

# Diffusion-based Vocoding for Real-Time Text-To-Speech

Lukas Gardberg

LTH

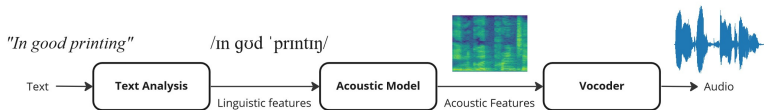
March 21, 2023

# Presentation Layout

- ▶ Problem Introduction
- ▶ Past Work
- ▶ Problem Statement / Goals
- ▶ Quick background
- ▶ TODO: Add more

# Typical TTS Pipeline

Text  $\longrightarrow$  Speech



Text Analysis

Acoustic Model

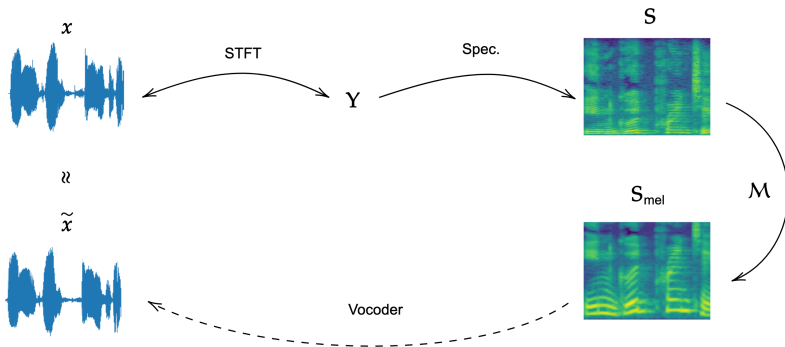
Vocoder

# The Phase Reconstruction Problem

- ▶ Goal: Reconstruct signal  $\mathbf{x}$  from its mel spectrogram  $\mathbf{S}_{\text{mel}}$

# The Phase Reconstruction Problem

- Goal: Reconstruct signal  $x$  from its mel spectrogram  $S_{\text{mel}}$



# The Phase Reconstruction Problem

How has this been done before?

# The Phase Reconstruction Problem

How has this been done before?

- ▶ Griffin-Lim Reconstruction

# The Phase Reconstruction Problem

How has this been done before?

- ▶ Griffin-Lim Reconstruction
- ▶ Autoregressive Neural Networks (WaveNet)



# The Phase Reconstruction Problem

How has this been done before?

- ▶ Griffin-Lim Reconstruction
- ▶ Autoregressive Neural Networks (WaveNet)
- ▶ GANs (HiFi-GAN)

# The Phase Reconstruction Problem

How has this been done before?

- ▶ Griffin-Lim Reconstruction
- ▶ Autoregressive Neural Networks (WaveNet)
- ▶ GANs (HiFi-GAN)
- ▶ Diffusion (DiffWave)

# The Phase Reconstruction Problem

How has this been done before?

- ▶ Griffin-Lim Reconstruction
- ▶ Autoregressive Neural Networks (WaveNet)
- ▶ GANs (HiFi-GAN)
- ▶ Diffusion (DiffWave)
- ▶ ...and more

# Diffusion

- ▶ Main idea
- ▶ Where do Neural Networks come in?
- ▶ How can it be used for vocoding?

# Diffusion

Collected training data  $\mathbb{X} = \{\mathbf{x}^{(i)}\}_{i=1}^N$  drawn as  $\mathbf{x} \sim q_{\text{data}}$ .

# Diffusion

Collected training data  $\mathbb{X} = \{\mathbf{x}^{(i)}\}_{i=1}^N$  drawn as  $\mathbf{x} \sim q_{\text{data}}$ .

Want to be able to generate *new data* from the "target" distribution!

# Diffusion

Collected training data  $\mathbb{X} = \{\mathbf{x}^{(i)}\}_{i=1}^N$  drawn as  $\mathbf{x} \sim q_{\text{data}}$ .

Want to be able to generate *new data* from the "target" distribution!

How?

# Diffusion

- ▶ Transform data from a simple distribution  $p_{\text{latent}}$  (e.g. noise) into the complex target distribution  $q_{\text{data}}$



# Diffusion

- ▶ Transform data from a simple distribution  $p_{\text{latent}}$  (e.g. noise) into the complex target distribution  $q_{\text{data}}$
- ▶ In one step is hard, in several steps easier

# Diffusion

- ▶ Transform data from a simple distribution  $p_{\text{latent}}$  (e.g. noise) into the complex target distribution  $q_{\text{data}}$
- ▶ In one step is hard, in several steps easier
- ▶ Also easy to *destroy* data into noise, hard to do the inverse

# Diffusion

- ▶ Transform data from a simple distribution  $p_{\text{latent}}$  (e.g. noise) into the complex target distribution  $q_{\text{data}}$
- ▶ In one step is hard, in several steps easier
- ▶ Also easy to *destroy* data into noise, hard to do the inverse
- ▶ Teach a model to perform each "denoising" step

# Forward Process

$T$ : Number of diffusion steps

# Forward Process

$T$ : Number of diffusion steps

$\mathbf{x}_0$ : Ground truth sample

# Forward Process

$T$ : Number of diffusion steps

$\mathbf{x}_0$ : Ground truth sample

$\mathbf{x}_t$ : Sample after  $t$  diffusion steps

# Forward Process

$T$ : Number of diffusion steps

$\mathbf{x}_0$ : Ground truth sample

$\mathbf{x}_t$ : Sample after  $t$  diffusion steps

$\mathbf{x}_T$ : Final sample, hopefully close to  $p_{\text{latent}}$

# Forward Process

$T$ : Number of diffusion steps

$\mathbf{x}_0$ : Ground truth sample

$\mathbf{x}_t$ : Sample after  $t$  diffusion steps

$\mathbf{x}_T$ : Final sample, hopefully close to  $p_{\text{latent}}$

$$\mathbf{x}_t \sim q(\mathbf{x}_t \mid \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I})$$

$$\mathbf{x}_t = \sqrt{1 - \beta_t} \mathbf{x}_{t-1} + \beta_t \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$



# Forward Process

$T$ : Number of diffusion steps

$\mathbf{x}_0$ : Ground truth sample

$\mathbf{x}_t$ : Sample after  $t$  diffusion steps

$\mathbf{x}_T$ : Final sample, hopefully close to  $p_{\text{latent}}$

$$\mathbf{x}_t \sim q(\mathbf{x}_t \mid \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I})$$

$$\mathbf{x}_t = \sqrt{1 - \beta_t} \mathbf{x}_{t-1} + \beta_t \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

Variance schedule  $\beta_t \in [0, 1]$ : At which "rate" we tend towards a unit Gaussian

# Forward Process

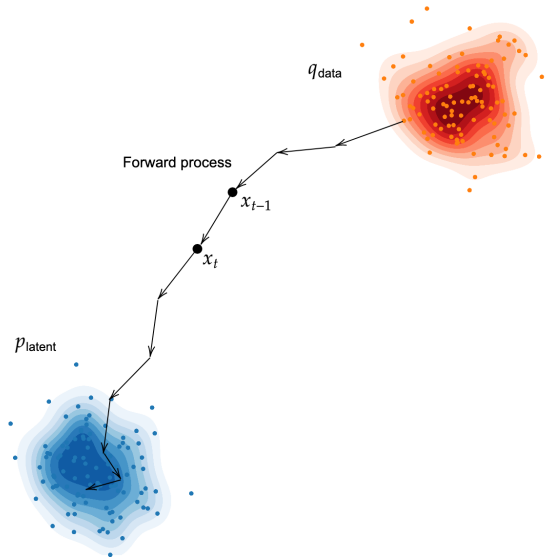
Quick way to get a noisy sample  $\mathbf{x}_t$ :

$$\bar{\alpha}_t = \prod_{i=1}^T (1 - \beta_i)$$

$$\mathbf{x}_t \sim q(\mathbf{x}_t \mid \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, \sqrt{1 - \bar{\alpha}_t} \mathbf{I}),$$

$$\mathbf{x} = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

# Forward Process



# Backward Process

- ▶ We want to be able to generate samples using the backward process

# Backward Process

- ▶ We want to be able to generate samples using the backward process

- ▶ Model each transition as

$$p_{\theta}(\mathbf{x}_{t-1} \mid \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_{\theta}(\mathbf{x}_t, t), \boldsymbol{\Sigma}_{\theta}(\mathbf{x}_t, t))$$

# Backward Process

