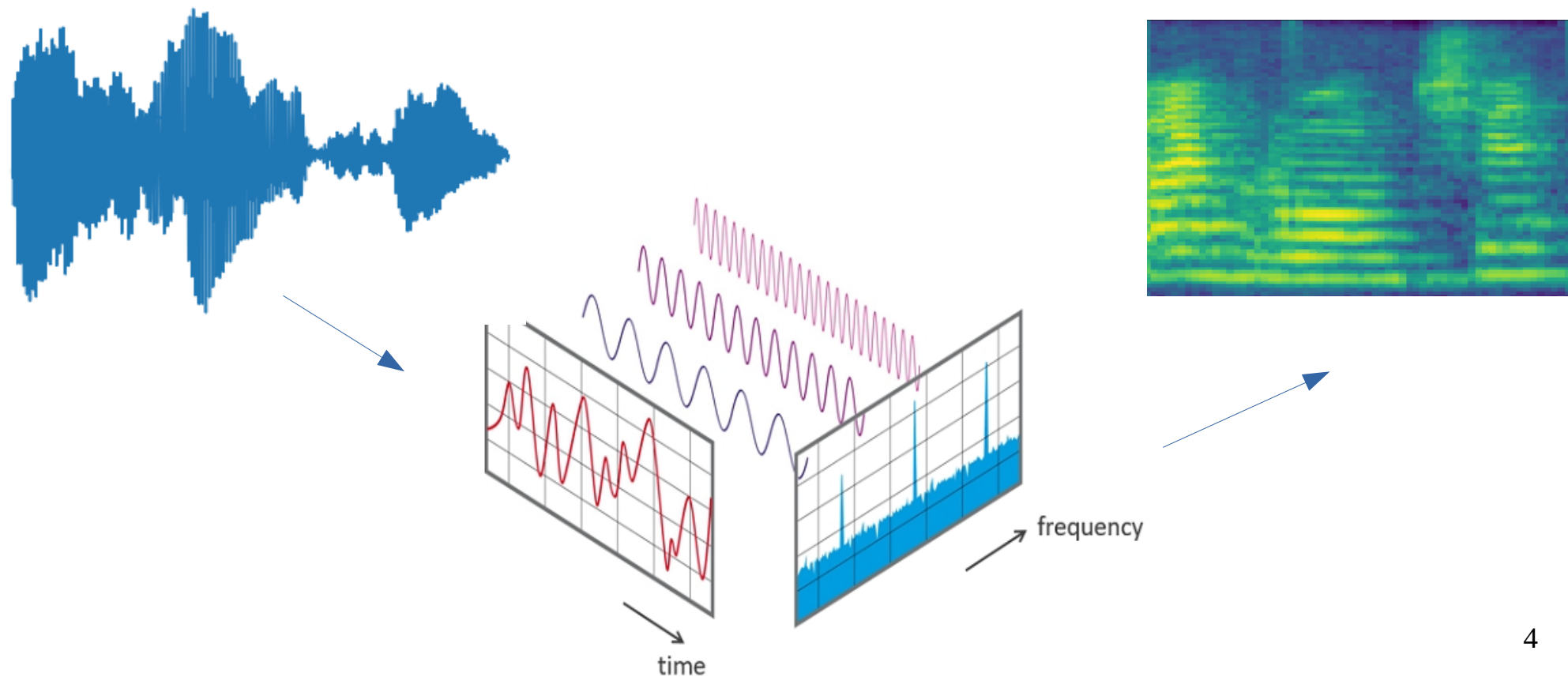
## Parallel WaveNet (2017)

- distilled WaveNet
- non autoregressive
- fast inference
- very small quality drop

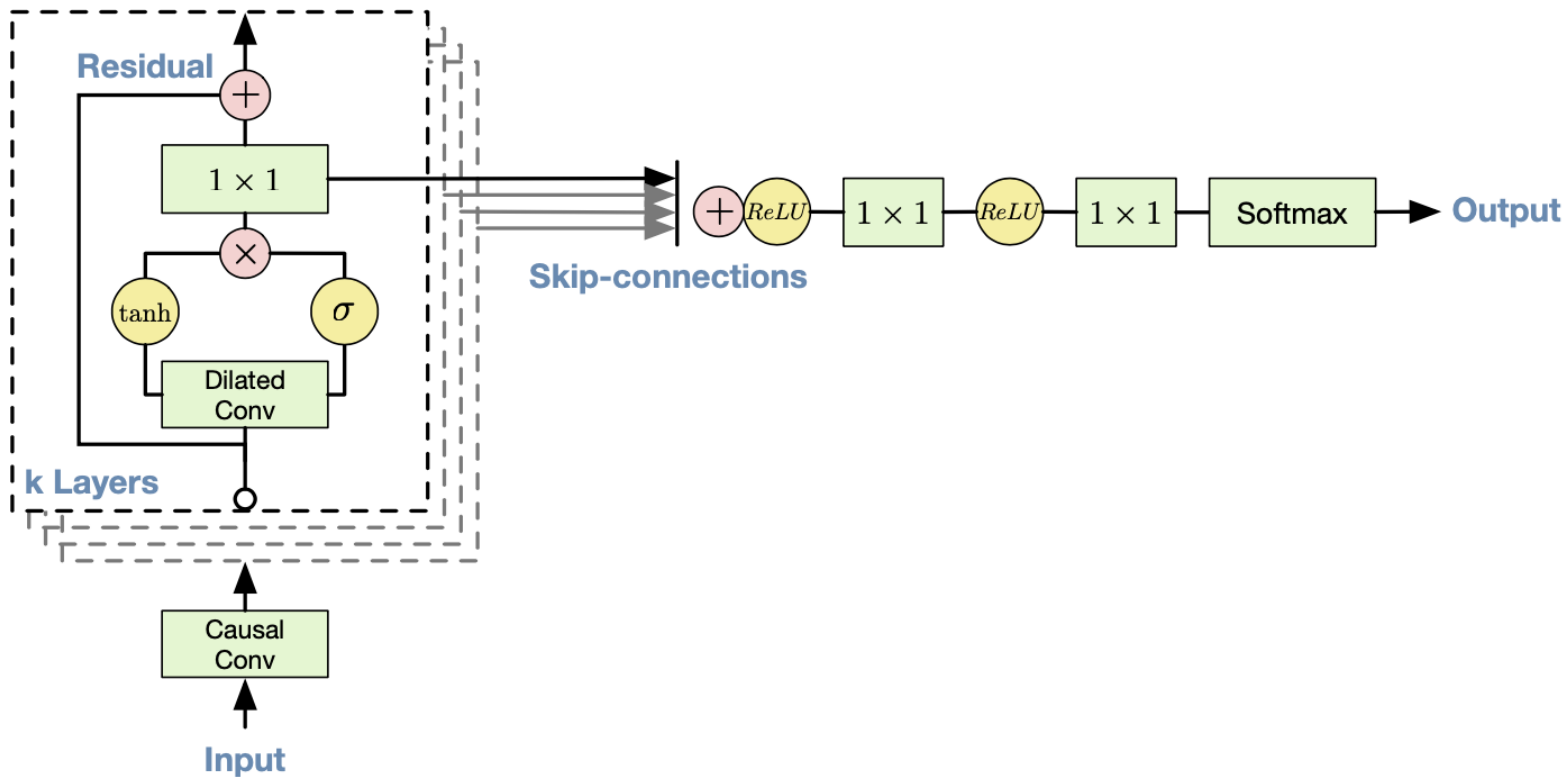
# Overall pipeline



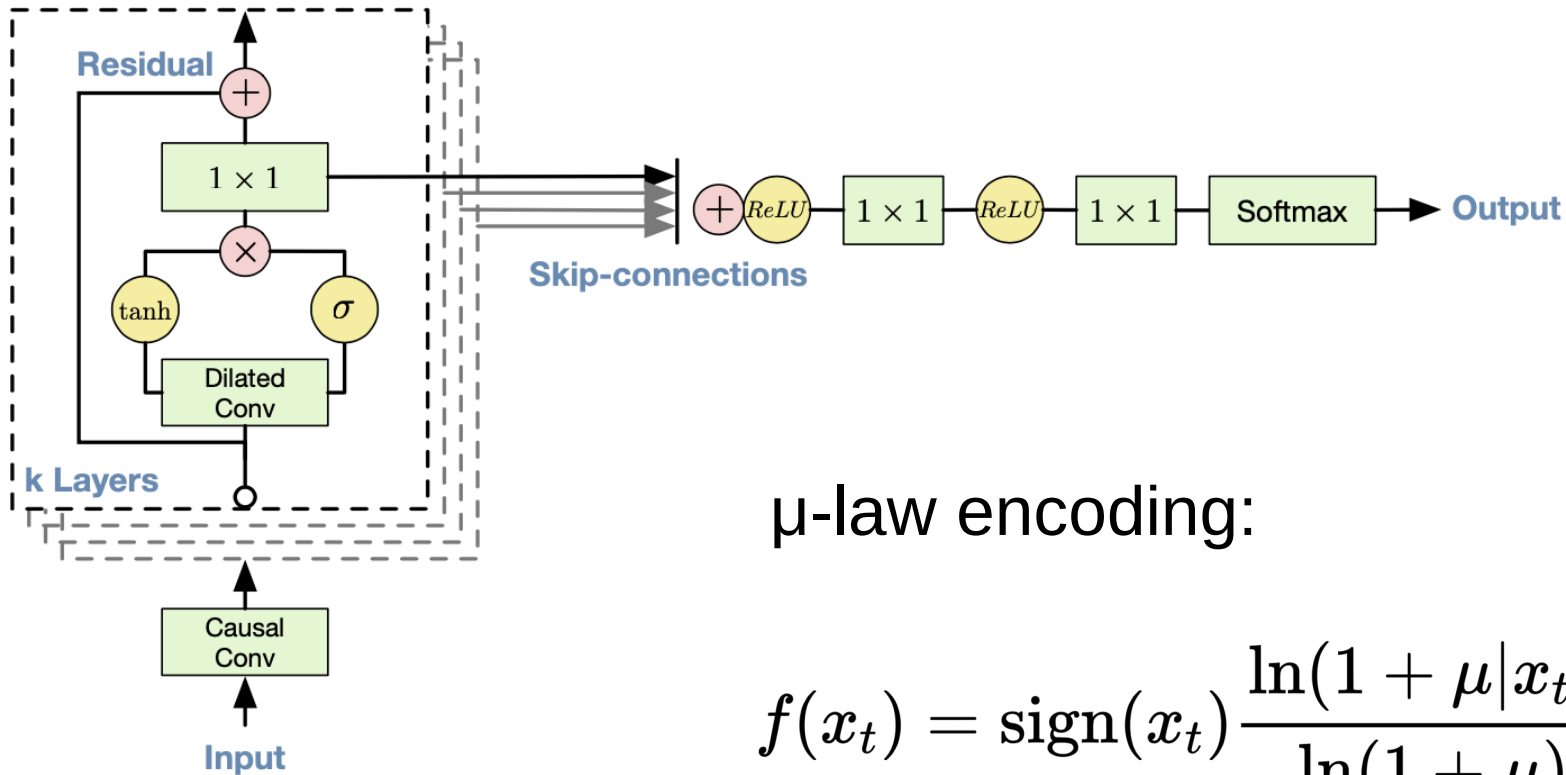
# Mel Spectrogram



# WaveNet



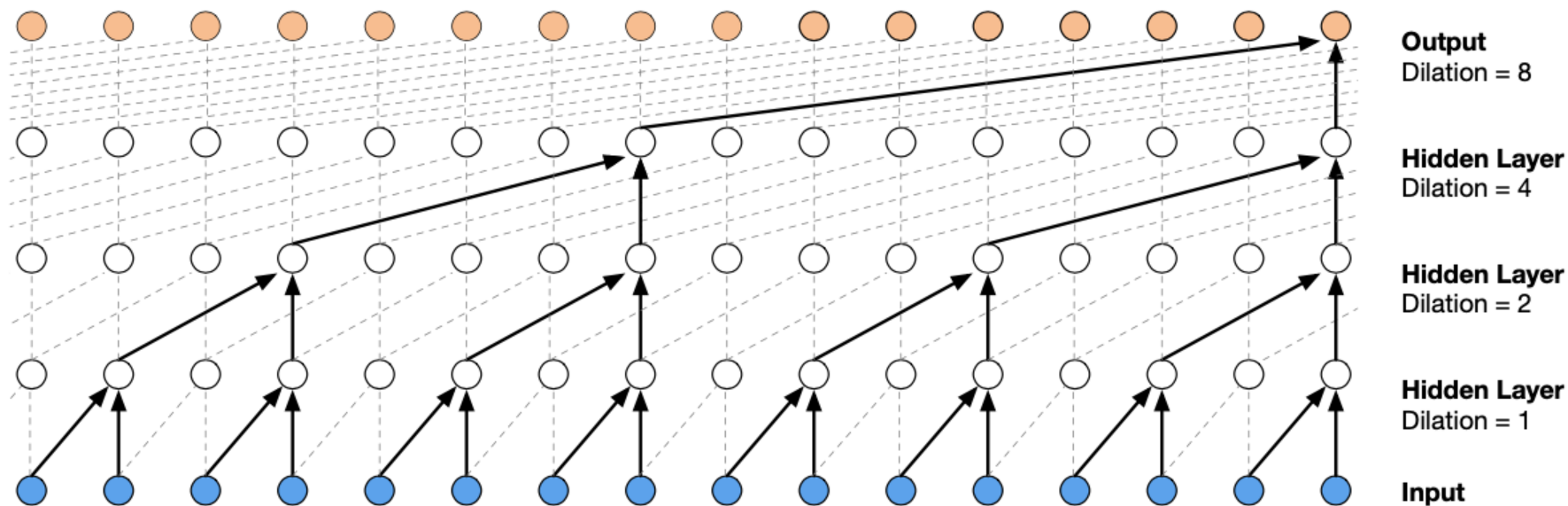
# WaveNet



$\mu$ -law encoding:

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

# WaveNet



- causal convolutions (mask or padding)
- dilated convolutions

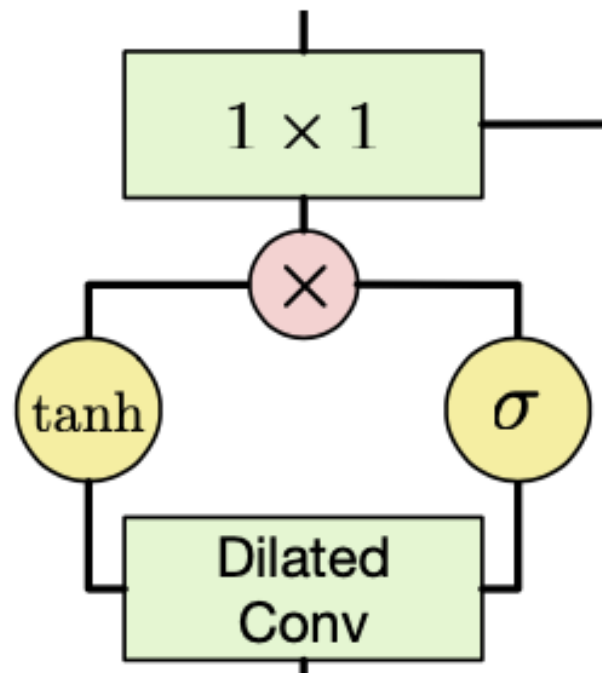
# WaveNet

Gated layer:

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

Conditioning on  $\mathbf{y}$ :

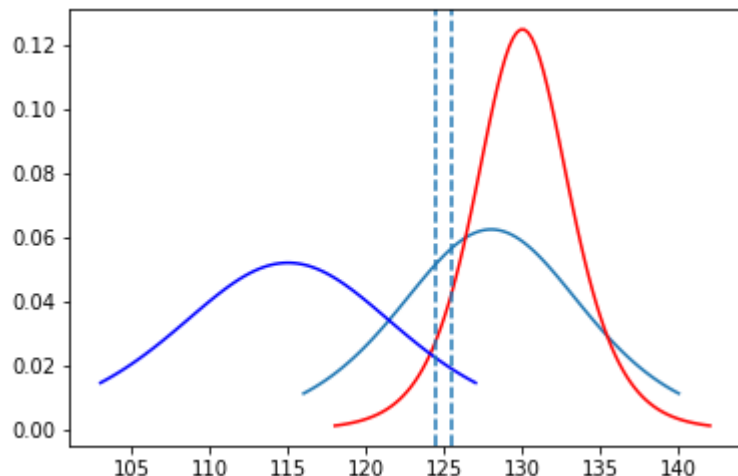
$$\mathbf{z} = \tanh(W_{f,k} * \mathbf{x} + V_{f,k} * \mathbf{y}) \odot \sigma(W_{g,k} * \mathbf{x} + V_{g,k} * \mathbf{y})$$





# High-fidelity WaveNet

1. 16-bit encoding (65536 values) instead of 8 bits (256 values).
2. Discretized mixture of logistics instead of softmax.
3. 24kHz instead of 16kHz (convolution filter size 2  $\rightarrow$  3)



$$\nu \sim \sum_{i=1}^K \pi_i \text{logistic}(\mu_i, s_i)$$
$$P(x|\pi, \mu, s) = \sum_{i=1}^K \pi_i [\sigma((x + 0.5 - \mu_i)/s_i) - \sigma((x - 0.5 - \mu_i)/s_i)],$$

# Normalizing flows

$$\mathbf{x} = f(\mathbf{z})$$

$$\log p_X(\mathbf{x}) = \log p_Z(\mathbf{z}) - \log \left| \frac{d\mathbf{x}}{d\mathbf{z}} \right|,$$

where  $\left| \frac{d\mathbf{x}}{d\mathbf{z}} \right|$  is the determinant of the Jacobian of  $f$

If  $f$  is invertible and its Jacobian determinant is easy to compute we can optimize maximum likelihood.

Composition of functions can be used for better quality:

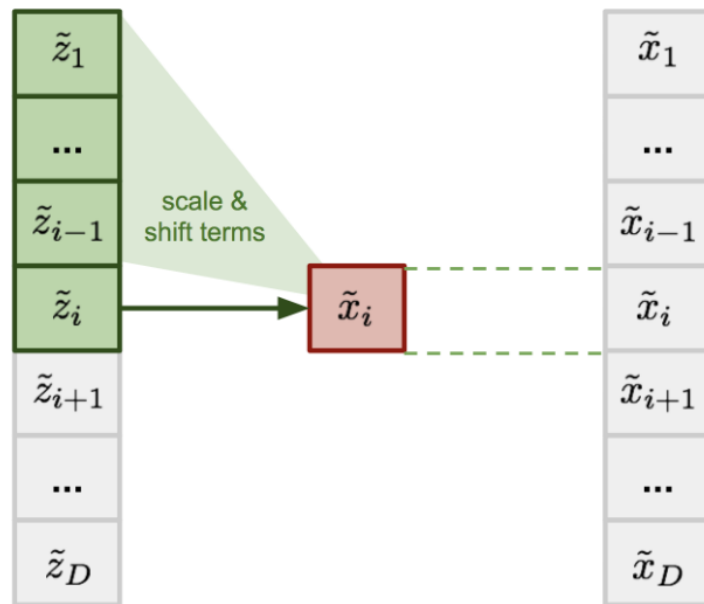
$$f_\theta = f^{(L)} \circ f^{(L-1)} \circ \dots \circ f^{(1)}$$

# Inverse Autoregressive Flows

$$x_t = f(\mathbf{z}_{\leq t})$$

$$\log \left| \frac{dx}{dz} \right| = \sum_t \log \frac{\partial f(\mathbf{z}_{\leq t})}{\partial z_t}$$

$$x_t = z_t \cdot s(\mathbf{z}_{<t}, \boldsymbol{\theta}) + \mu(\mathbf{z}_{<t}, \boldsymbol{\theta})$$



# Parallel WaveNet

$$\mathbf{z} \sim \text{Logistic}(0, I)$$

$$x_t = z_t \cdot s(\mathbf{z}_{<t}, \boldsymbol{\theta}) + \mu(\mathbf{z}_{<t}, \boldsymbol{\theta})$$

After N iterations:

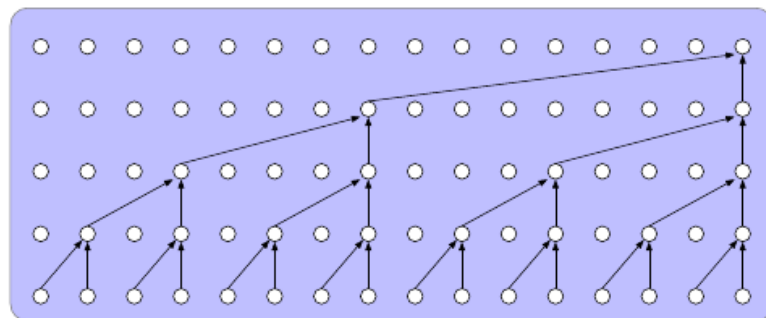
$$\mu_{\text{tot}} = \sum_i^N \mu^i \left( \prod_{j>i}^N s^j \right)$$

$$\mathbf{s}_{\text{tot}} = \prod_i^N \mathbf{s}_i$$

# Training

## WaveNet Teacher

Linguistic features  $\dashrightarrow$



Teacher Output

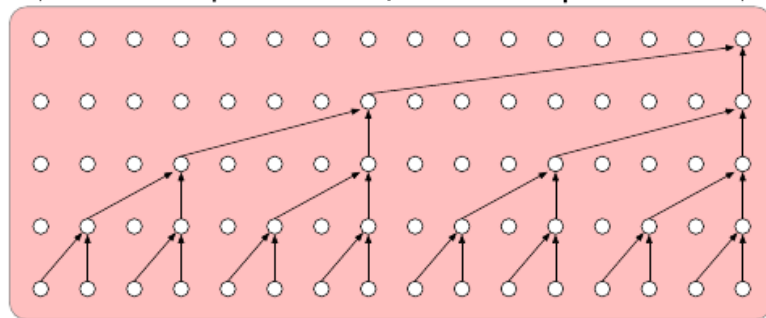
$$P(x_i | x_{<i})$$

Generated Samples

$$x_i = g(z_i | z_{<i})$$

## WaveNet Student

Linguistic features  $\dashrightarrow$



Student Output

$$P(x_i | z_{<i})$$

Input noise

$$z_i$$

# Loss function

$$D_{\text{KL}}(P_S || P_T) = H(P_S, P_T) - H(P_S)$$

$$\begin{aligned} H(P_S, P_T) &= \int_{\mathbf{x}} p_S(\mathbf{x}) \ln p_T(\mathbf{x}) \\ &= \sum_{t=1}^T \mathbb{E}_{p_S(\mathbf{x}_{<t})} H(p_S(x_t | \mathbf{x}_{<t}), p_T(x_t | \mathbf{x}_{<t})). \end{aligned}$$
$$\begin{aligned} H(P_S) &= \mathbb{E}_{z \sim L(0,1)} \left[ \sum_{t=1}^T -\ln p_S(x_t | \mathbf{z}_{<t}) \right] \\ &= \mathbb{E}_{z \sim L(0,1)} \left[ \sum_{t=1}^T \ln s(\mathbf{z}_{<t}, \boldsymbol{\theta}) \right] + 2T, \end{aligned}$$

# Additional loss terms

1) Power loss:  $\|\phi(g(\mathbf{z}, \mathbf{c})) - \phi(\mathbf{y})\|^2$

$$\phi(\mathbf{x}) = |\text{STFT}(\mathbf{x})|^2$$

2) Perceptual loss

3) Contrastive loss:

$$D_{\text{KL}}(P_S(\mathbf{c}_1) \| P_T(\mathbf{c}_1)) - \gamma D_{\text{KL}}(P_S(\mathbf{c}_1) \| P_T \mathbf{c}_2))$$

# Experiments: fidelity, speed

| Method  | Subjective 5-scale MOS |
|---|------------------------|
| <b>16kHz, 8-bit <math>\mu</math>-law, 25h data:</b> |                        |
| LSTM-RNN parametric [27]                            | $3.67 \pm 0.098$       |
| HMM-driven concatenative [27]                       | $3.86 \pm 0.137$       |
| WaveNet [27]  | $4.21 \pm 0.081$       |
| <b>24kHz, 16-bit linear PCM, 65h data:</b>          |                        |
| HMM-driven concatenative                            | $4.19 \pm 0.097$       |
| Autoregressive WaveNet                              | $4.41 \pm 0.069$       |
| Distilled WaveNet                                   | $4.41 \pm 0.078$       |

Table 1: Comparison of WaveNet distillation with the autoregressive teacher WaveNet, unit-selection (concatenative), and previous results from [27]. MOS stands for Mean Opinion Score.

WaveNet speed – 172 timesteps/second

Parallel WaveNet speed – over 500000 timesteps/second



# Experiments: multispeaker

|  | Parametric | Concatenative | Distilled WaveNet |
|--|------------|---------------|-------------------|
| <b>English speaker 1</b> (female - 65h data) | 3.88       | 4.19          | 4.41              |
| <b>English speaker 2</b> (male - 21h data)   | 3.96       | 4.09          | 4.34              |
| <b>English speaker 3</b> (male - 10h data)   | 3.77       | 3.65          | 4.47              |
| <b>English speaker 4</b> (female - 9h data)  | 3.42       | 3.40          | 3.97              |
| <b>Japanese speaker</b> (female - 28h data)  | 4.07       | 3.47          | 4.23              |

Table 2: Comparison of MOS scores on English and Japanese with multi-speaker distilled WaveNets. Note that some speakers sounded less appealing to people and always get lower MOS, however distilled parallel WaveNet always achieved significantly better results.

# Ablation studies

| <b>Method</b>                                     | <b>Preference Scores</b>   |
|---|--|
|   | <b>versus baseline concatenative system</b><br><b>Win - Lose - Neutral</b> |
| <b>Losses used</b>                                |  |
| KL + Power  | 60% - 15% - 25%  |
| KL + Power + Perceptual                           | 66% - 10% - 24%  |
| KL + Power + Perceptual + Contrastive (= default) | 65% - 9% - 26%   |

Table 3: Performance with respect to different combinations of loss terms. We report preference comparison scores since their mean opinion scores tend to be very close and inconclusive.

# Conclusions

1. WaveNet is an autoregressive model for generating waveforms from text; based on causal and dilated convolutions.
2. Parallel WaveNet is a model based on Inverse Autoregressive Flows. Model itself is not autoregressive, hence much faster generation.
3. Parallel WaveNet is slow to train, therefore it is trained using distillation from WaveNet.
4. The result is x1000 speedup without any drop in quality.

# References

1. Parallel WaveNet: Fast High-Fidelity Speech Synthesis  
<https://arxiv.org/abs/1711.10433>
2. WaveNet: A Generative Model for Raw Audio  
<https://arxiv.org/abs/1609.03499>
3. Improving Variational Inference with Inverse Autoregressive Flow
4. PixelCNN++: Improving the PixelCNN with Discretized Logistic Mixture Likelihood and Other Modifications  
<https://arxiv.org/abs/1701.05517>

# Вопросы:

1. Опишите архитектуру модели WaveNet (какие данные подаются на вход, как выглядит блок модели, какие свёртки используются, какая функция потерь).
2. Опишите процедуру обучения. Выпишите функцию потерь для Probability density distillation.
3. Выпишите формулу для contrastive loss. Зачем он используется?