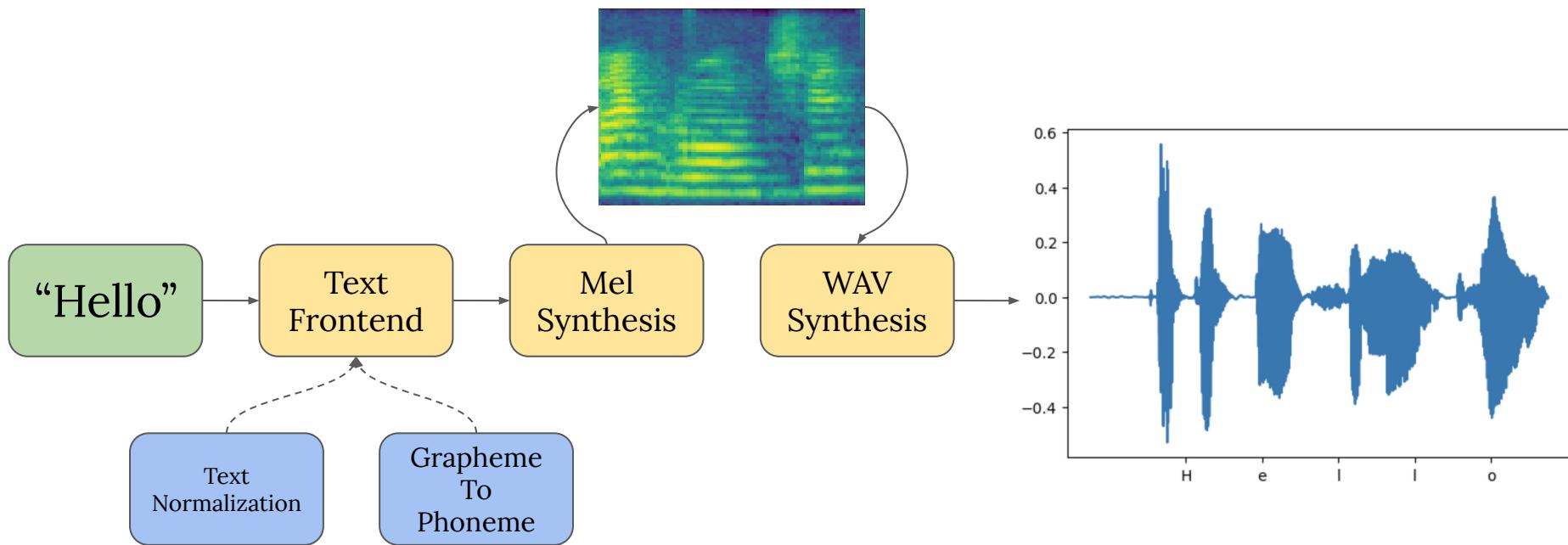


Speech Synthesis

What is a TTS in general?



Grapheme VS Phoneme

- A phoneme is simply a sound
- A grapheme is the smallest fundamental unit in written language
- There isn't bijection because depend on context

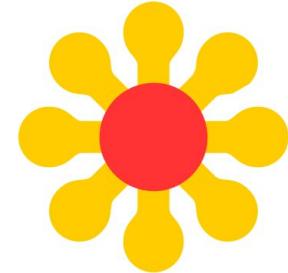
Phoneme	Example	Translation
AA	odd	AA D
AE	at	AE T
AH	hut	HH AH T
AO	ought	AO T
AW	cow	K AW

What About Public Datasets?

- [LJSpeech](#)
- [LibtiTTS](#)
- [CommonVoice](#)
- [OpenTTS \(Russian\)](#)

How to Measure Quality?

- There's no correct answer
- Subjective perception
- A lot of types of mistakes
- One solution is MOS
- Another is Side-by-side (SBS)



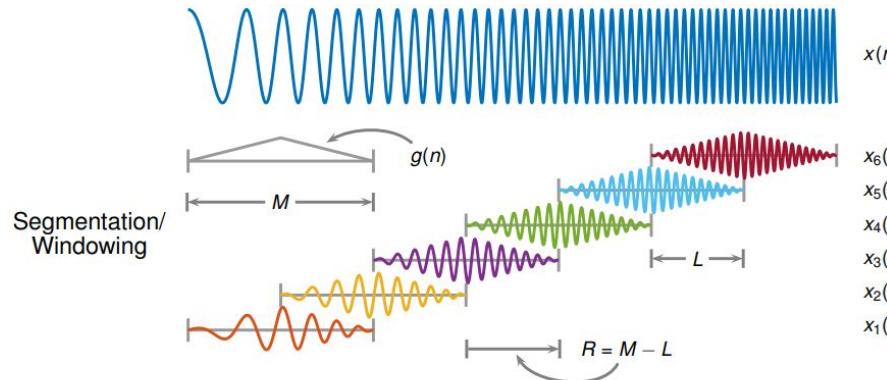
4b34552a61a45e2a5ede53f98225bb10
Где звучит лучше? — 2021-02-09

Last Thursday at 5:03 PM

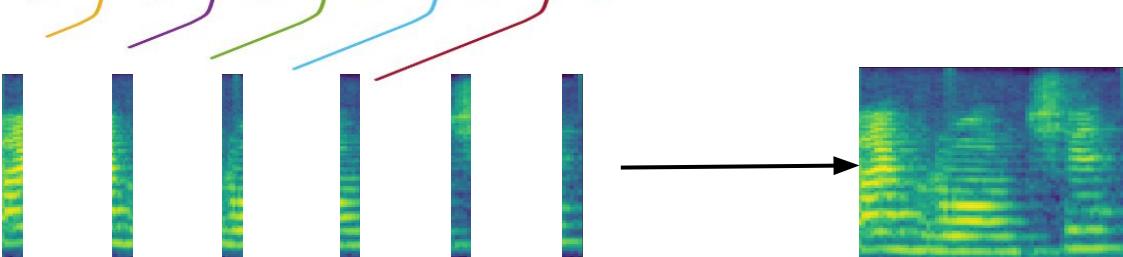
hello, i just wanted to do this task every day and night, this is one of my favorite task, so kindly give me unlimited task with best rate.

thanks

What is MelSpectrogram?



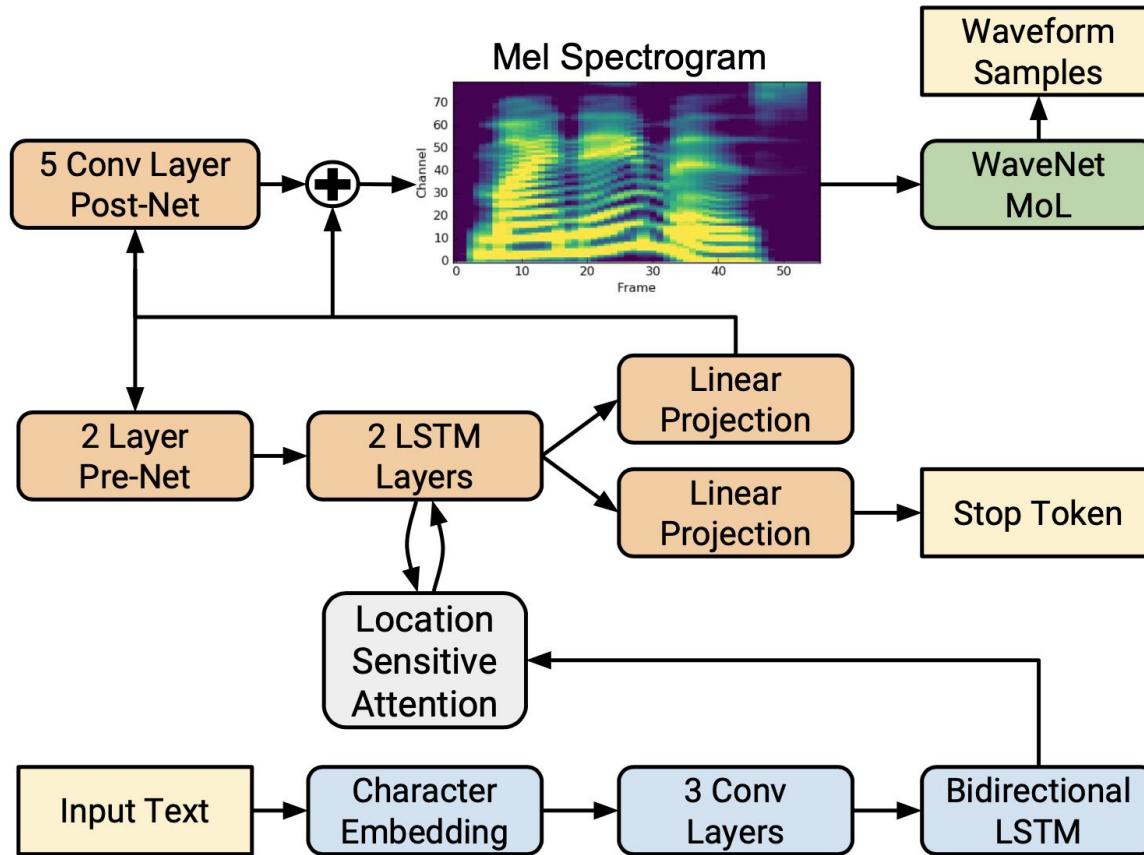
$\log(\text{MelScale}(x))$



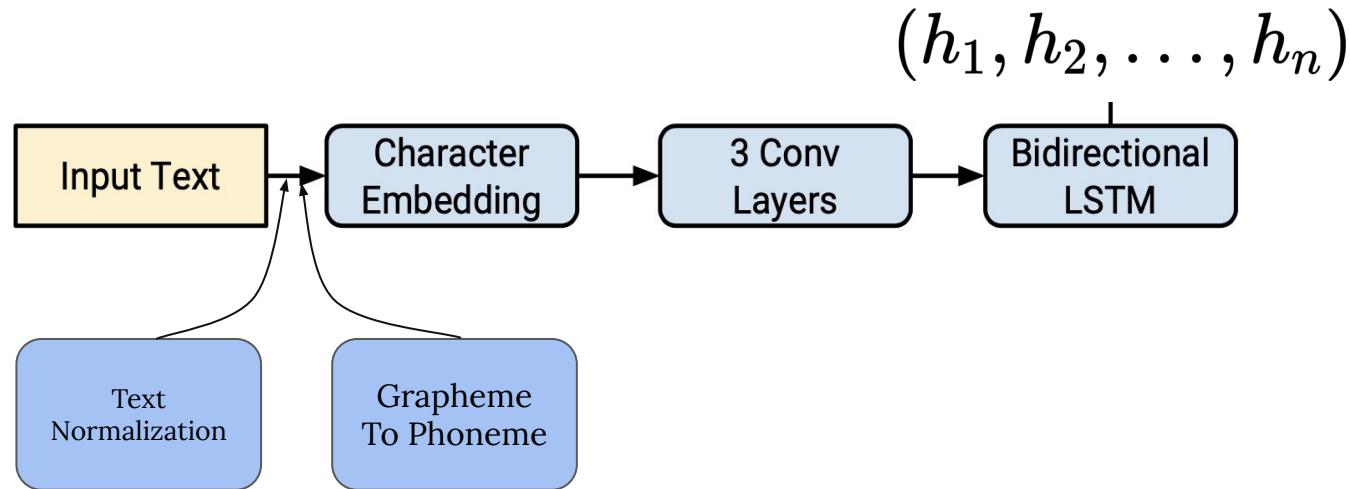
Wait! Why are we introducing MelSpectrogram?

- Train two component separately
 - Faster and easier
- MelSpectrogram is smoother than waveform
- Easier to train using a MSE

Tacotron 2



Text Encoder

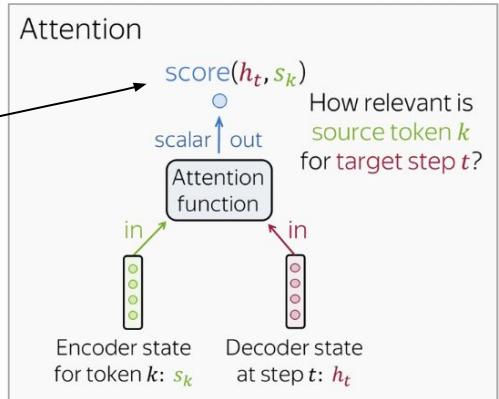


Attention

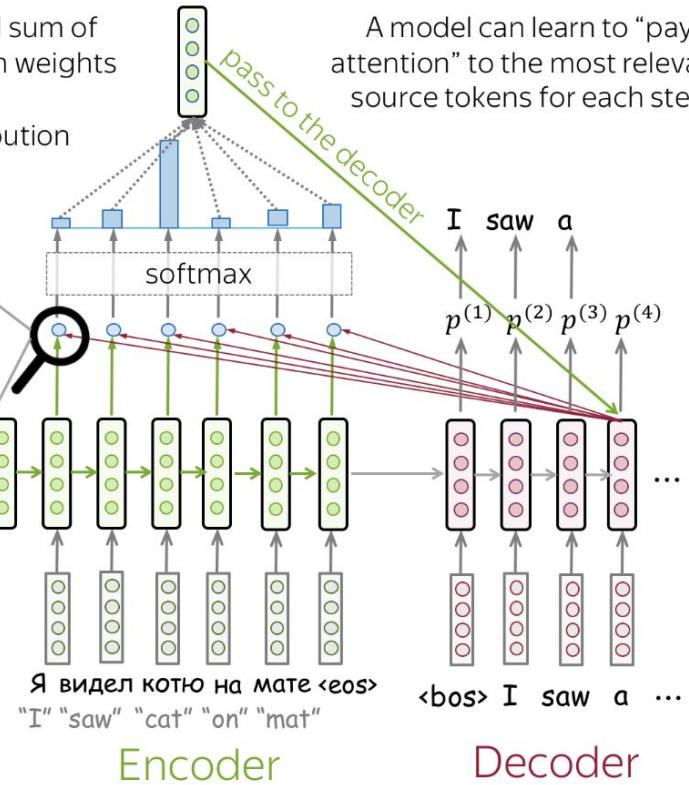
Attention output: weighted sum of encoder states with attention weights

Attention weights: distribution over source tokens

A model can learn to “pay attention” to the most relevant source tokens for each step



$$\text{score}_{t,k} = w^T \tanh(Ws_t + Vh_k + b)$$



Decoder

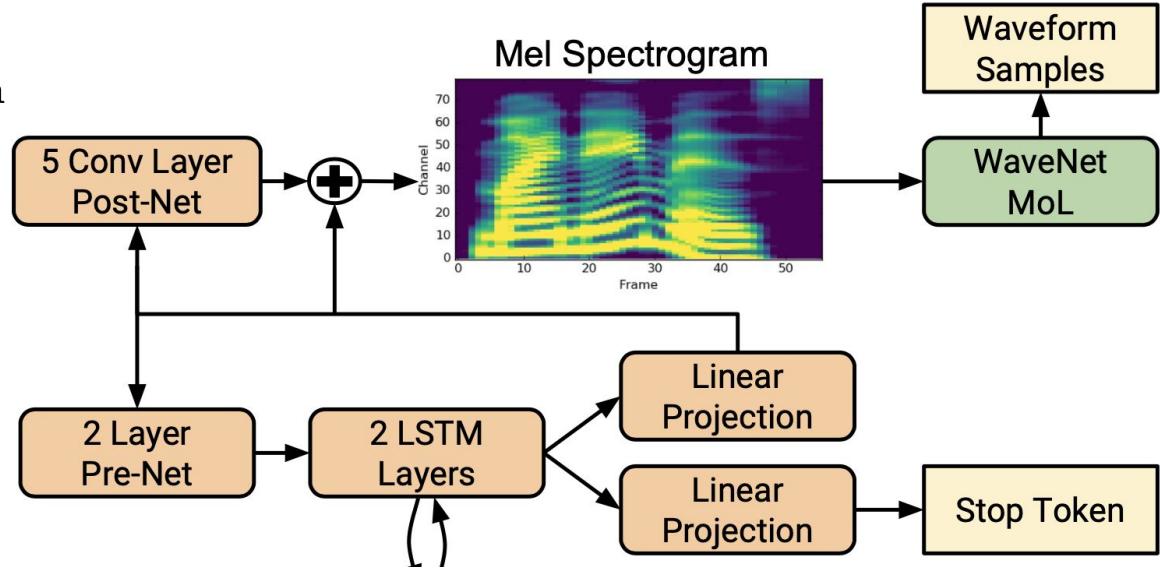
$$\mathcal{L} = \mathcal{L}_{\text{pre}} + \mathcal{L}_{\text{post}} + \text{StopToken}$$

$$\mathcal{L}_{\text{pre}} = \text{MSE}(x, \hat{x}_{\text{pre}})$$

$$\mathcal{L}_{\text{pre}} = \text{MSE}(x, \hat{x}_{\text{post}})$$

$$\text{StopToken} = \text{CE}(h, \mathbb{I}[h = \text{stop}])$$

- Neighboring predicted mels are correlated
- Turn on dropout on inference in PreNet

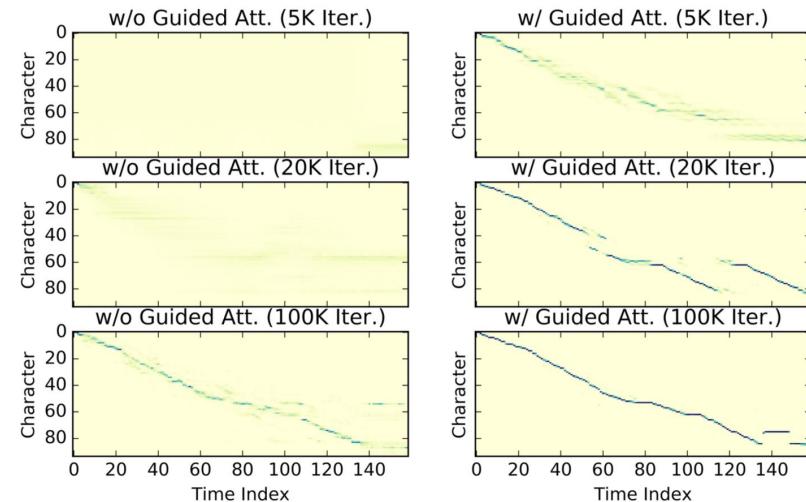


Guided Attention

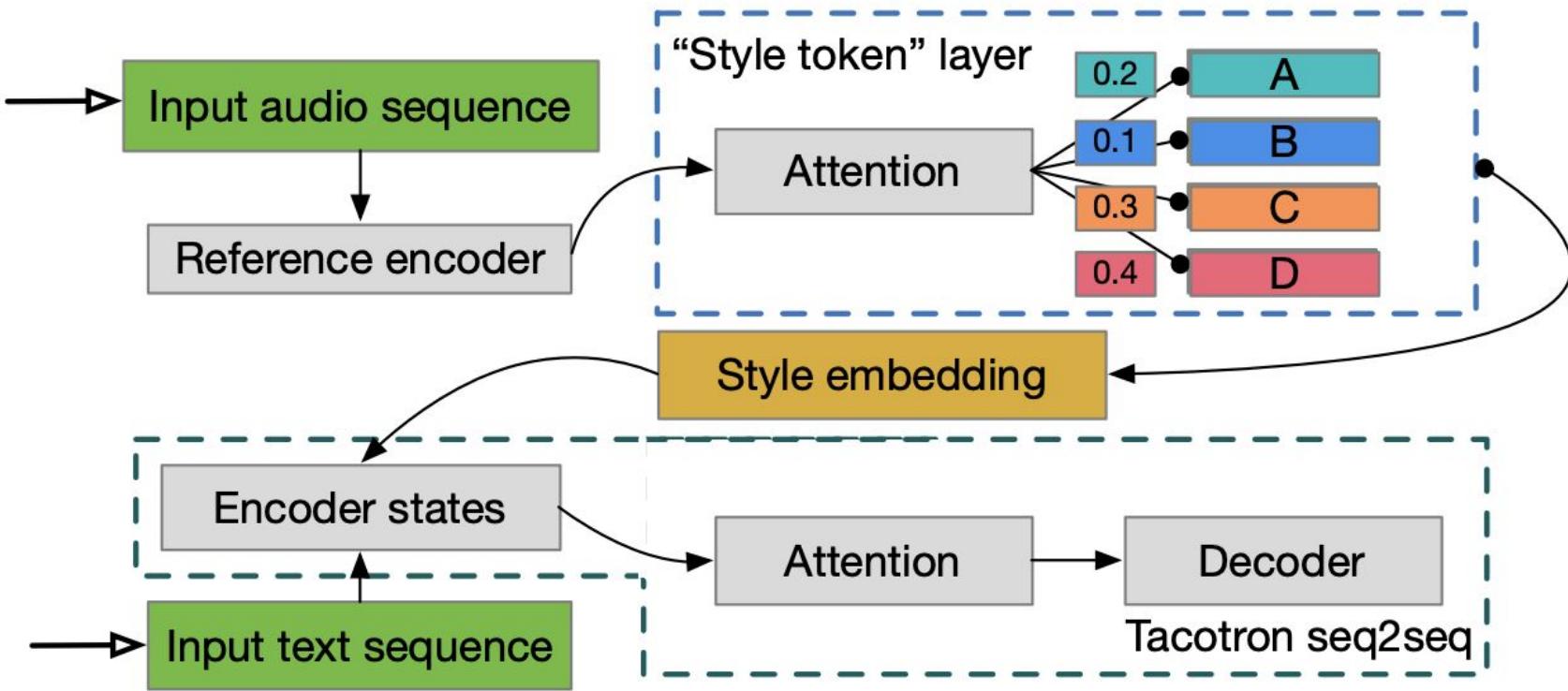
- Hard to learn attention from scratch
- Text position n progresses linearly to time t :

$$n \sim at, \text{ where } a \sim N/T$$

- Prompts attention matrix A to be “diagonal”
- $\mathcal{L}_{\text{att}}(A) = \mathbb{E}_{nt}[A_{nt}W_{nt}]$, where $W_{nt} = 1 - \exp\{-(n/N - t/T)^2/2g^2\}$



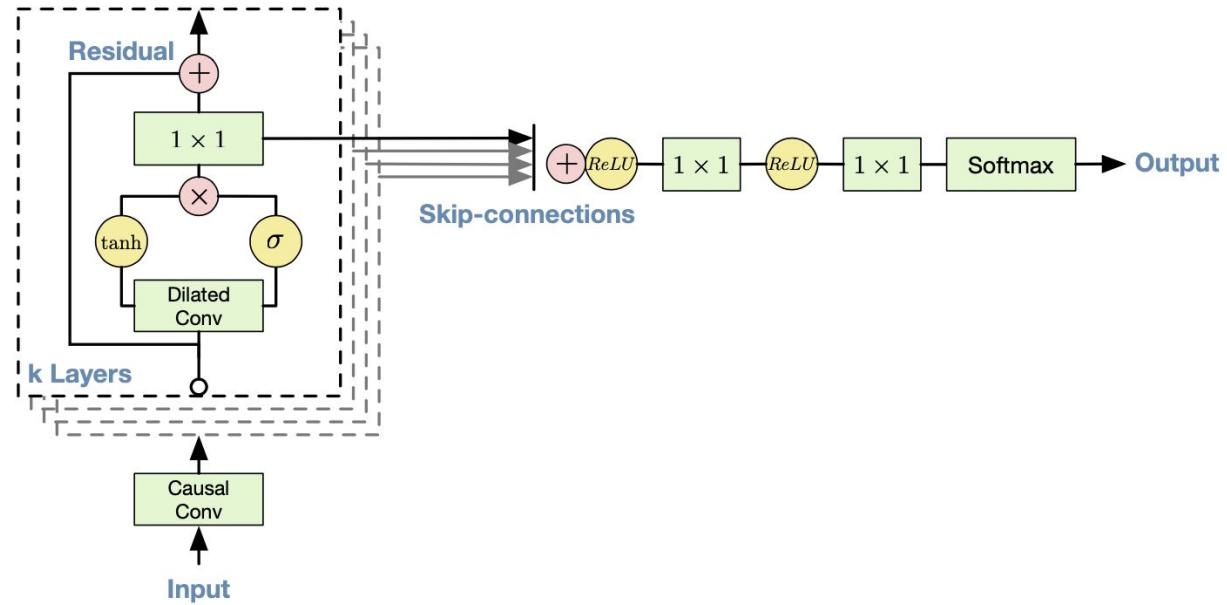
Global Style Token (GST)



WaveNet

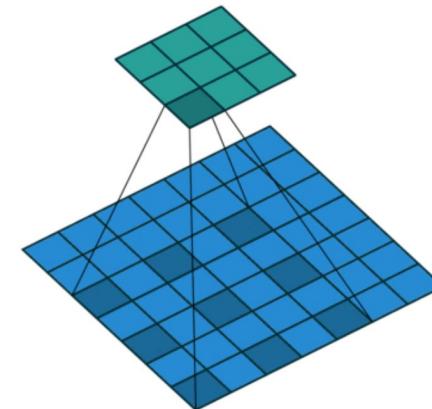
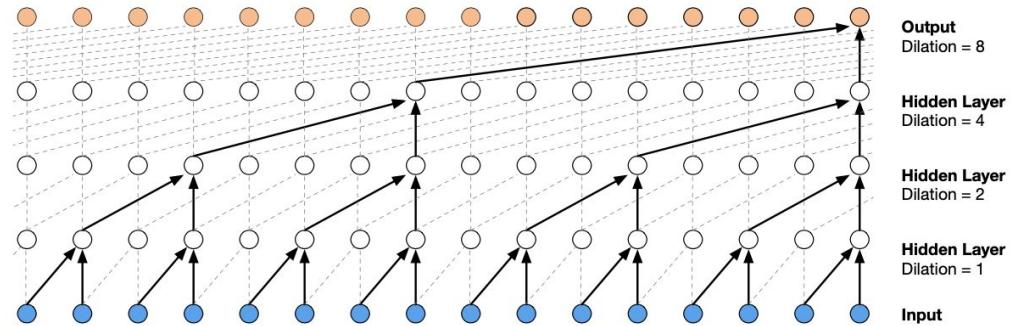
$$p(\mathbf{x}) = \prod_{t=1}^T p(x_t \mid x_1, \dots, x_{t-1})$$

$$p(x_t \mid x_i, \dots, x_{t-1}) \sim \text{Cat}(\pi_\theta)$$



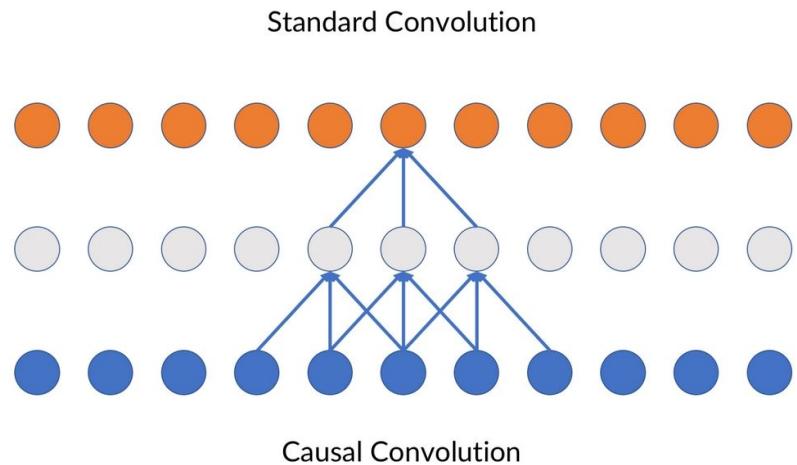
Dilated Convolution

- Increase receptive field
- Allow modeling long time dependencies



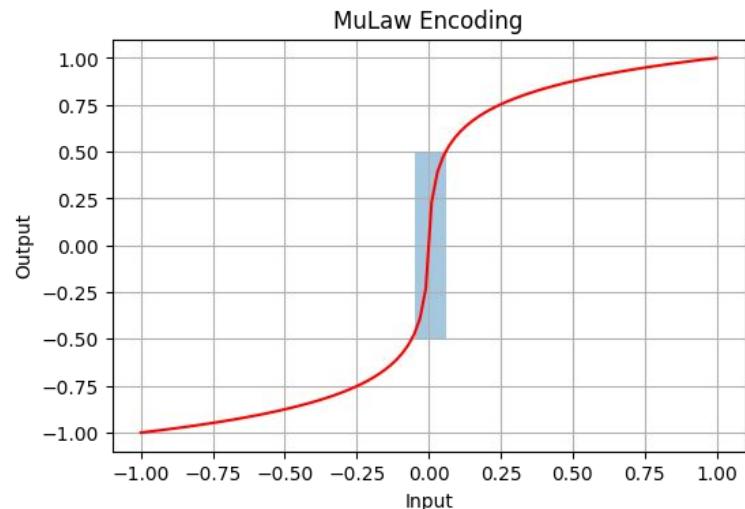
Causal Convolution

- $p(x_{t+1} \mid x_1, \dots, x_t)$
- Don't use padding in Conv
- Use a separate Pad



Mu Law Encoding

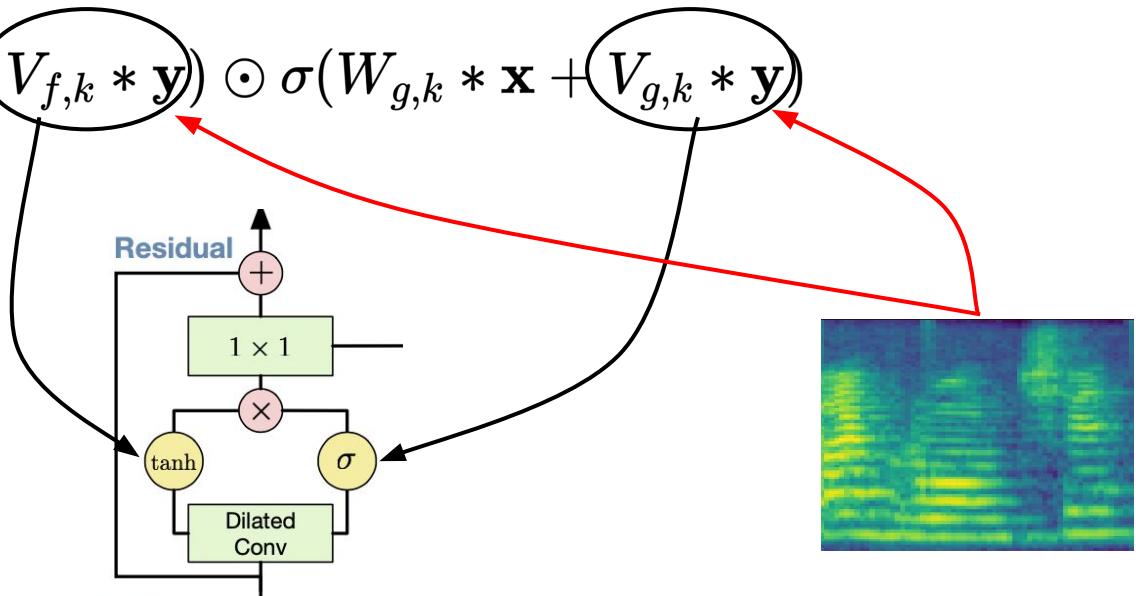
- $f(x_t) = \text{sign}(x_t) \frac{\ln(1 + \mu|x_t|)}{\ln(1 + \mu)}$
- 16-bit WAV contain 2^{16} values
- Softmax will die :(
- Human hearing on a **logarithmic** scale
- **Low-amplitude** sounds in **high** resolution
- **High-amplitude** sounds in **low** resolution



(Condition) Gated Mechanism

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

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

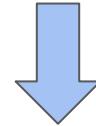


But how do we align the WAV and Mel?



But how do we align the WAV and Mel?

$$p(x_i | x_1, \dots, x_{i-1})$$

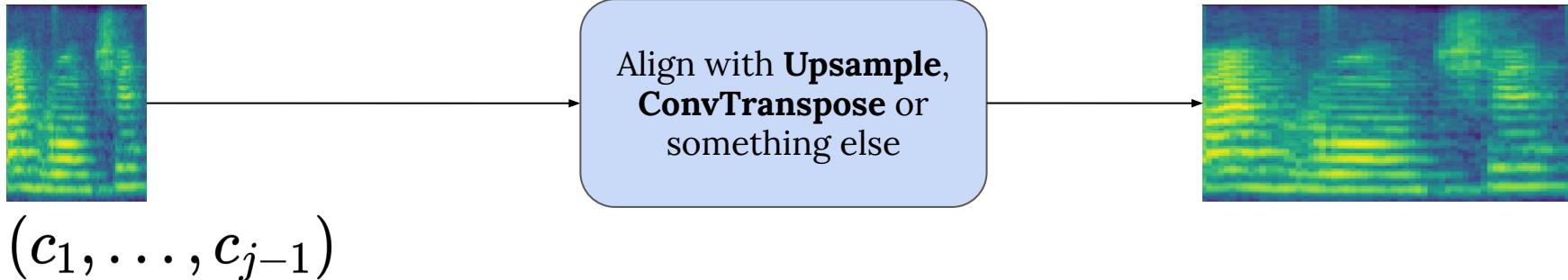
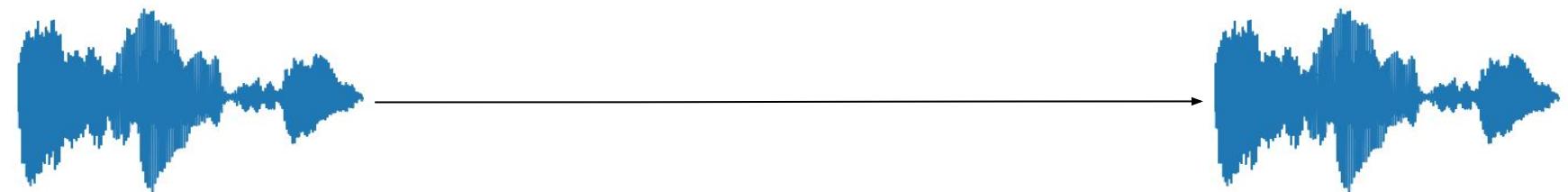


add a condition using **MelSpec**

$$p(x_i | x_1, \dots, x_{i-1}; c_1, \dots, c_{j-1})$$

Samples and MelSpec are **not**
time-aligned!

But how do we align the WAV and Mel?



Upsample is our everything!



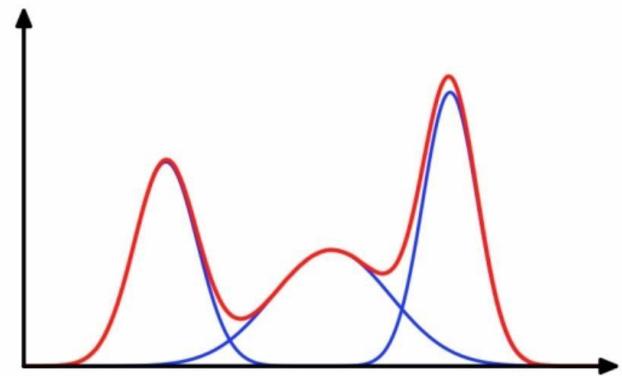
But how do we align the WAV and Mel?

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

The diagram illustrates the alignment process between a waveform and a spectrogram. On the left, a vertical blue bar represents the waveform segment x_i . In the center, a blue waveform plot shows the sequence (x_1, \dots, x_{i-1}) . On the right, a spectrogram plot shows the sequence (c_1, \dots, c_{i-1}) . Blue arrows point from the waveform segments to the first term of the sequence. Green arrows point from the spectrogram segments to the second term of the sequence. The mathematical formula at the top describes the computation of a latent variable \mathbf{z} based on the concatenated features \mathbf{x} and \mathbf{y} from both the waveform and spectrogram paths.

What about a loss function?

- Categorical distribution
- Normal distribution
- Logistic distribution
- Mixture of Normals or Logistics
- Use `torch.distributions :)`

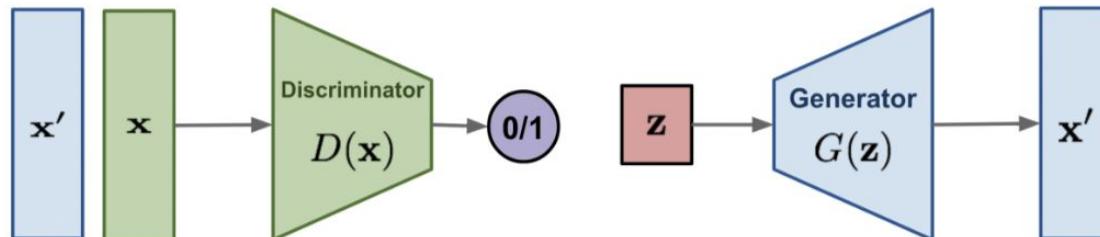


What if... we just learn to map Mels to WAVs?

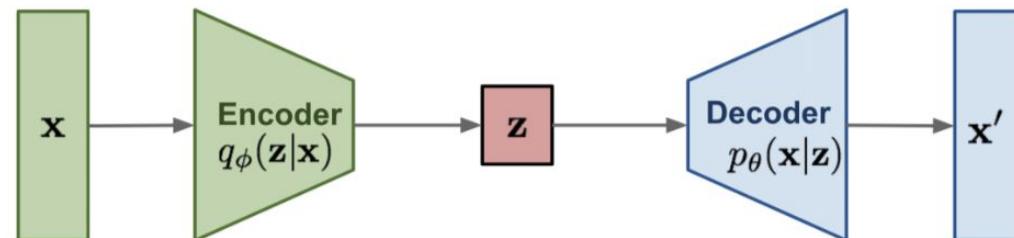


Generative models

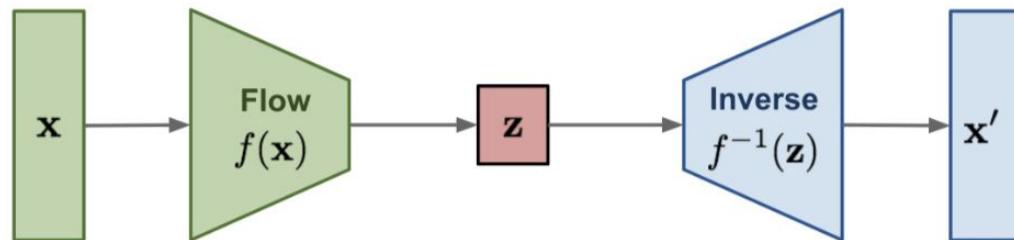
GAN: minimize the classification error loss.



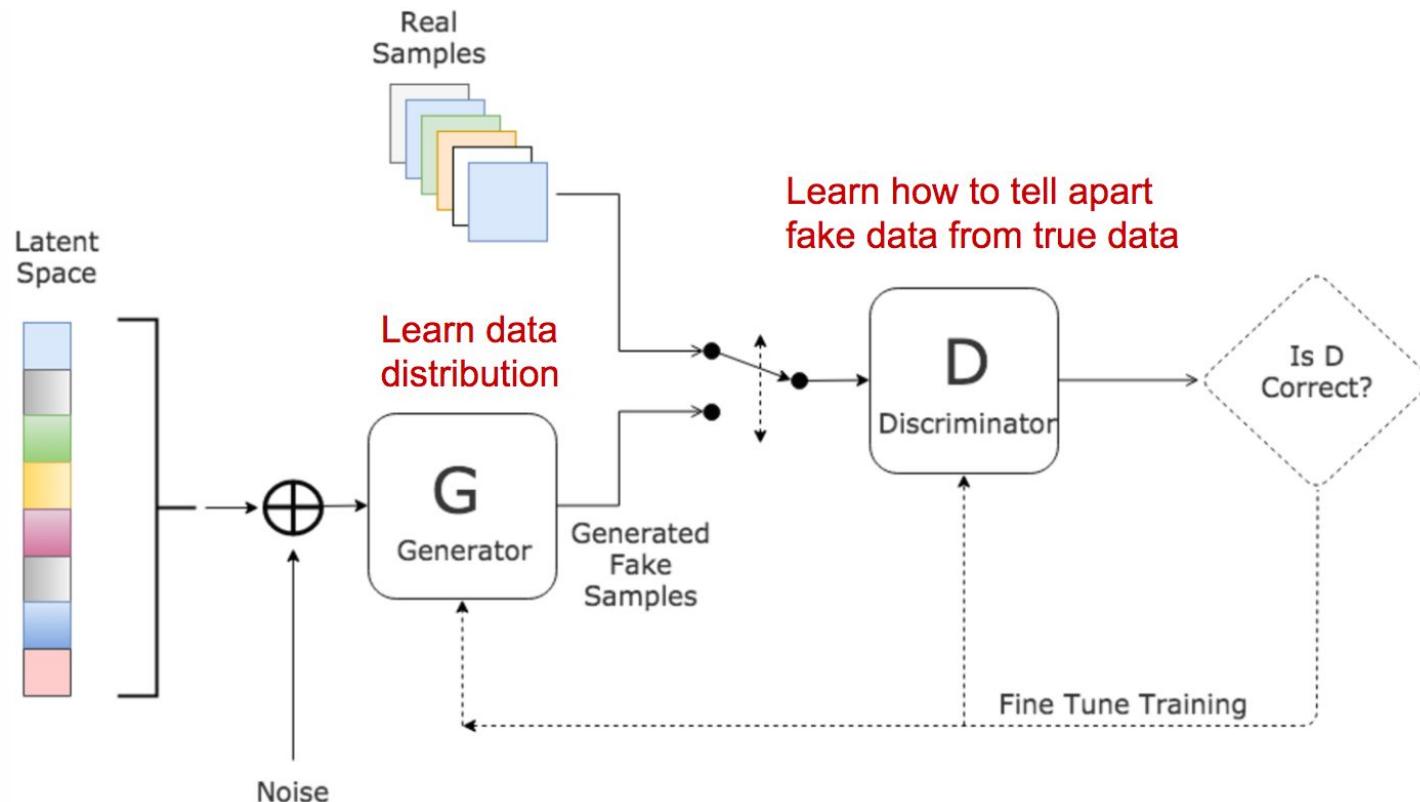
VAE: maximize ELBO.



Flow-based generative models: minimize the negative log-likelihood



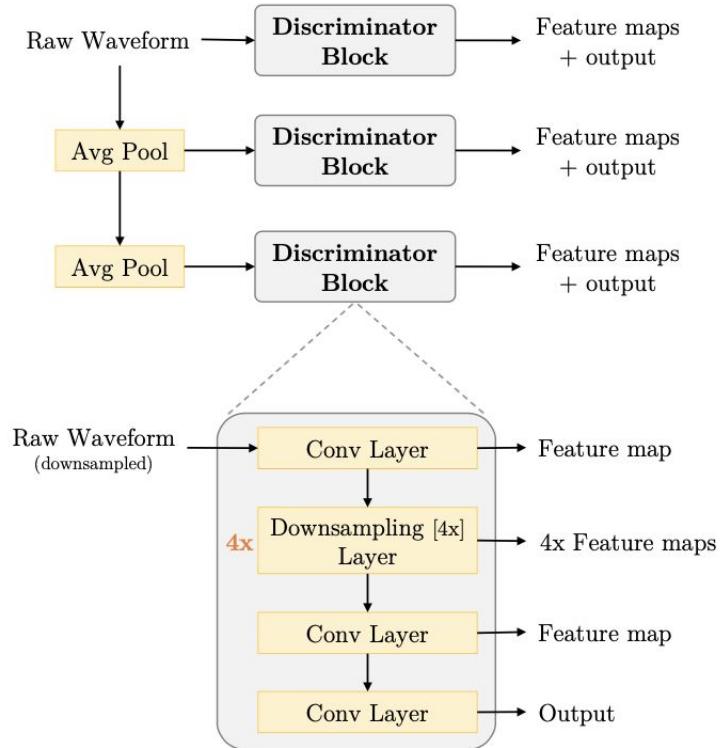
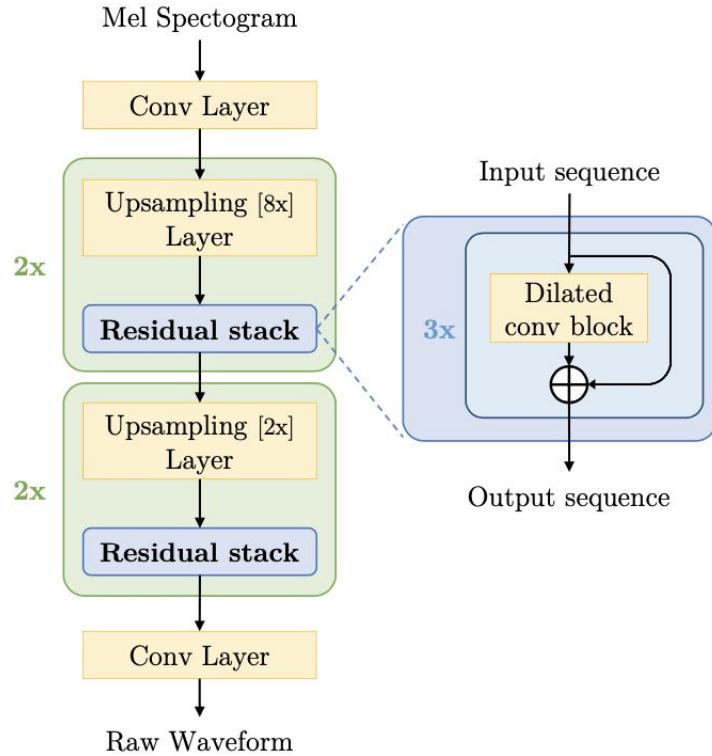
Generative Adversarial Networks



Generative Adversarial Networks

- A discriminator **D** estimates the probability of a given sample coming from the real or synthetic dataset
- A generator **G** outputs synthetic samples given a noise variable input **z**
- $\min_G \max_D L(D, G) = \mathbb{E}_{x \sim p_r(x)}[\log D(x)] + \mathbb{E}_{z \sim p_z(z)}[\log(1 - D(G(z)))]$
 $= \mathbb{E}_{x \sim p_r(x)}[\log D(x)] + \mathbb{E}_{x \sim p_g(x)}[\log(1 - D(x))]$

MelGAN



MelGAN

- Multiscale Discriminators
- Weight Normalization
- Markovian Discriminator
- Hinge Loss
- Feature Matching

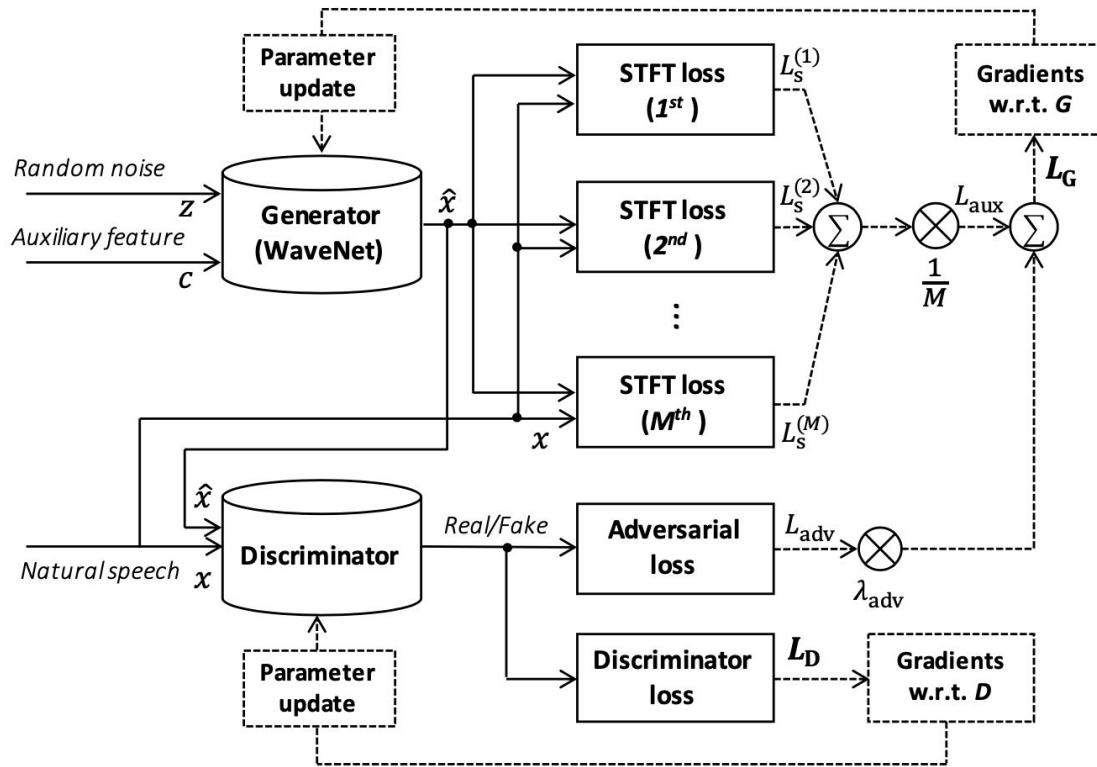
$$\mathbf{w} = \frac{g}{\|\mathbf{v}\|} \mathbf{v}$$

$$\mathcal{L}_{\text{FM}}(G, D_k) = \mathbb{E}_{x,s \sim p_{\text{data}}} \left[\sum_{i=1}^T \frac{1}{N_i} \left\| D_k^{(i)}(x) - D_k^{(i)}(G(s)) \right\|_1 \right]$$

$$\begin{aligned} & \min_{D_k} \mathbb{E}_x [\min(0, 1 - D_k(x))] + \mathbb{E}_{s,z} [\min(0, 1 + D_k(G(s, z)))] , \forall k = 1, 2, 3 \\ & \min_G \mathbb{E}_{s,z} \left[\sum_{k=1,2,3} -D_k(G(s, z)) \right] \end{aligned}$$

Parallel WaveGAN

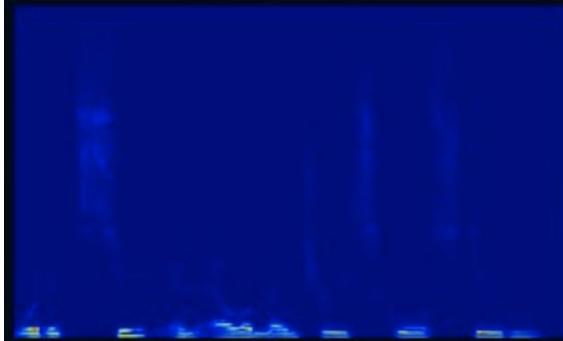
- Use WaveNet based Generator
- Additionally use multi-STFT loss
- LSGAN



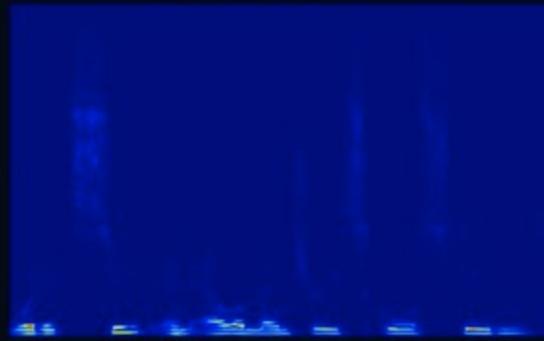
STFT LOSS

Spectral Convergence part

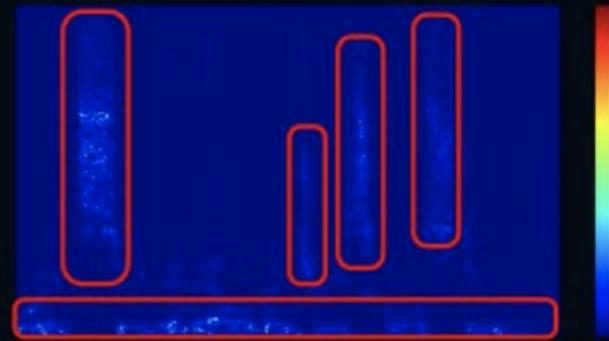
$|\text{STFT}(x)|$



$|\text{STFT}(\hat{x})|$



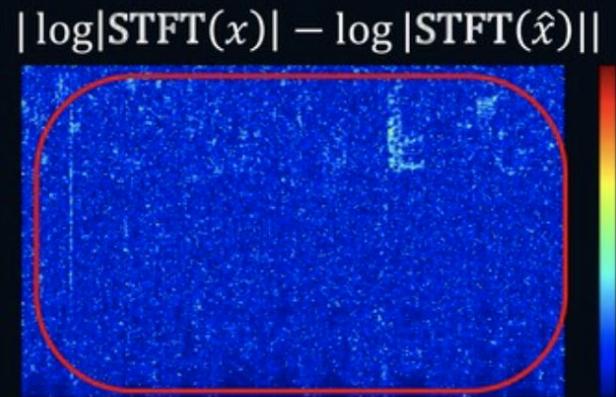
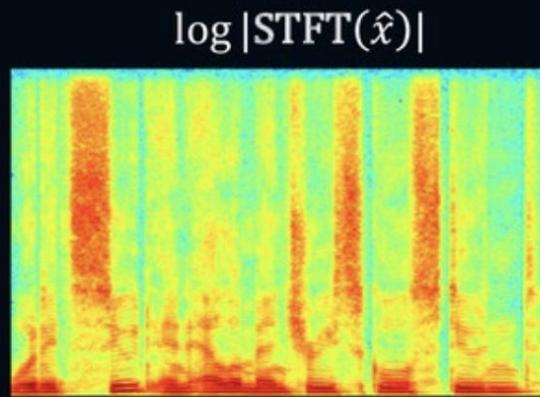
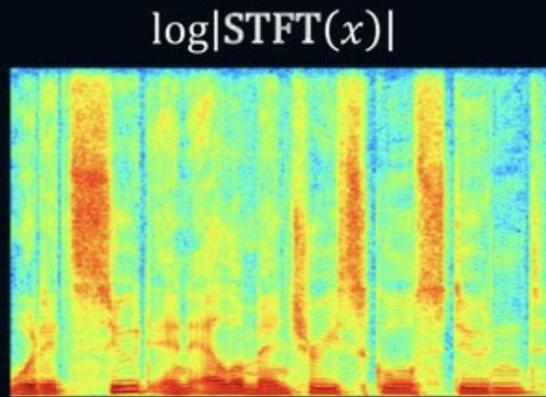
$\|\|\text{STFT}(x)| - |\text{STFT}(\hat{x})\|\|$



$$L_{sc} = \frac{\|\|\text{STFT}(x)| - |\text{STFT}(\hat{x})\|\|_F}{\|\|\text{STFT}(x)\|\|_F}$$

STFT LOSS

Log scale STFT magnitude part

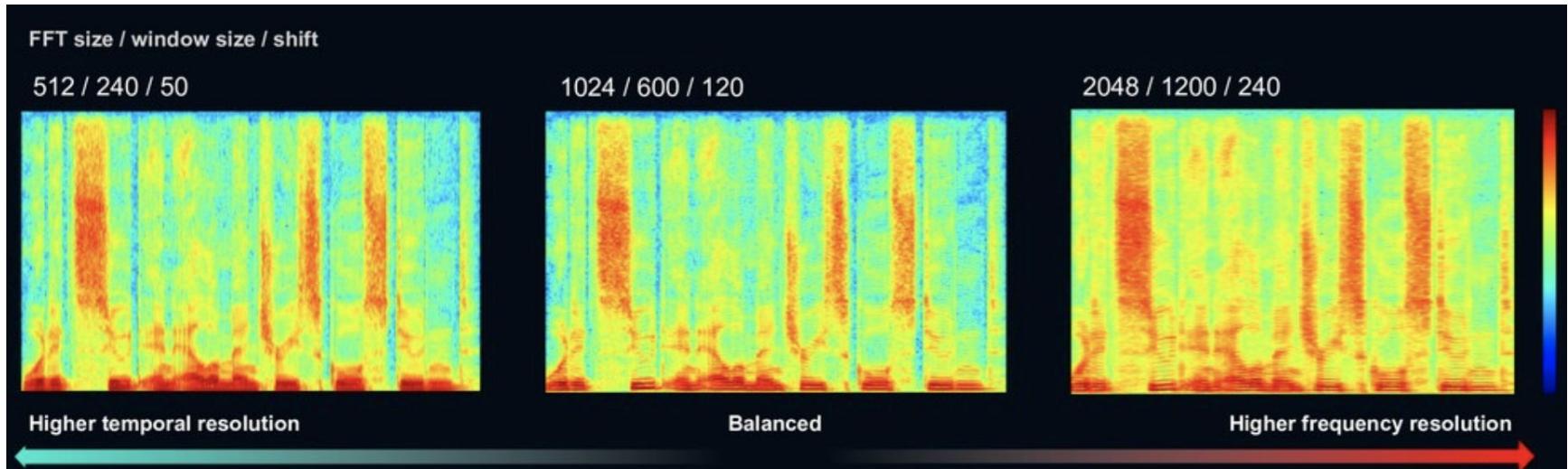


$$L_{\text{mag}} = \frac{1}{N} \|\log|\text{STFT}(x)| - \log |\text{STFT}(\hat{x})||\|_1$$

N : number of elements in the STFT magnitude

STFT LOSS

Multi Resolution



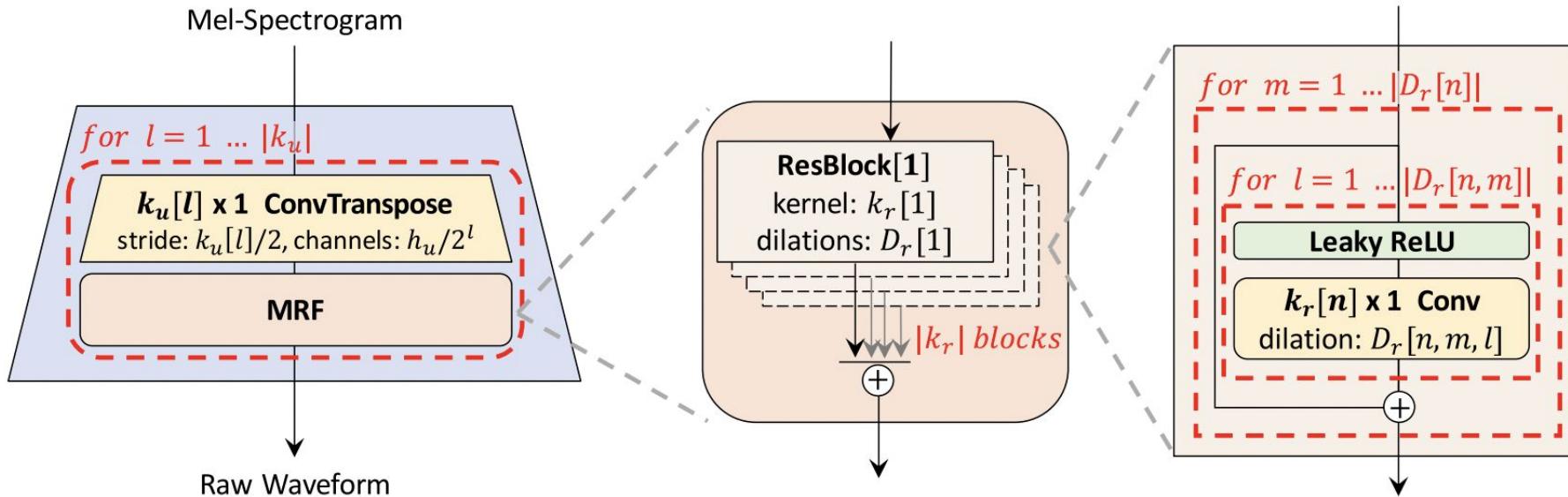
$$L_{\text{aux}}(G) = \frac{1}{M} \sum_{m=1}^M L_s^{(m)}(G)$$

$$L_s(G) = \mathbb{E}_{z \sim p(z), x \sim p_{\text{data}}} [L_{\text{sc}}(x, \hat{x}) + L_{\text{mag}}(x, \hat{x})]$$

M : number of STFT losses

HiFi GAN

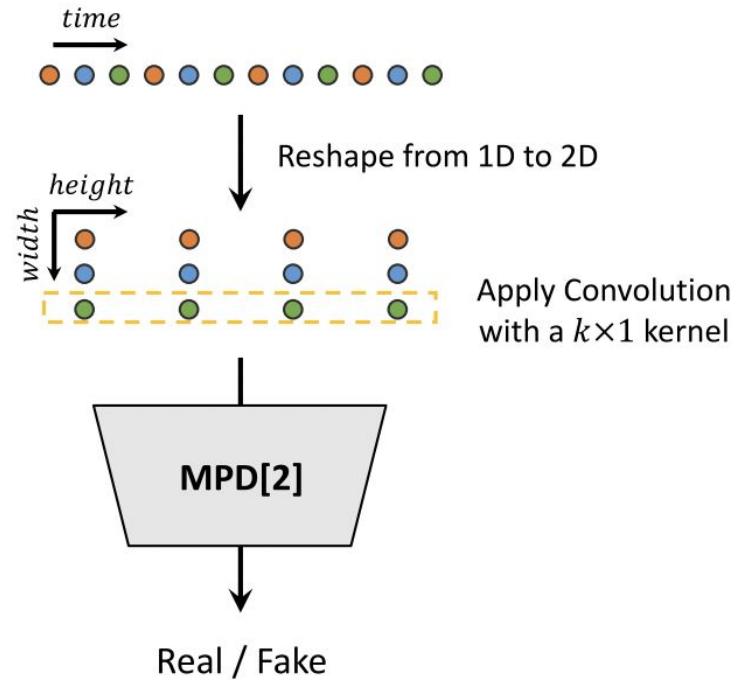
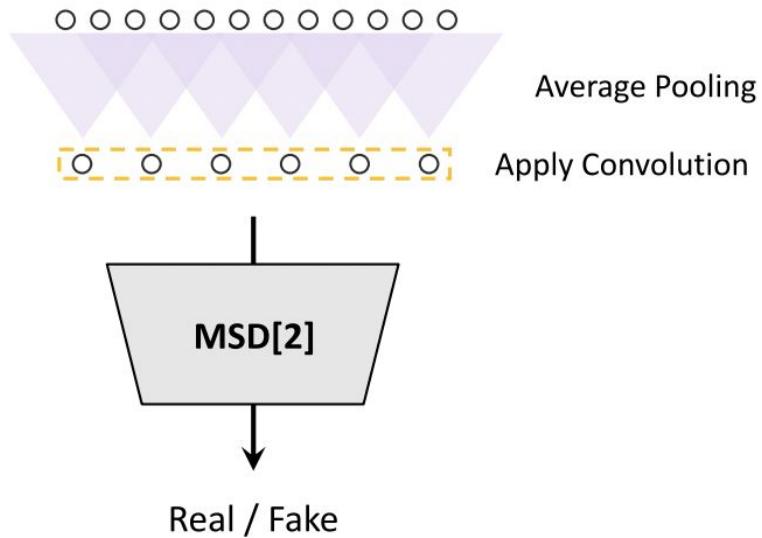
Generator



h_u	k_u	k_r	D_r
512	[16, 16, 4, 4]	[3, 7, 11]	[[1, 1], [3, 1], [5, 1]] $\times 3$

HiFi GAN

Discriminator



HiFi GAN

Criterions

$$\begin{aligned}\mathcal{L}_{Adv}(D; G) &= \mathbb{E}_{(x,s)}[(D(x) - 1)^2 + (D(G(s)))^2] \\ \mathcal{L}_{Adv}(G; D) &= \mathbb{E}_s[(D(G(s))) - 1]^2\end{aligned}$$

$$\mathcal{L}_{Mel}(G) = \mathbb{E}_{(x,s)}[\|\phi(x) - \phi(G(s))\|_1]$$

$$\mathcal{L}_{FM}(G; D) = \mathbb{E}_{(x,s)} \left[\sum_{i=1}^T \frac{1}{N_i} \|D^i(x) - D^i(G(s))\|_1 \right]$$

Flow-based models

Basic idea:

- ▶ Use an invertible transformation $z = f_\theta(x)$ to transform x into z ;
- ▶ Suppose that z has a simple distribution $p(z)$ (e.g. standard normal);
- ▶ Use maximum likelihood to fit θ .

$$\sum_{i=1}^n \log p_\theta(x_i) = \sum_{i=1}^n \log p(z_i) + \log \left| \det \left(\frac{dz_i}{dx_i} \right) \right| \rightarrow \max_{\theta}.$$

- ▶ $p(z)$ is easy to compute. But what about log-abs-det-Jacobian?

Flow-based models

- ▶ We suppose that $f_\theta = f^{(L)} \circ f^{(L-1)} \circ \dots \circ f^{(1)}$.
- ▶ If each of $f^{(j)}$ is invertible, then

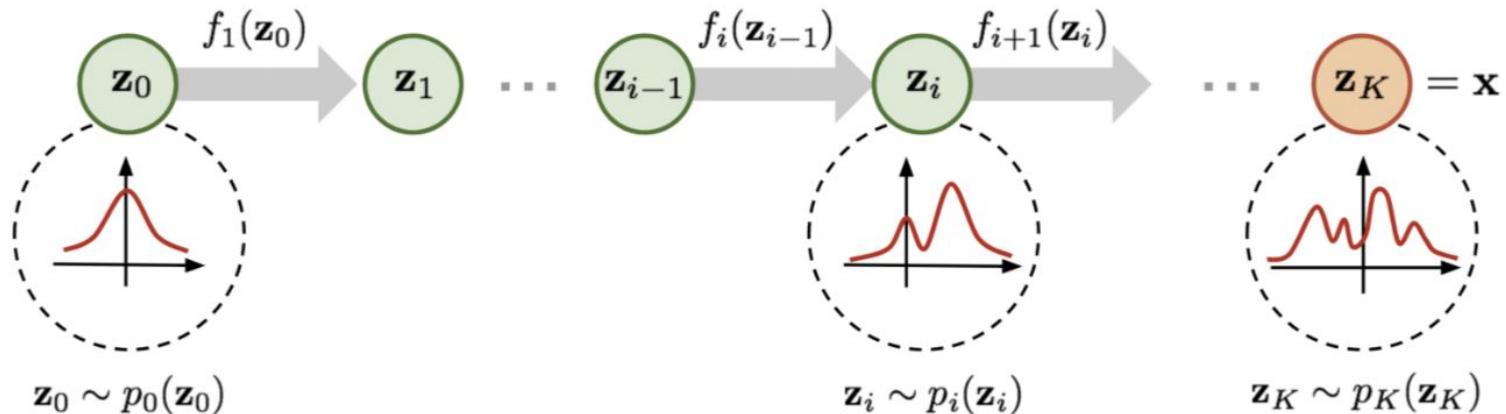
$$\log \left| \det \left(\frac{dz}{dx} \right) \right| = \sum_{j=1}^L \log \left| \det \left(\frac{dh_j}{dh_{j-1}} \right) \right|,$$

where $h_j = f_j(x)$, $h_0 = x$, and $h_L = z$.

- ▶ We now have to come up with invertible parametric transformations such that they have easily computable log-abs-det-Jacobian and their composition is expressive enough.
- ▶ STACK MORE LAYERS.

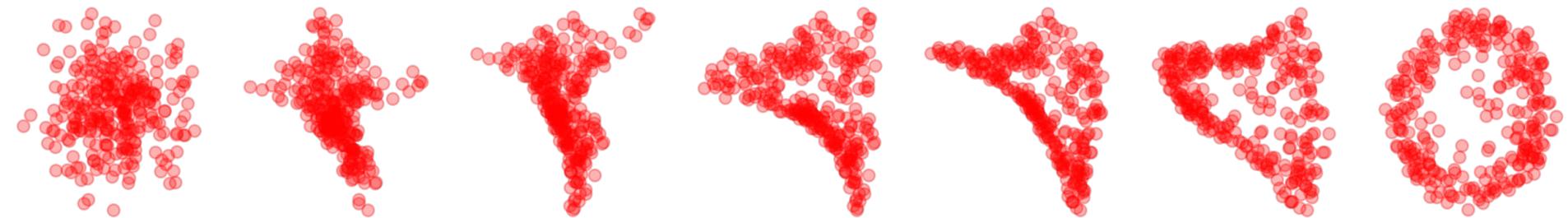
Flow-based models

- ▶ The resulting transformation f is called a *normalizing flow* or just *flow*.
- ▶ We can sample from the model just by applying the inverse transform $x = f_{\theta}^{-1}(z)$:



Flow-based models

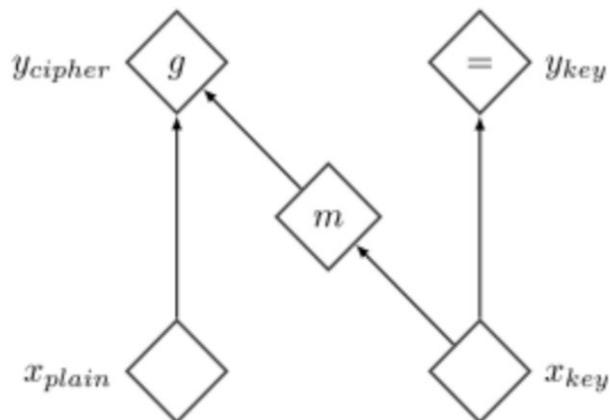
$$z \sim p(z)$$
$$x = f^{-1}(z)$$



Additive coupling

$$y_{I_1} = x_{I_1}$$

$$y_{I_2} = g(x_{I_2}, m(x_{I_1})) \cdot$$



$$\frac{\partial y}{\partial x} = \begin{bmatrix} I_d & 0 \\ \frac{\partial y_{I_2}}{\partial x_{I_1}} & \frac{\partial y_{I_2}}{\partial x_{I_2}} \end{bmatrix}$$

$$\det \left(\frac{\partial y}{\partial x} \right) = \prod_{i=1}^n \left(\frac{\partial y_{I_2}}{\partial x_{I_2}} \right)_{ii}$$

Additive coupling

- ▶ [Dinh et al. (2016)] A simple extension to the previous approach allows the transform to have non-trivial Jacobian:

$$y_{1:d} = x_{1:d},$$

$$y_{d+1:D} = \exp(s(x_{1:d})) \odot x_{d+1:D} + t(x_{1:d}).$$

The Log-abs-det-Jacobian is $\text{sum}(s(x_{1:d}))$ and inverse is also easy to compute.

- ▶ In case of image data, it makes sense to use CNNs for s and t and split dimensions over channels.

Autoregressive Normalizing Flows

Masked Autoregressive Flow

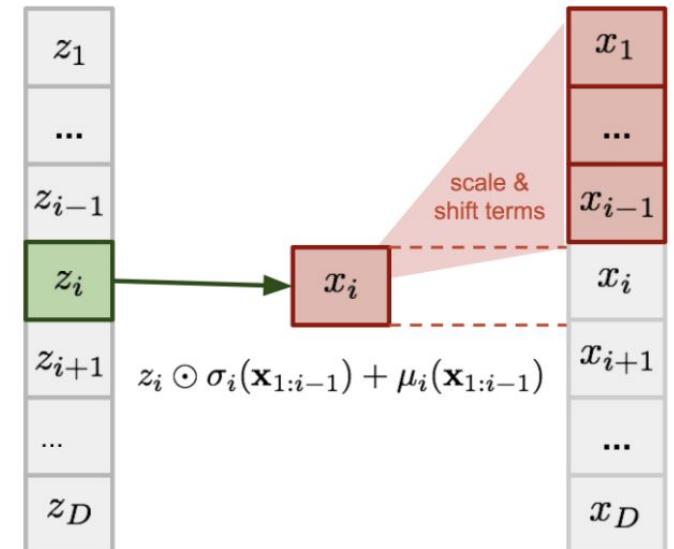
$$z_i \sim \mathcal{N}(0, 1)$$

$$\mu_i = f_{\mu_i}(\mathbf{x}_{1:i-1})$$

$$\alpha_i = f_{\alpha_i}(\mathbf{x}_{1:i-1})$$

$$x_i = z_i \exp \alpha_i + \mu_i$$

$$\left| \det \left(\frac{\partial f^{-1}}{\partial \mathbf{x}} \right) \right| = \exp \left(- \sum_i \alpha_i \right)$$



$\mathbf{z} \sim \pi(\mathbf{z}) \xrightarrow{\text{(known)}} ? \xrightarrow{\text{(unknown)}} \mathbf{x} \sim p(\mathbf{x})$

Autoregressive Normalizing Flows

Inverse Autoregressive Flow

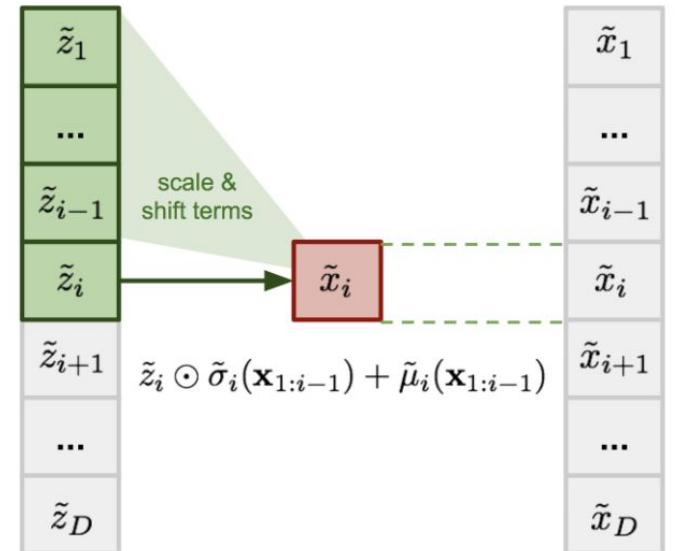
$$z_i \sim \mathcal{N}(0, 1)$$

$$\mu_i = f_{\mu_i}(\mathbf{z}_{1:i-1})$$

$$\alpha_i = f_{\alpha_i}(\mathbf{z}_{1:i-1})$$

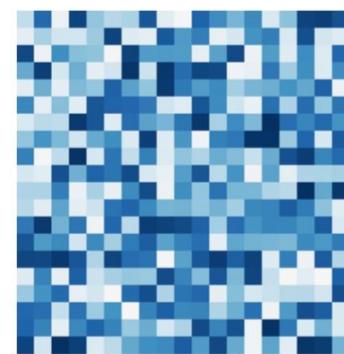
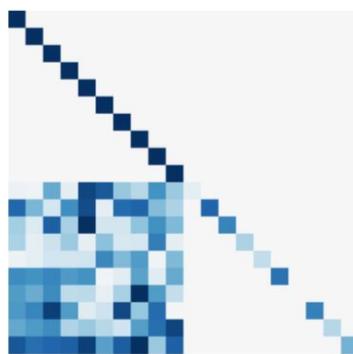
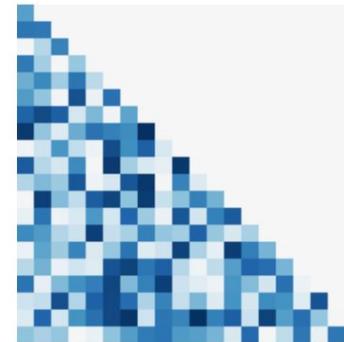
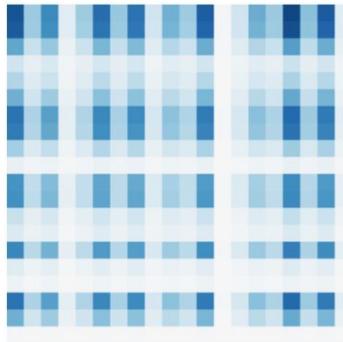
$$x_i = z_i \exp \alpha_i + \mu_i$$

$$\left| \det \left(\frac{\partial f^{-1}}{\partial \mathbf{x}} \right) \right| = \exp \left(- \sum_i \alpha_i \right)$$



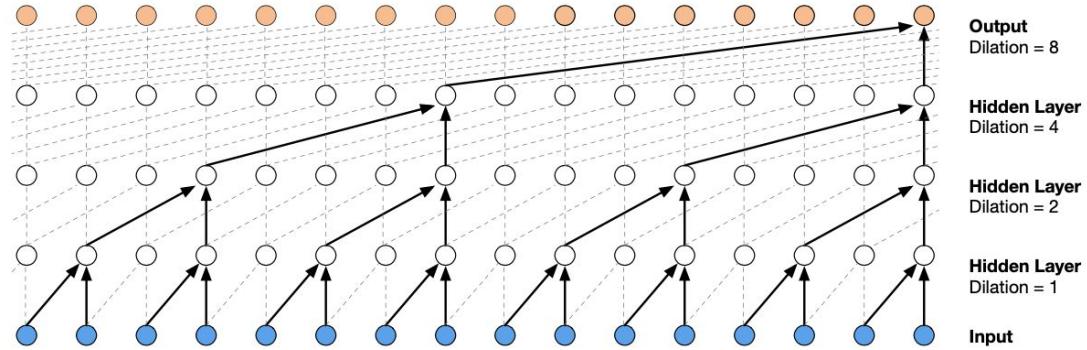
$$\tilde{\mathbf{z}} \sim \tilde{\pi}(\tilde{\mathbf{z}}) \xrightarrow[\text{(known)}]{} ? \xrightarrow{} \tilde{\mathbf{x}} \sim \tilde{p}(\tilde{\mathbf{x}}) \quad \text{(unknown)}$$

Jacobians



Recap WaveNet

- Dilated Convolution
- Casual Convolution
- Mu Law Encoding
- Gated Mechanism
- Global & Local Conditioning



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

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

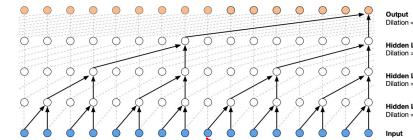
$$\mathbf{z} = \tanh(W_{f,k} * \mathbf{x} + V_{f,k} * \mathbf{y}) \odot \sigma(W_{g,k} * \mathbf{x} + V_{g,k} * \mathbf{y})$$

Parallel WaveNet

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

$$x_t = z_t \cdot s(z_{<t}, \theta) + \mu(z_{<t}, \theta)$$

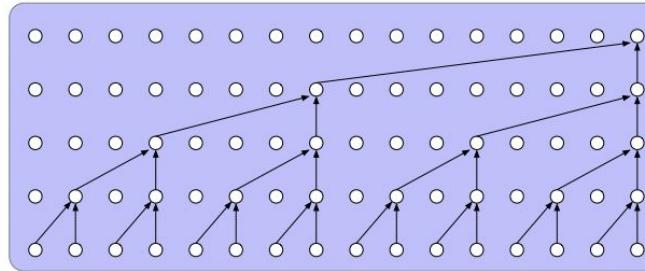
$$p(x_t \mid z_{<t}, \theta) = \mathbb{L}(x_t \mid \mu(z_{<t}, \theta), s(z_{<t}, \theta))$$



Parallel WaveNet

WaveNet Teacher

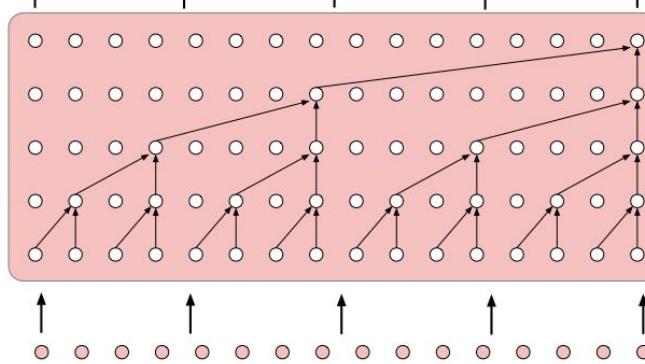
Linguistic features



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

WaveNet Student

Linguistic features



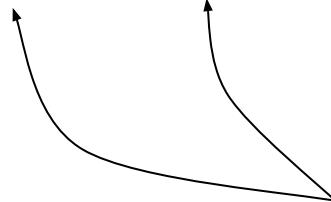
Generated Samples
 $x_i = g(z_i|z_{<i})$

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

Input noise
 z_i

Parallel WaveNet

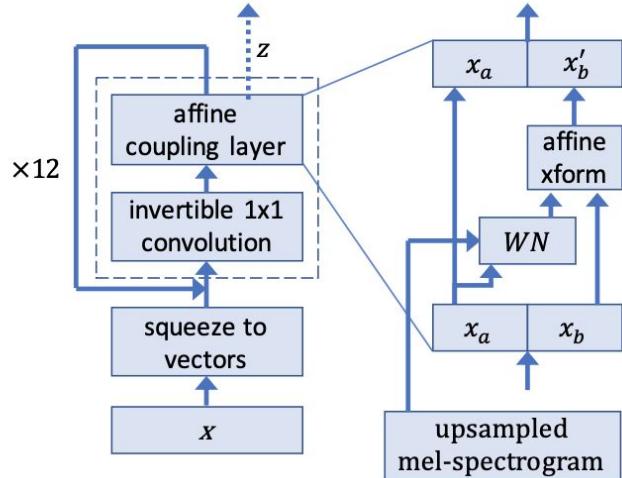
- Probability Density Distillation $D_{\text{KL}}(P_S \| P_T) = H(P_S, P_T) - H(P_S)$
- Power Loss $\|\phi(g(\mathbf{z}, \mathbf{c})) - \phi(\mathbf{y})\|^2$
- Perceptual Loss
- Contrastive Loss $D_{\text{KL}}(P_S(\mathbf{c}_1) \| P_T(\mathbf{c}_1)) - \gamma D_{\text{KL}}(P_S(\mathbf{c}_1) \| P_T(\mathbf{c}_2))$



Linguistic features, speaker ID

WaveGlow

- Combine insights from Glow and WaveNet
- Squeeze Operation
- Affine Coupling Layer
- 1x1 Invertible Convolution
- Early outputs



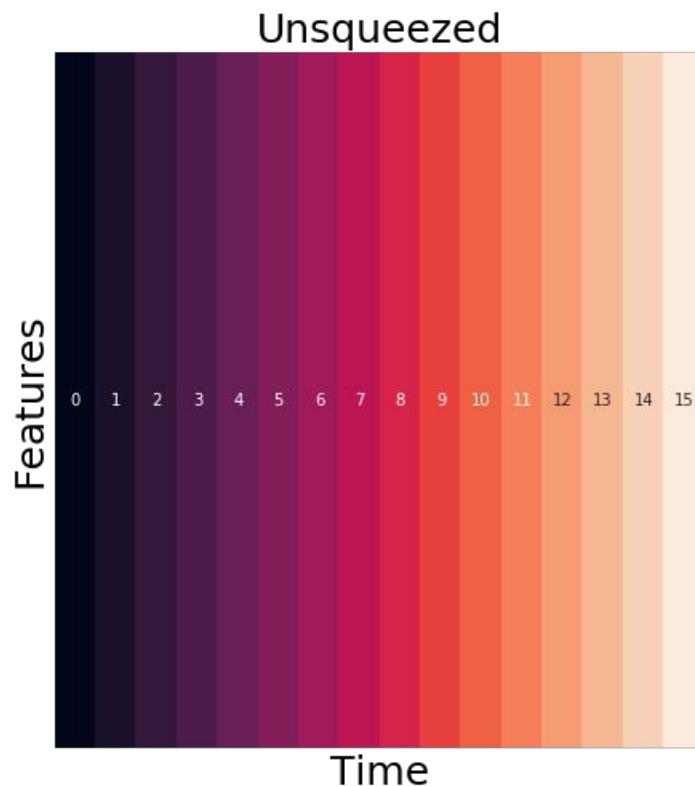
$$\mathbf{x}_a, \mathbf{x}_b = \text{split}(\mathbf{x})$$

$$(\log \mathbf{s}, \mathbf{t}) = \text{WaveNet}(\mathbf{x}_a, \text{MelSpectrogram})$$

$$\mathbf{x}_{b'} = \mathbf{s} \odot \mathbf{x}_b + \mathbf{t}$$

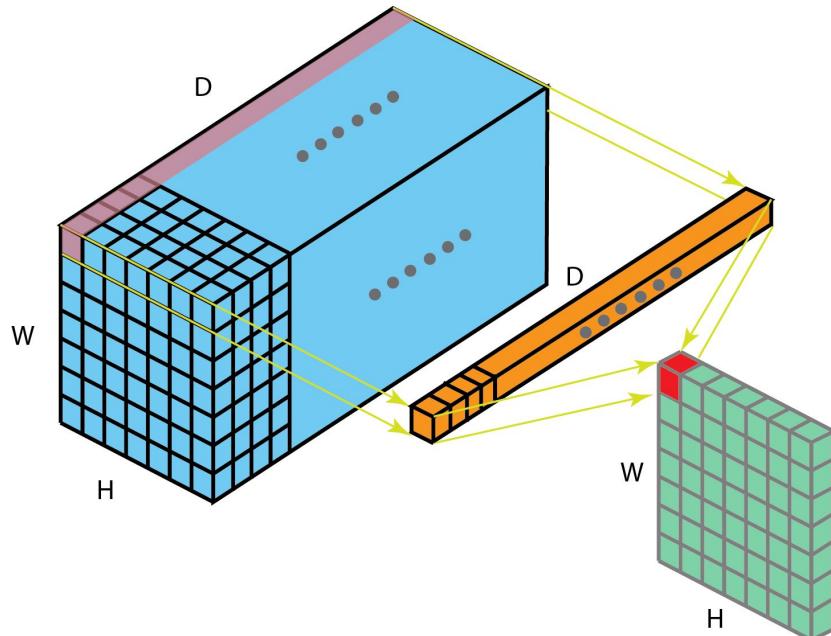
$$\mathbf{f}_{\text{coupling}}^{-1}(\mathbf{x}) = \text{concat}(\mathbf{x}_a, \mathbf{x}_{b'})$$

Squeeze Operation



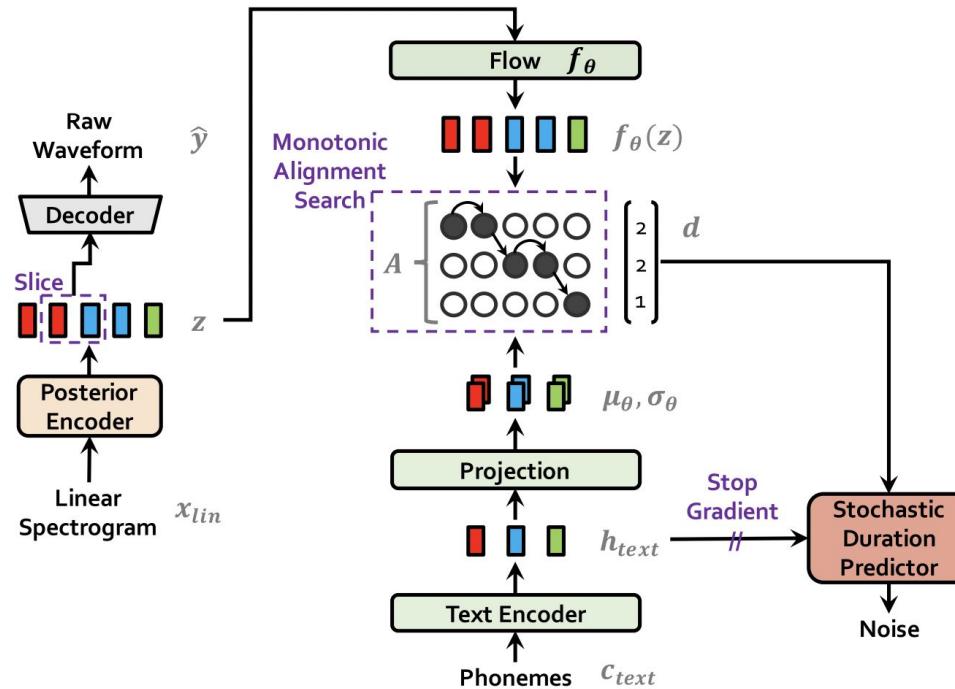
1x1 Invertible Convolution

- We want to permute **squeezed** channels
- Initialize weights as random **rotation** matrix



VITS

Training



$$L_{vae} = L_{recon} + L_{kl} + L_{dur} + L_{adv}(G) + L_{fm}(G)$$

VITS

Inference

