

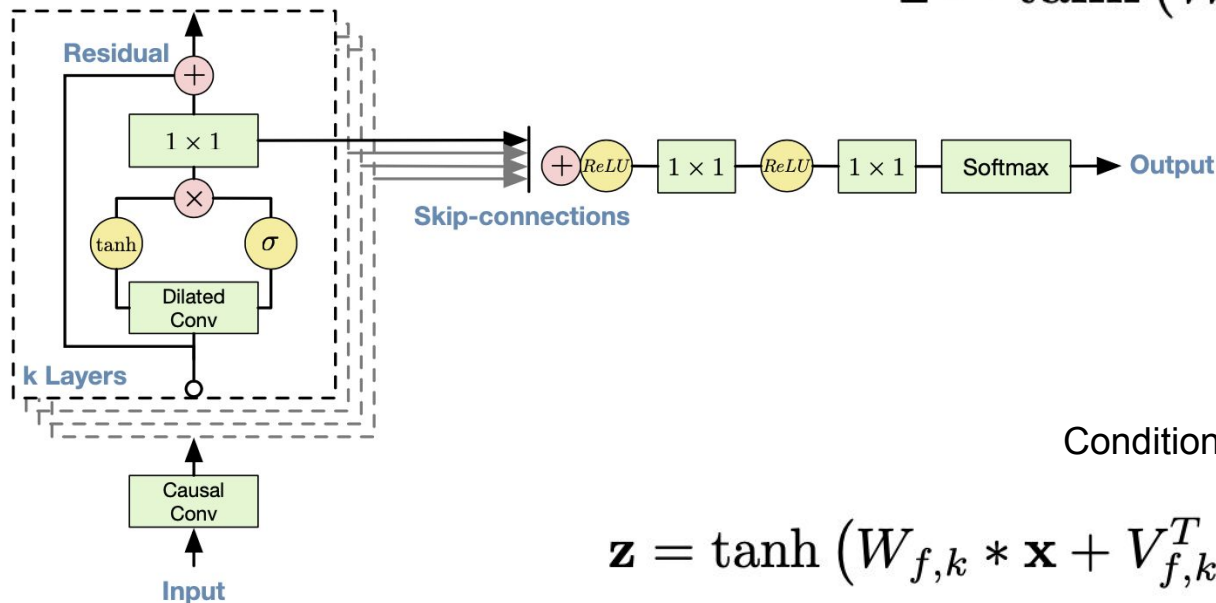
TTS: Vocoders

History lesson aka WaveNet

autoregressive model i.e.
$$p(x) = \prod_{t=1}^T p(x_t | x_1 \dots x_{t-1})$$

Gated Activation Unit

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



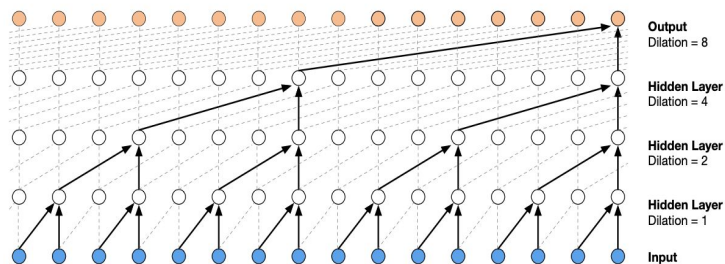
Conditional Gated Activation Unit

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

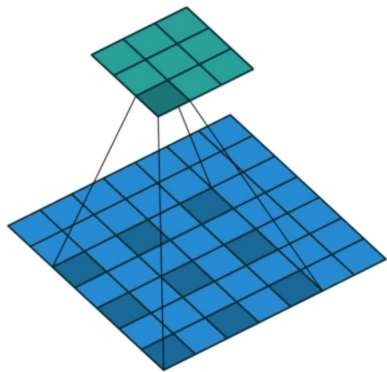
Dilated Convolution:

increase receptive field to better process long-term dependencies

1D

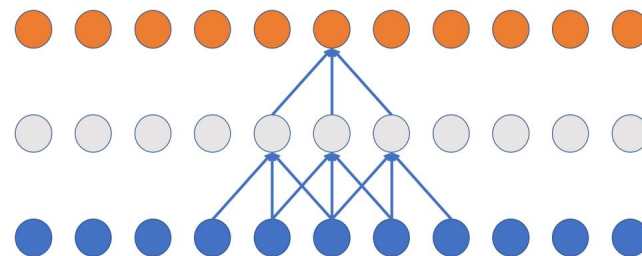


2D

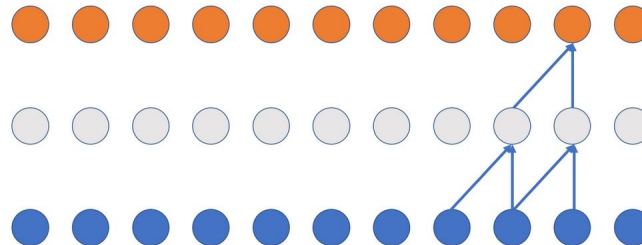


Causal Convolution:

Standard Convolution



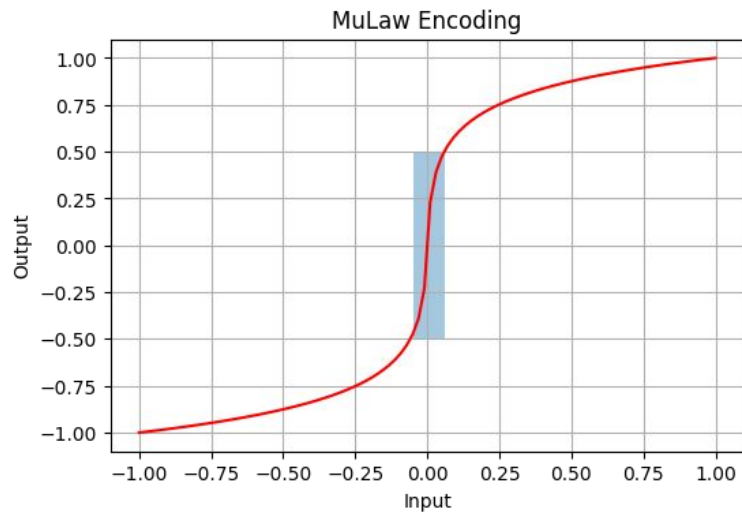
Causal Convolution



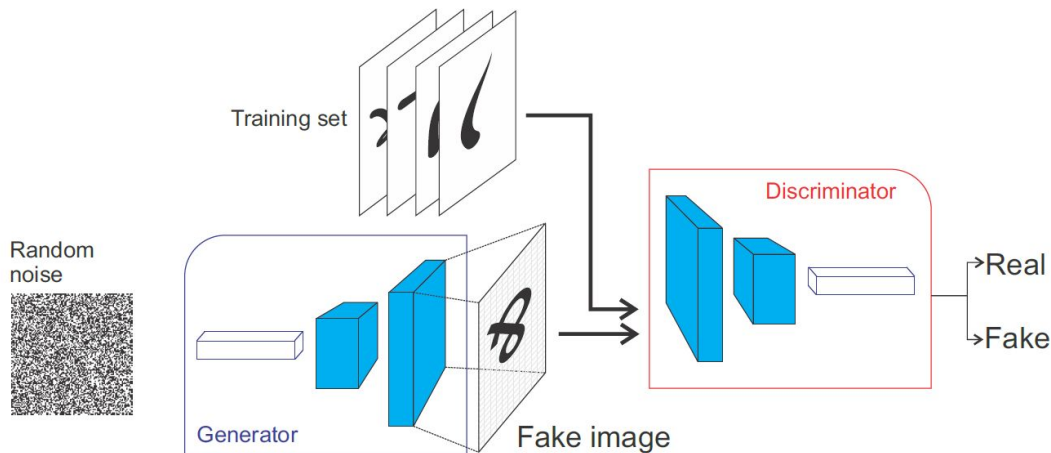
Mu Law Encoding

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

$$\mu = 255 \quad -1 \leq x_t \leq 1$$



Generative Adversarial Networks (GANs)



1) Vanilla adversarial loss:

$$G - \log(D(G(z_i)))$$

$$D - \log D(x_i) + \log(1 - D(G(z_i)))$$

2) Least Squares loss: (LS-GAN)

$$G - (D(G(z_i)) - 1)^2$$

$$D - (D(x_i) - 1)^2 + (D(G(z_i)))^2$$

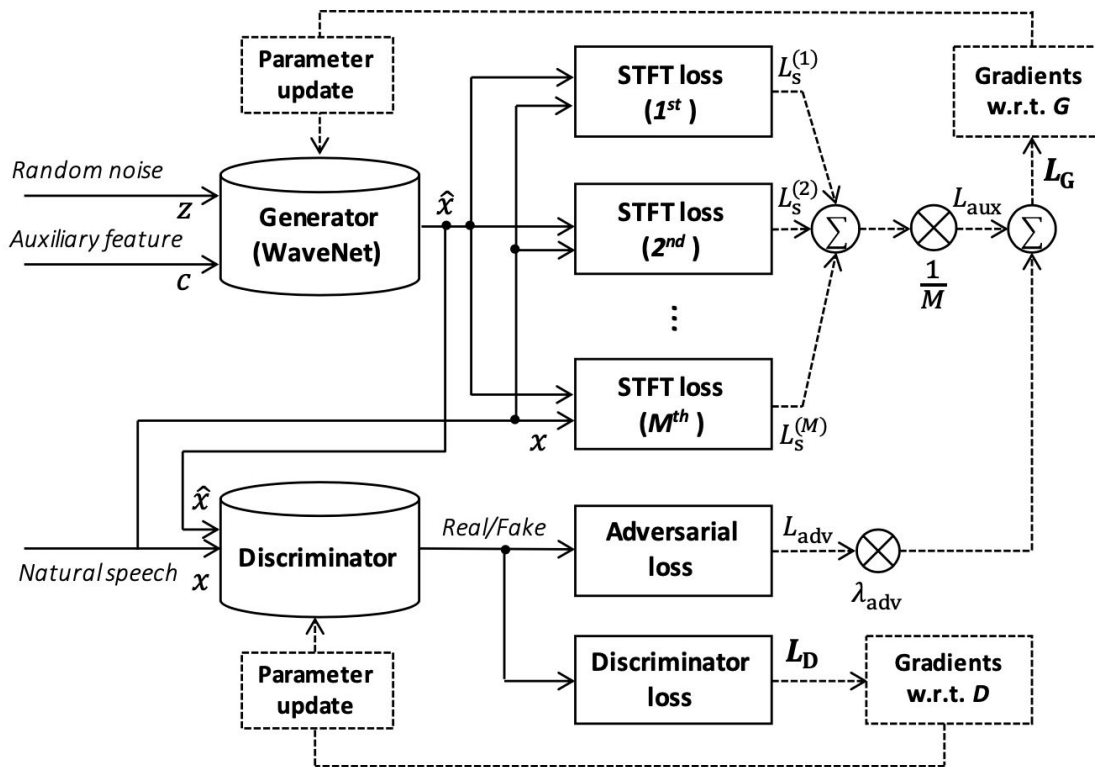
3) Feature Matching loss:

$$G - \sum_{i=1}^T \frac{1}{N_i} \|D_k^{(i)}(x_j) - D_k^{(i)}(G(z_j))\|_1$$

4) Markovian loss (aka PatchGan)

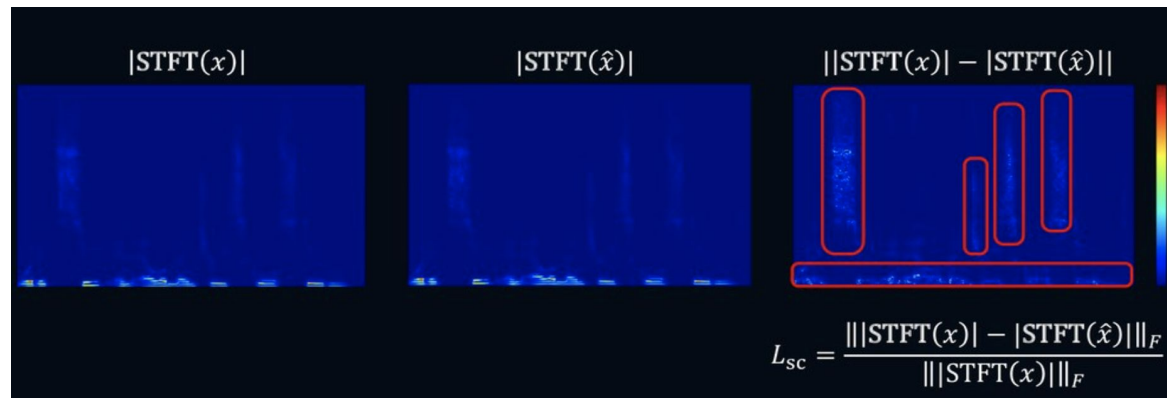
Parallel WaveGAN

- 2 stage training
- WaveNet generator
- multi resolution stft loss

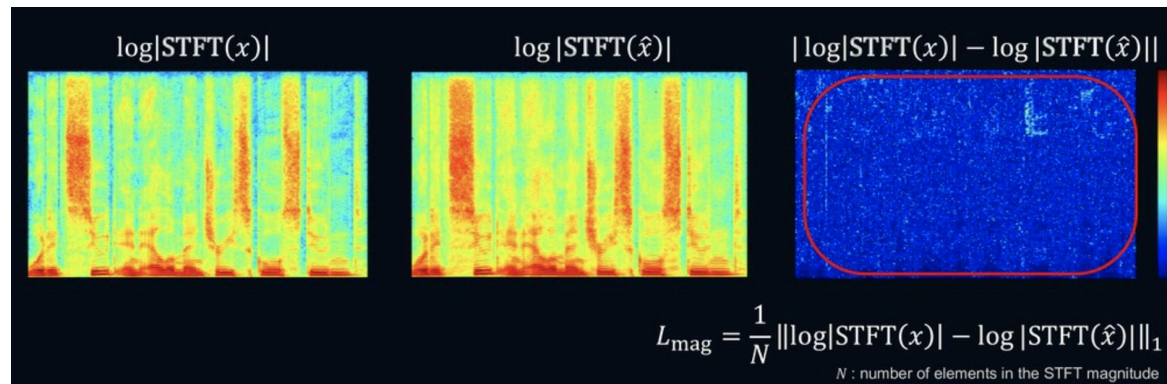


STFT Loss

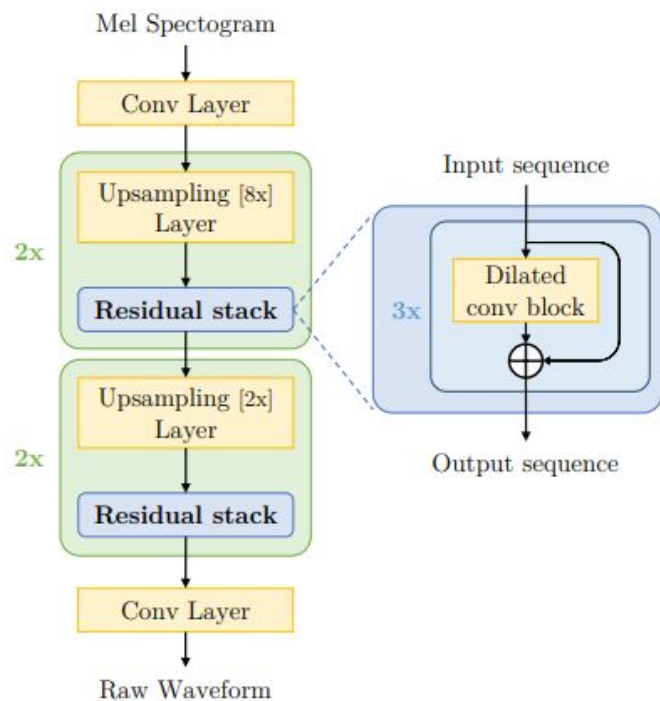
Spectral Convergence part



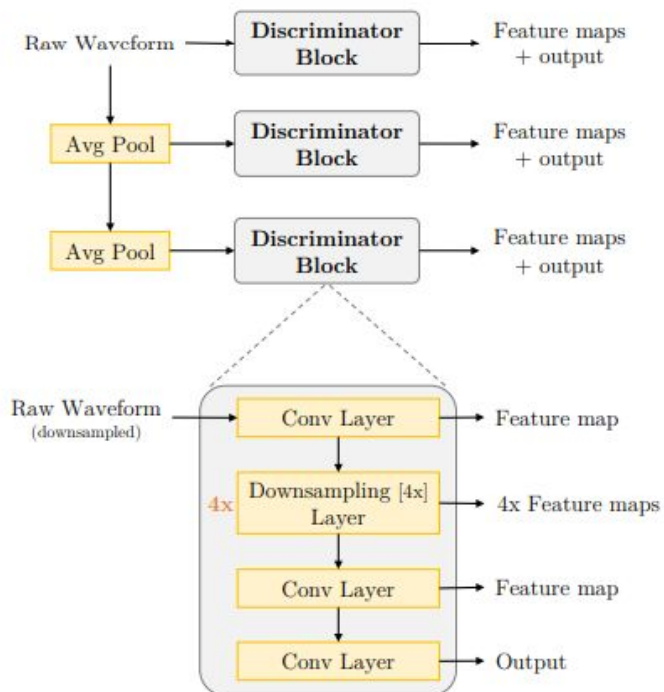
Log scale STFT magnitude part



MelGan



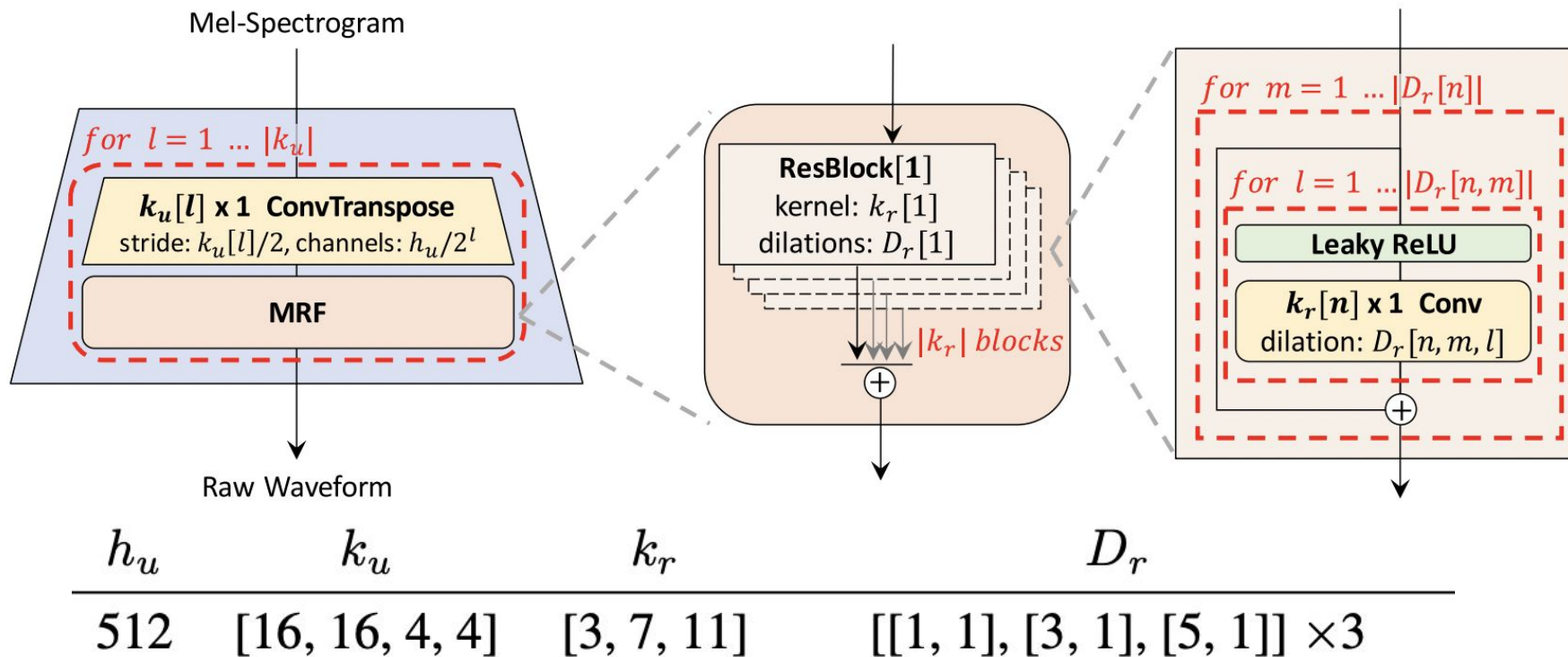
(a) Generator



(b) Discriminator

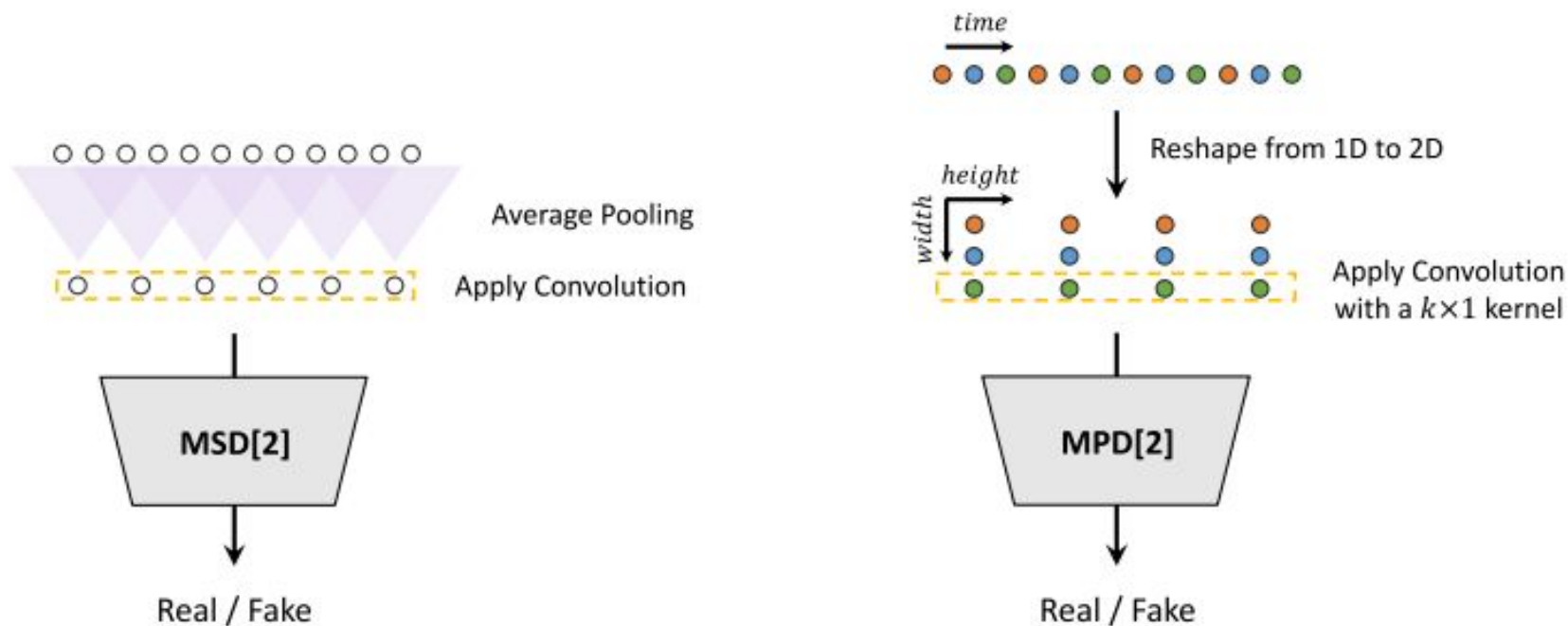
Hifi-GAN

Generator:

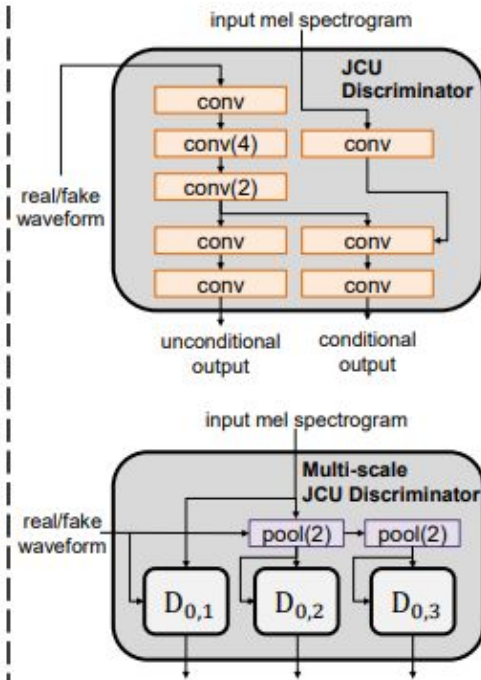
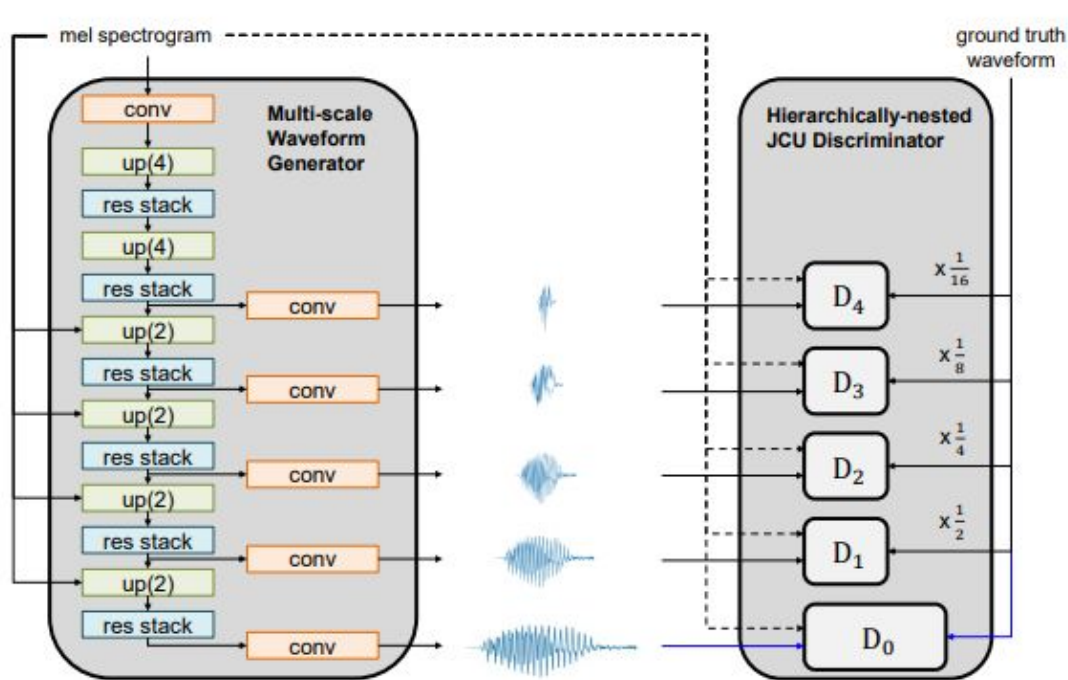


Hifi-GAN

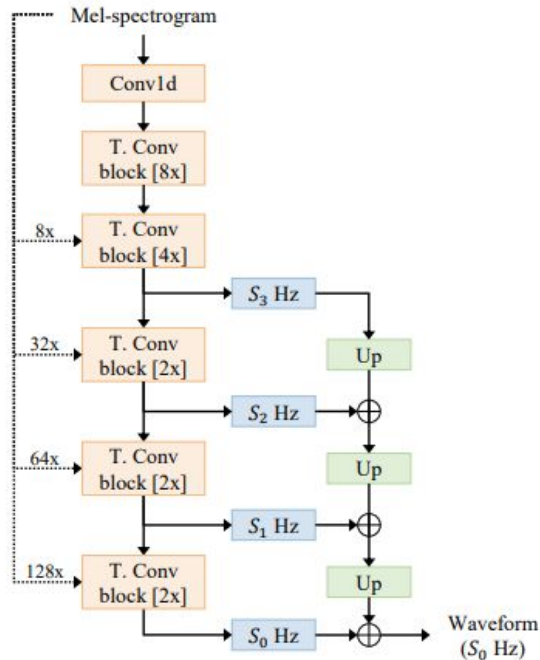
Discriminator:



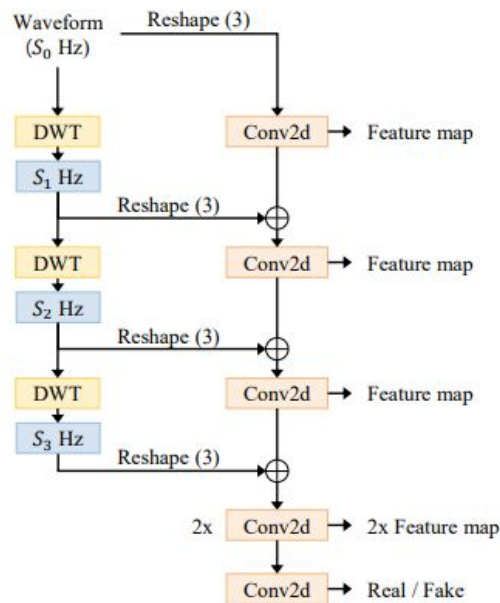
VocGAN



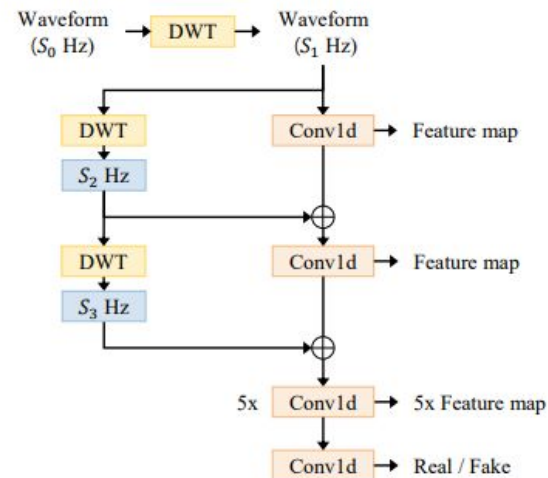
Fre-GAN



(a) RCG



(b) RPD



(c) RSD

Discrete Wavelet Transform (DWT)

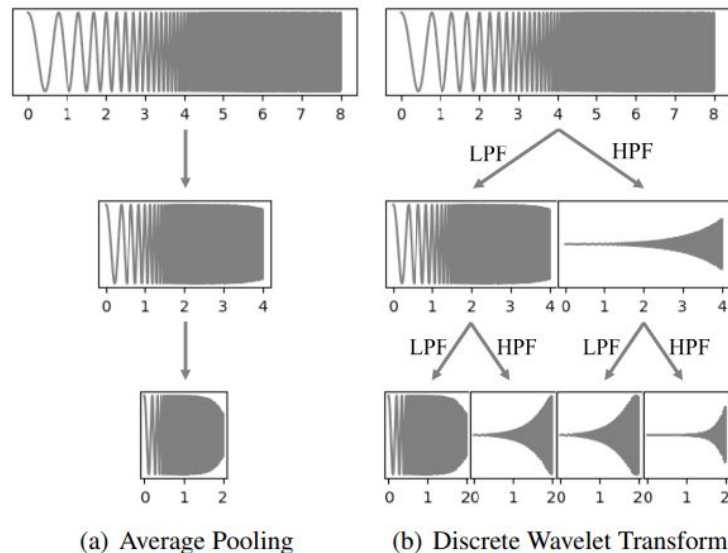
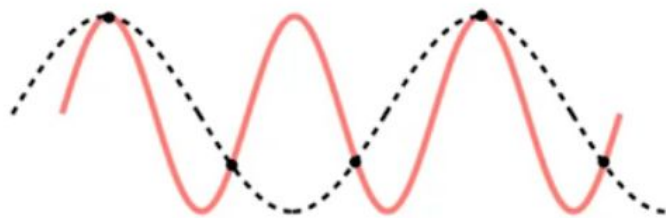


Figure 1: Comparison of Average Pooling (AP) and Discrete Wavelet Transform (DWT). Here, LPF and HPF refer to Low-Pass and High-Pass Filter, respectively. In this example, an up-chirp signal whose frequency increases from 0 Hz at time $t = 0$ to 150 Hz at time $t = 8$ is downsampled by AP and DWT.

Other vocoders:

FLOWs:

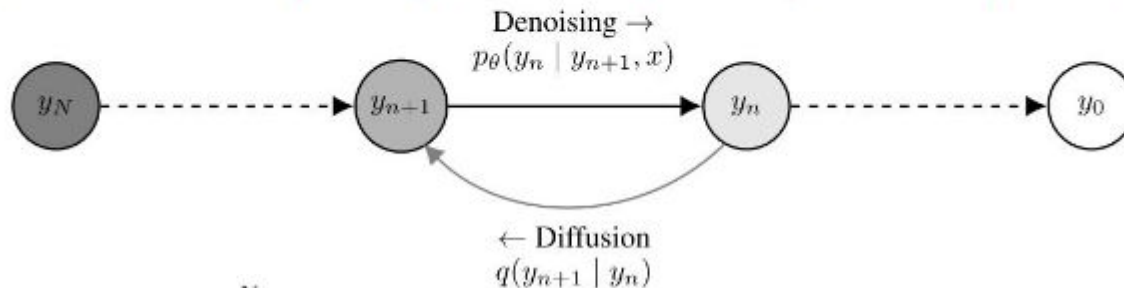
- [Parallel WaveNet](#)
- [WaveGlow](#)
- [WaveFlow](#)
- [NanoFlow](#)

Denoising Diffusion Probabilistic

Models (DDPM):

- [Diff-TTS](#)
- [more on DDPM](#)

Denoising Diffusion Probabilistic Models



$$q(y_1, \dots, y_N | y_0) := \prod_{n=1}^N q(y_n | y_{n-1})$$

$$q(y_n | y_{n-1}) := \mathcal{N}(y_t; \sqrt{1 - \beta_t} x_{n-1}, \beta_t I)$$

[more on DDPMs](#)

