


Diffusion-Based Voice Conversion with **FAST** Maximum Likelihood Sampling Scheme



ICLR 2022

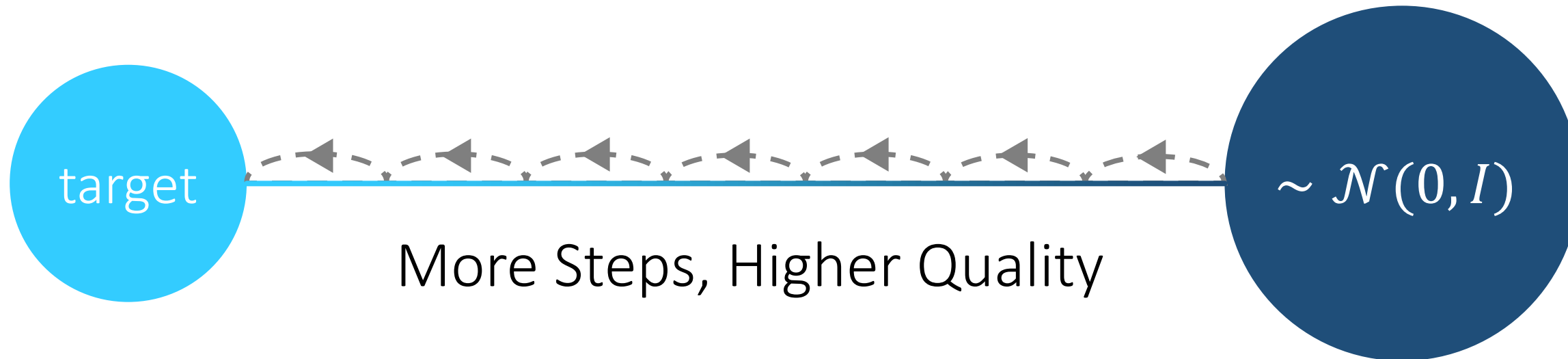
*Vadim Popov, Ivan Vovk, Vladimir Gogoryan, Tasnima Sadekova,
Mikhail Sergeevich Kudinov, Jiansheng Wei*

Why Use Diffusion Models?

- More Stable than GAN
- Higher Quality than VAE
- Easier to Design than Flow Models

But **Slower** than Them

Diffusion Steps

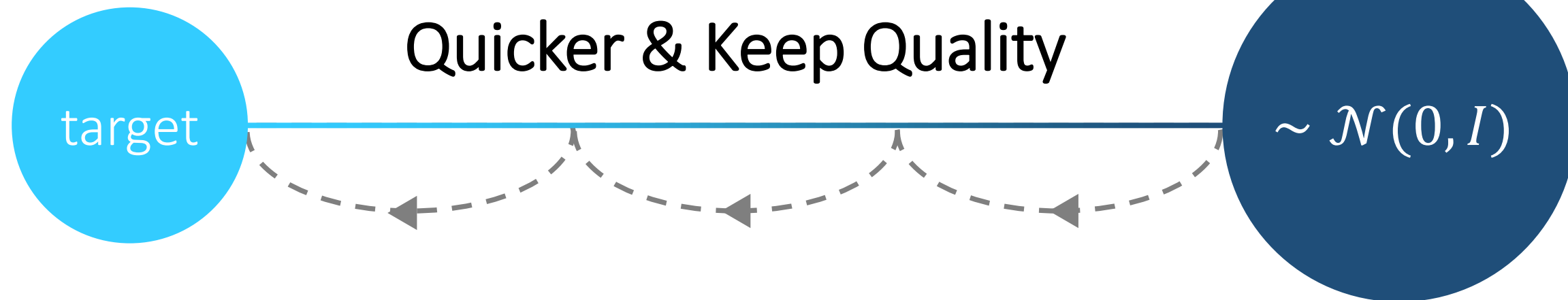


| CIFAR10 (32×32) | | | | | CelebA (64×64) | | | | |
|----------------------------|--------|-------|------|-------------|---------------------------|--------|-------|-------|-------------|
| 10 | 20 | 50 | 100 | 1000 | 10 | 20 | 50 | 100 | 1000 |
| 367.43 | 133.37 | 32.72 | 9.99 | 3.17 | 299.71 | 183.83 | 71.71 | 45.20 | 3.26 |

CIFAR10 and CelebA image generation measured in FID.

Image source: [Song et al., 2020](#)

Speed Up



Propose a **FAST** Sampling Scheme

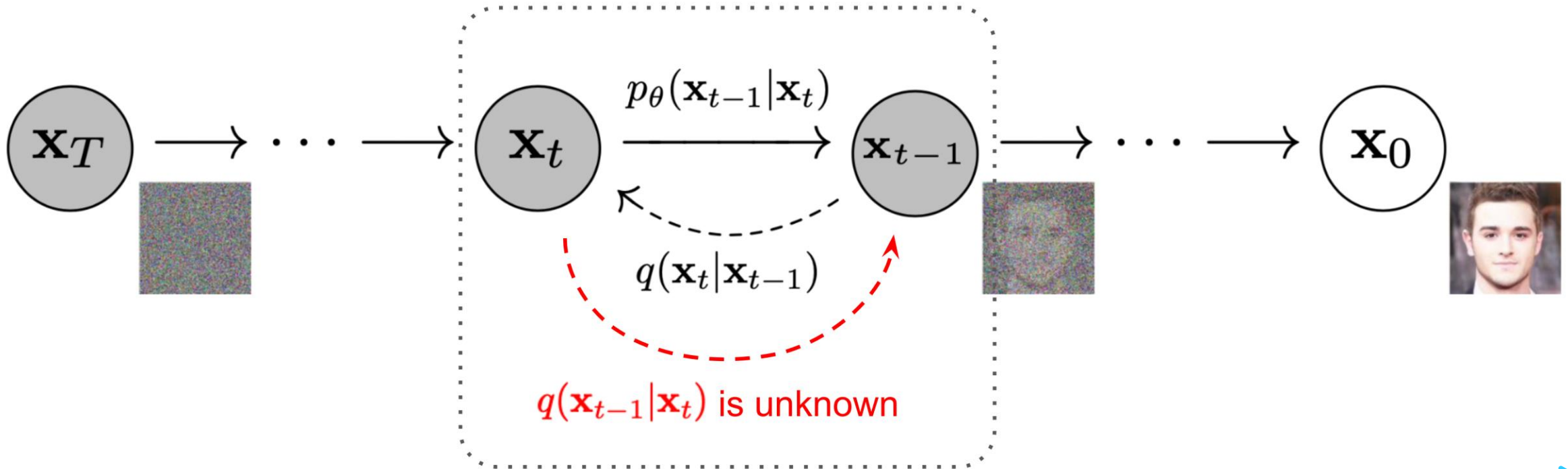
+

Average Voice Encoder

= SOTA Any to Any VC

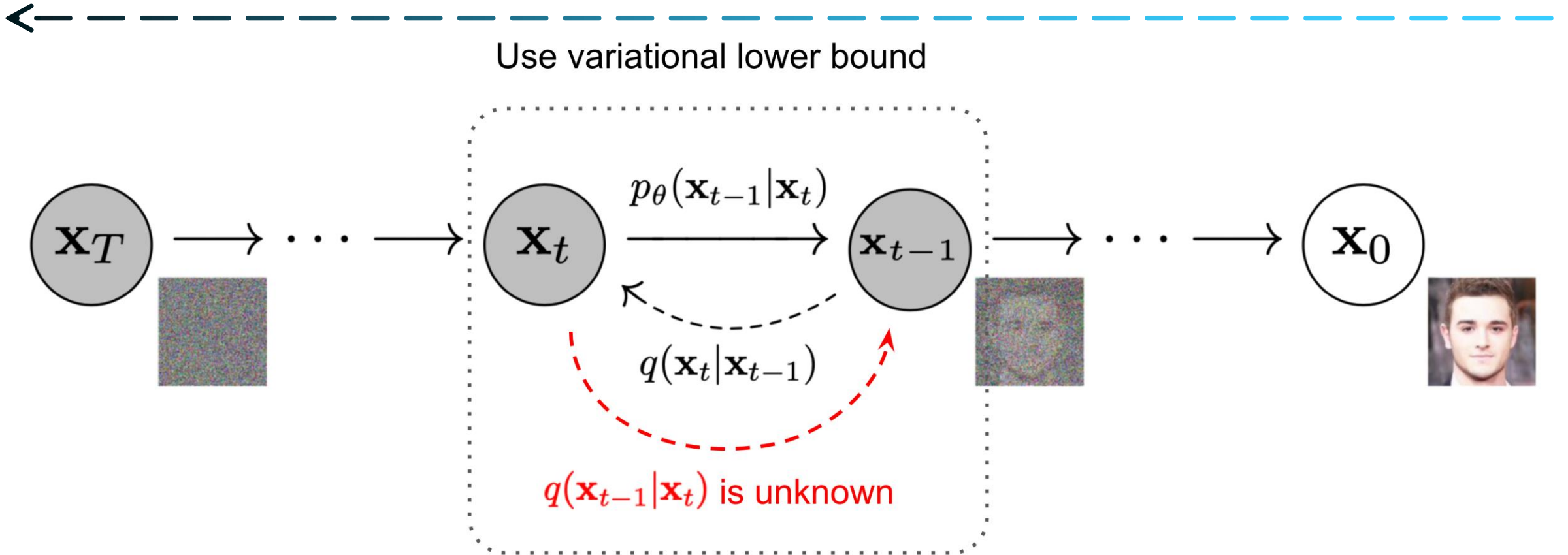
Forward Diffusion (Training)

Use variational lower bound



Reverse Diffusion (Inference)

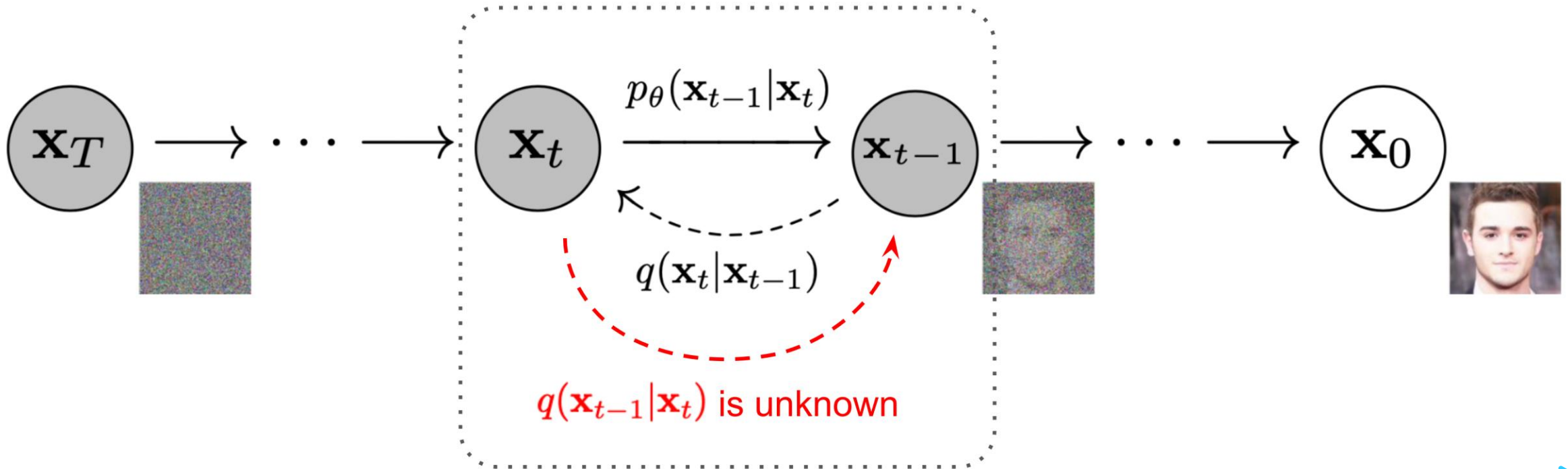
Forward Diffusion



Mix \mathbf{x}_0 with noise to get \mathbf{x}_t , and train the model to estimate noise.

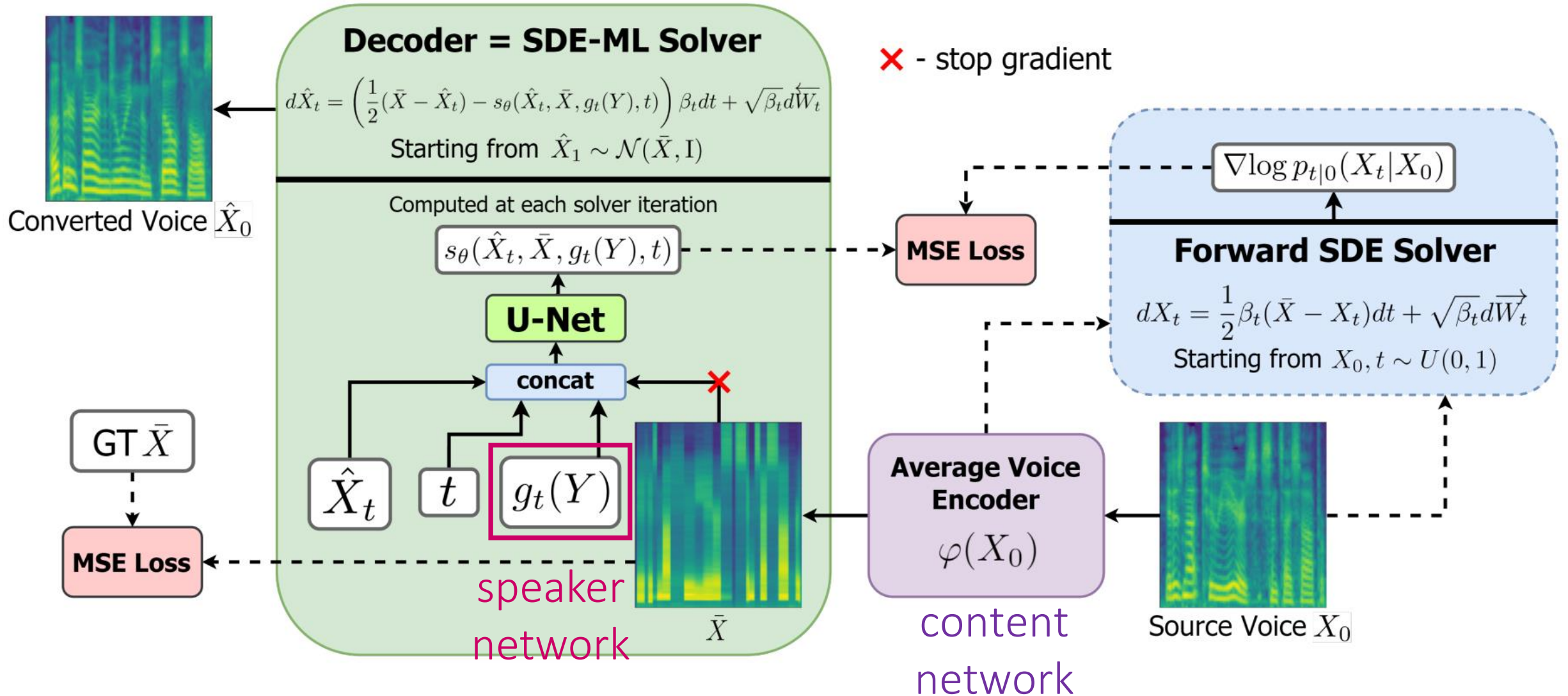
Reverse the diffusion process and sampling from it, you can generate real samples from Gaussian noise.

Use variational lower bound



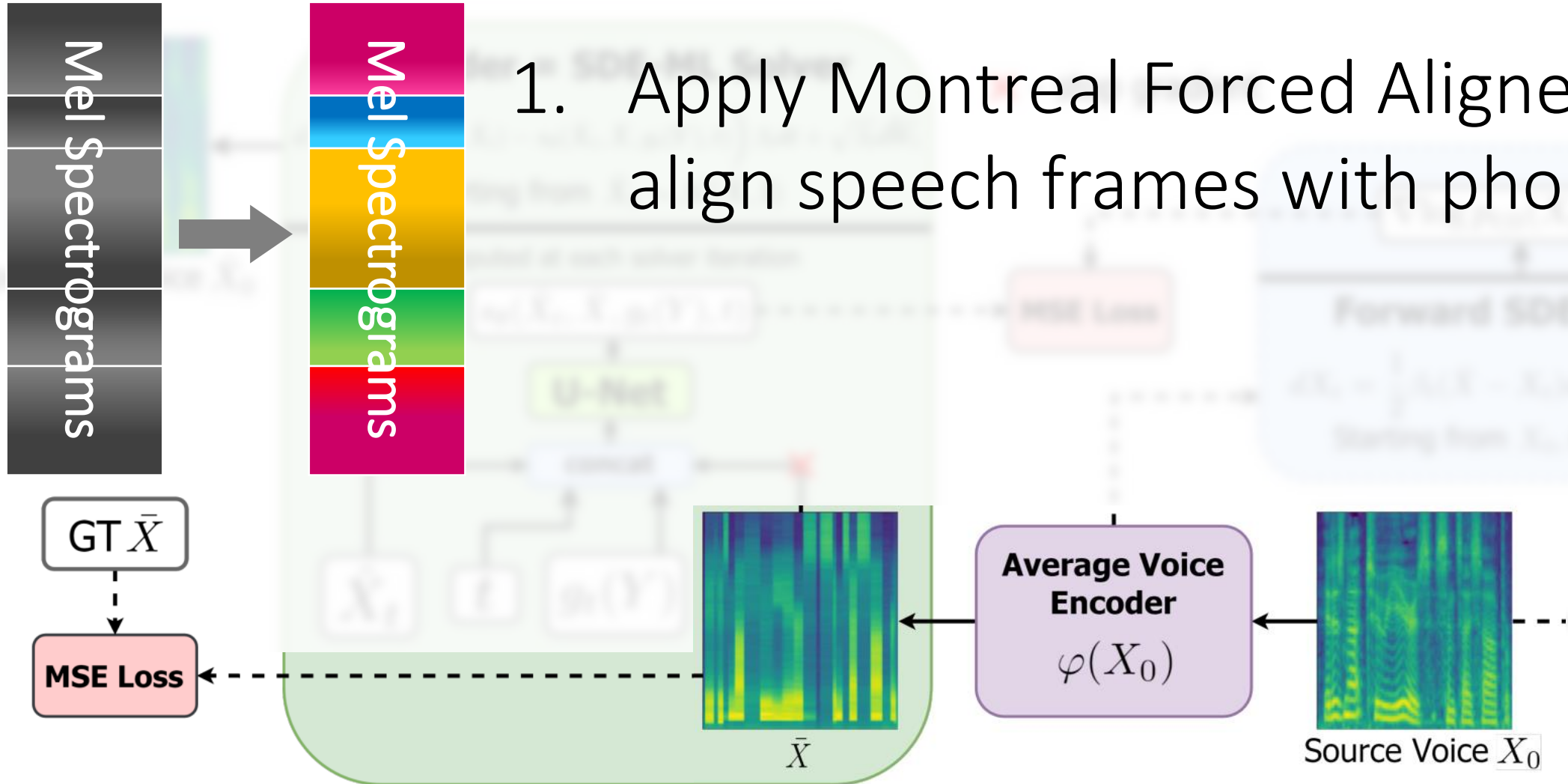
Reverse Diffusion (Inference)

Voice Conversion Diffusion Model



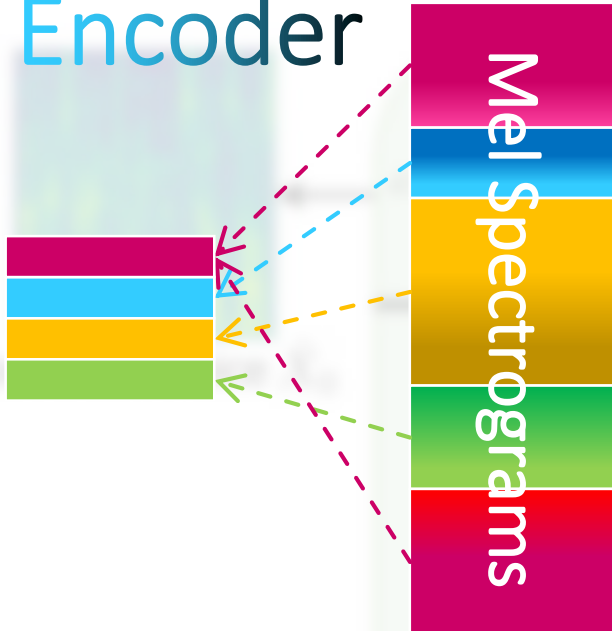
Average Voice Encoder

1. Apply Montreal Forced Aligner to align speech frames with phonemes.

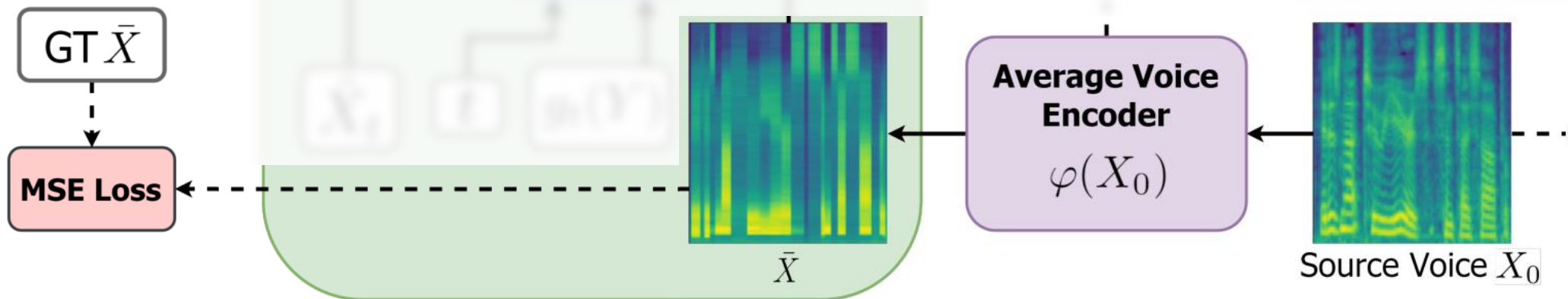


Average Voice

Encoder



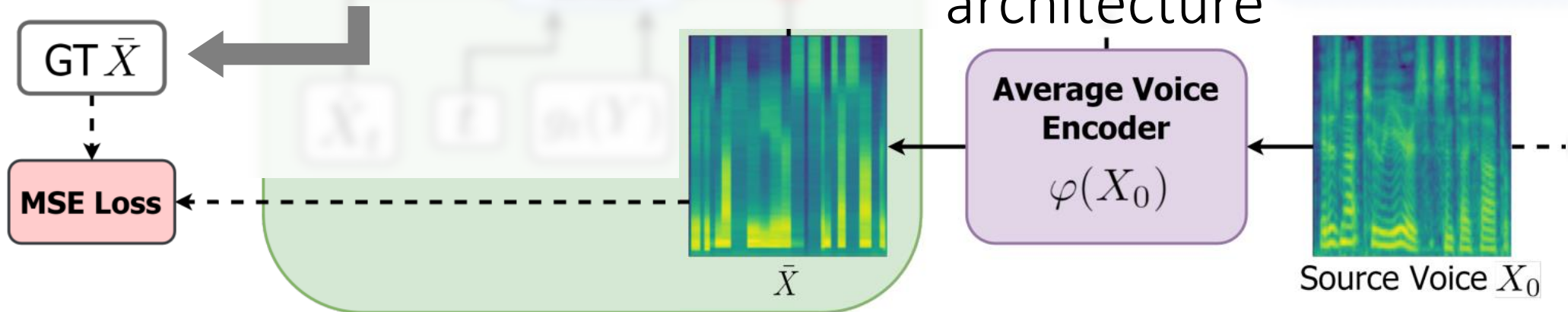
2. Calculate the average Mel feature for each phoneme across **the whole LibriTTS dataset**.



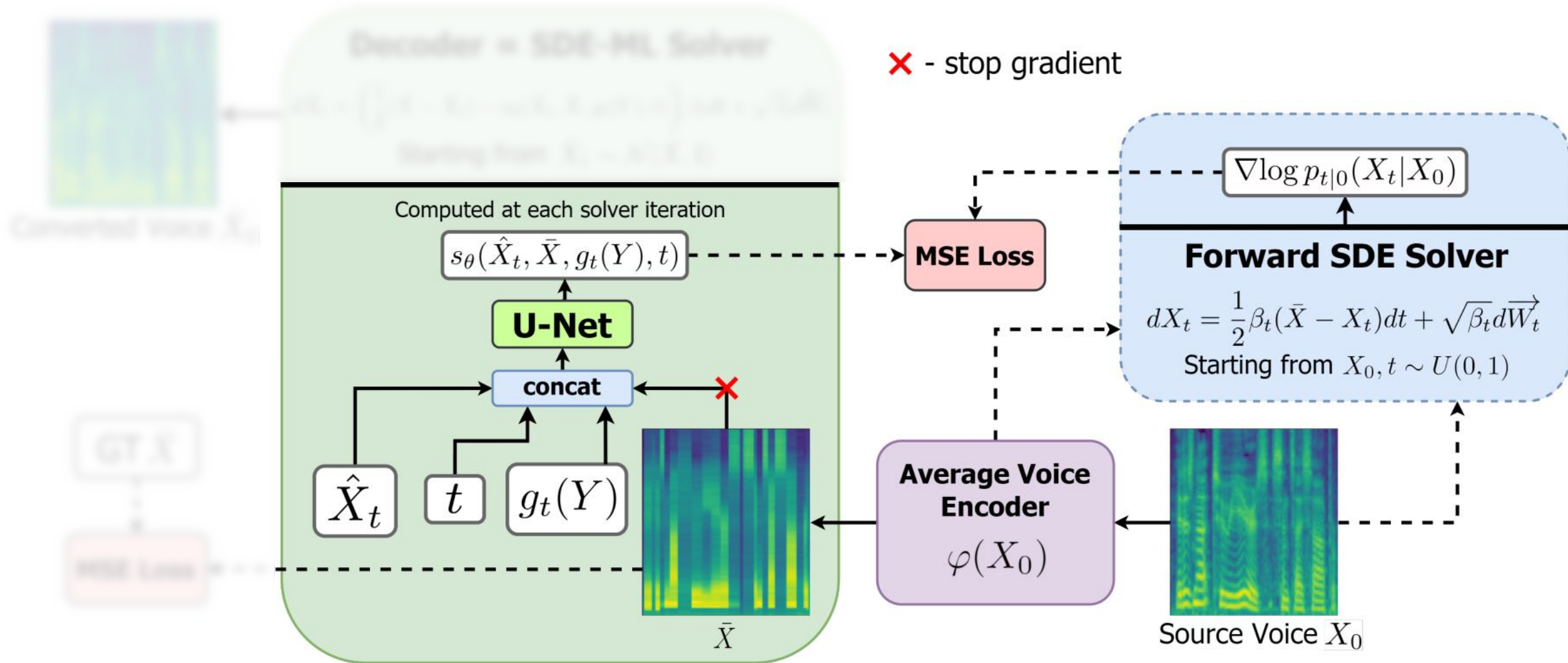
Average Voice Encoder

3. Encoder is trained with “Average Voice” Mel Spectrogram as the target.

Transformer-based architecture



Forward Diffusion



Forward Diffusion: Sample X_t

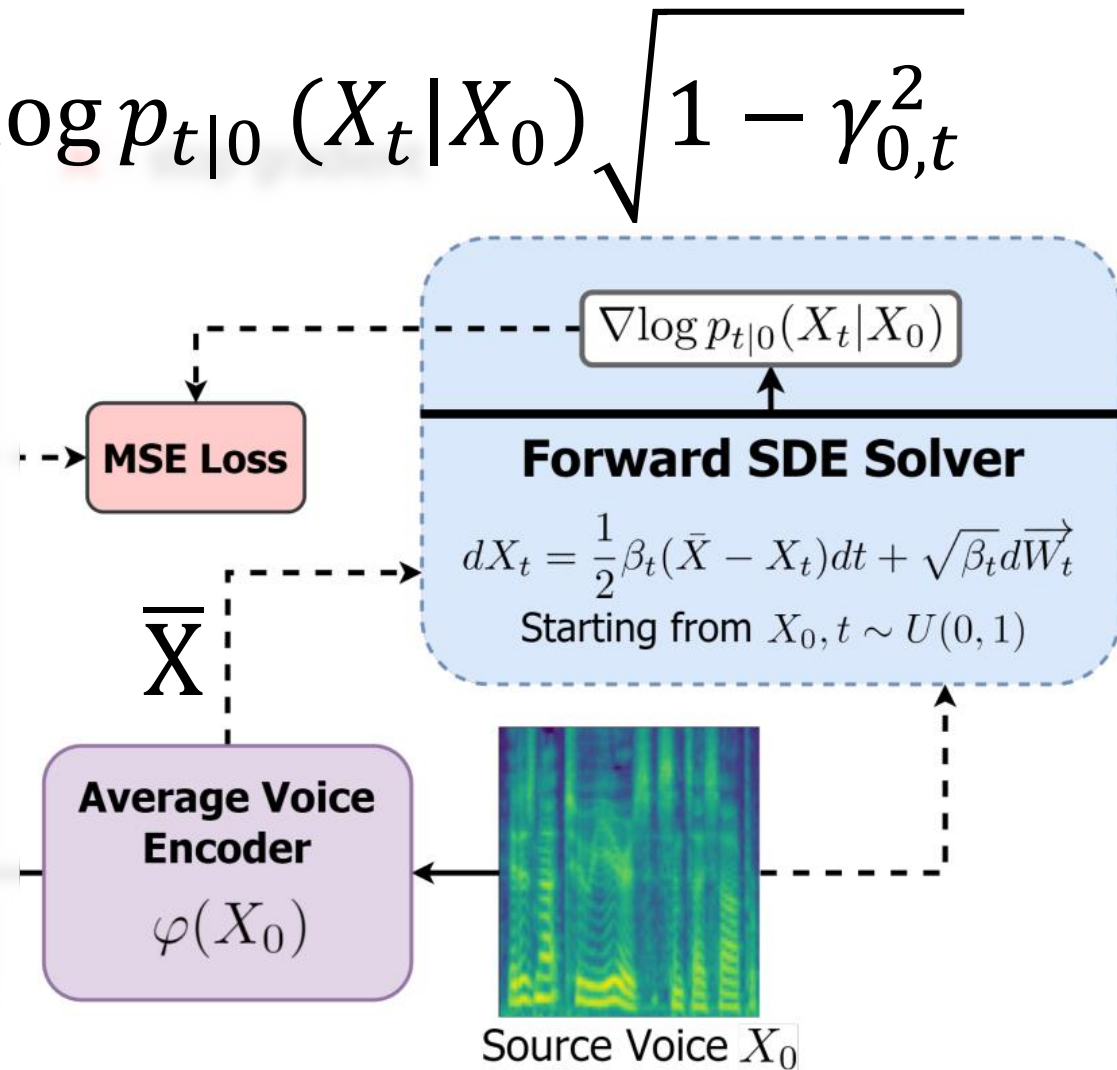
$$X_t = \gamma_{0,t}X_0 + (1 - \gamma_{0,t})\bar{X} + \nabla \log p_{t|0}(X_t|X_0)\sqrt{1 - \gamma_{0,t}^2}$$

$$\nabla \log p_{t|0}(X_t|X_0) \sim \mathcal{N}(0, I)$$

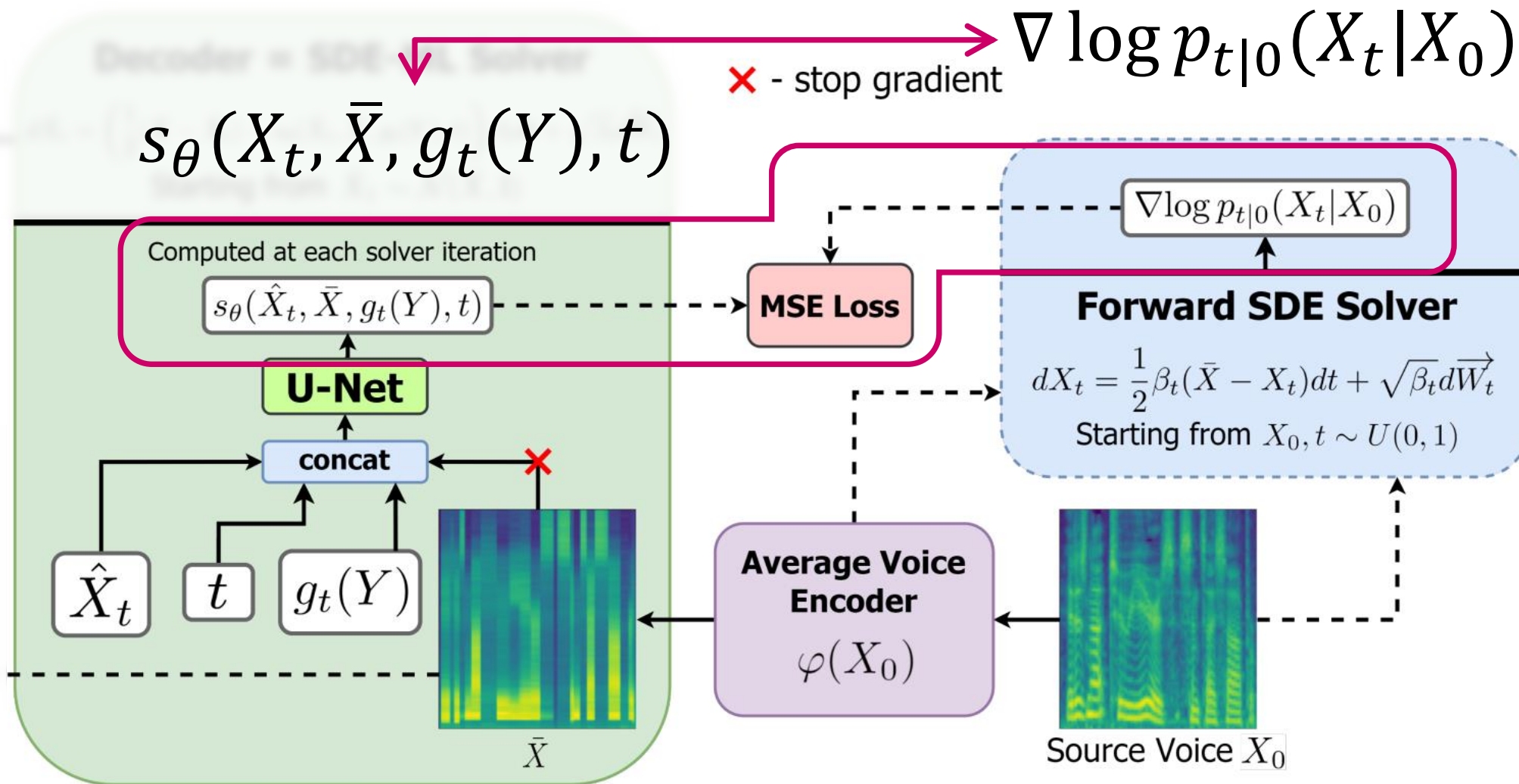
$$\gamma_{s,t} = e^{-\frac{1}{2} \int_s^t \beta_u du}$$

$$\beta_t = \beta_0 + t(\beta_1 - \beta_0)$$

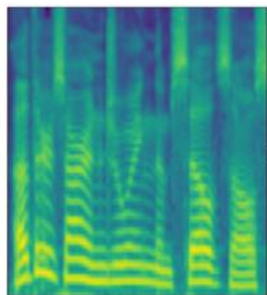
Hyper parameters



Forward Diffusion: Loss



Reverse Diffusion: Euler-Maruyama



Converted Voice \hat{X}_0

Decoder = SDE-ML Solver

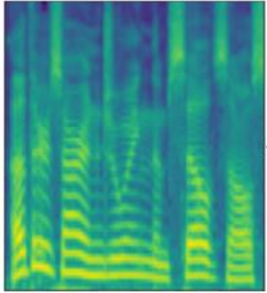
$$d\hat{X}_t = \left(\frac{1}{2}(\bar{X} - \hat{X}_t) - s_\theta(\hat{X}_t, \bar{X}, g_t(Y), t) \right) \beta_t dt + \sqrt{\beta_t} d\tilde{W}_t$$

Starting from $\hat{X}_1 \sim \mathcal{N}(\bar{X}, I)$

$$\begin{aligned} \hat{X}_{t-h} &= \hat{X}_t + \hat{\sigma}_{t,h} \nabla \log p_{t|0}(X_t | X_0) \\ &\quad + \beta_t h \left(\left(\frac{1}{2} + \hat{\omega}_{t,h} \right) (\hat{X}_t - \bar{X}) + (1 + \hat{\omega}_{t,h}) s_\theta(X_t, \bar{X}, g_t(Y), t) \right) \end{aligned}$$

step size

Reverse Diffusion: Maximum Likelihood



Converted Voice \hat{X}_0

Decoder = SDE-ML Solver

$$d\hat{X}_t = \left(\frac{1}{2}(\bar{X} - \hat{X}_t) - s_\theta(\hat{X}_t, \bar{X}, g_t(Y), t) \right) \beta_t dt + \sqrt{\beta_t} d\hat{W}_t$$

Starting from $\hat{X}_1 \sim \mathcal{N}(\bar{X}, I)$

By Theorem 1.

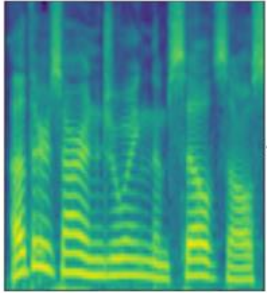
$$\hat{\sigma}_{t,h} = \sigma_{t,h}^*, \hat{\omega}_{t,h} = \omega_{t,h}^*, \hat{\kappa}_{t,h} = \kappa_{t,h}^*$$

$$\hat{X}_{t-h}$$

$$= \hat{X}_t + \hat{\sigma}_{t,h} \nabla \log p_{t|0}(X_t | X_0)$$

$$+ \beta_t \underset{\text{step size}}{h} \left(\left(\frac{1}{2} + \hat{\omega}_{t,h} \right) (\hat{X}_t - \bar{X}) + (1 + \hat{\kappa}_{t,h}) s_\theta(X_t, \bar{X}, g_t(Y), t) \right)$$

Reverse Diffusion: Maximum Likelihood



Converted Voice \hat{X}_0

Decoder = SDE-ML Solver

$$d\hat{X}_t = \left(\frac{1}{2}(\bar{X} - \hat{X}_t) - s_\theta(\hat{X}_t, \bar{X}, g_t(Y), t) \right) \beta_t dt + \sqrt{\beta_t} d\hat{W}_t$$

Starting from $\hat{X}_1 \sim \mathcal{N}(\bar{X}, I)$

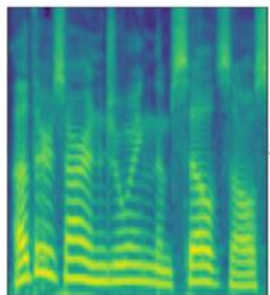
By Theorem 1.

$$\hat{\sigma}_{t,h} = \sigma_{t,h}^*, \quad \hat{\omega}_{t,h} = \omega_{t,h}^*, \quad \hat{\kappa}_{t,h} = \kappa_{t,h}^*$$

$$\kappa_{t,h}^* = \frac{\nu_{t-h,t}(1 - \gamma_{0,t}^2)}{\gamma_{0,t}\beta_t h} - 1, \quad \omega_{t,h}^* = \frac{\mu_{t-h,t} - 1}{\beta_t h} + \frac{1 + \kappa_{t,h}^*}{1 - \gamma_{0,t}^2} - \frac{1}{2},$$

$$(\sigma_{t,h}^*)^2 = \sigma_{t-h,t}^2 + \frac{1}{n} \nu_{t-h,t}^2 \mathbb{E}_{X_t} [\text{Tr}(\text{Var}(X_0|X_t))],$$

Reverse Diffusion: Maximum Likelihood



Converted Voice \hat{X}_0

Decoder = SDE-ML Solver

$$d\hat{X}_t = \left(\frac{1}{2}(\bar{X} - \hat{X}_t) - s_\theta(\hat{X}_t, \bar{X}, g_t(Y), t) \right) \beta_t dt + \sqrt{\beta_t} d\hat{W}_t$$

Starting from $\hat{X}_1 \sim \mathcal{N}(\bar{X}, I)$

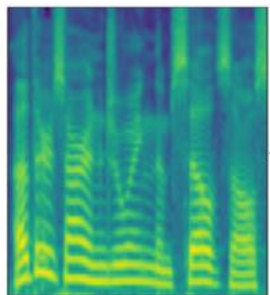
$$\mu_{s,t} = \gamma_{s,t} \frac{1 - \gamma_{0,s}^2}{1 - \gamma_{0,t}^2}, \quad \nu_{s,t} = \gamma_{0,s} \frac{1 - \gamma_{s,t}^2}{1 - \gamma_{0,t}^2}, \quad \sigma_{s,t}^2 = \frac{(1 - \gamma_{0,s}^2)(1 - \gamma_{s,t}^2)}{1 - \gamma_{0,t}^2},$$

$$\kappa_{t,h}^* = \frac{\nu_{t-h,t}(1 - \gamma_{0,t}^2)}{\gamma_{0,t}\beta_t h} - 1, \quad \omega_{t,h}^* = \frac{\mu_{t-h,t} - 1}{\beta_t h} + \frac{1 + \kappa_{t,h}^*}{1 - \gamma_{0,t}^2} - \frac{1}{2},$$

$$(\sigma_{t,h}^*)^2 = \sigma_{t-h,t}^2 + \frac{1}{n} \nu_{t-h,t}^2 \mathbb{E}_{\hat{X}_t} [\text{Tr}(\text{Var}(X_0 | X_t))],$$

Without in source code?

Reverse Diffusion: Maximum Likelihood



Converted Voice \hat{X}_0

Decoder = SDE-ML Solver

$$d\hat{X}_t = \left(\frac{1}{2}(\bar{X} - \hat{X}_t) - s_\theta(\hat{X}_t, \bar{X}, g_t(Y), t) \right) \beta_t dt + \sqrt{\beta_t} d\tilde{W}_t$$

Starting from $\hat{X}_1 \sim \mathcal{N}(\bar{X}, I)$

$$\hat{X}_1 \sim \mathcal{N}(\bar{X}, I)$$

for $i = 0$ **to** $N - 1$ **do**

$$t \leftarrow i \times h$$

$$X'_{t-h} \leftarrow \hat{X}_t + \beta_t h \left(\left(\frac{1}{2} + \hat{\omega}_{t,h} \right) (\hat{X}_t - \bar{X}) + (1 + \hat{\kappa}_{t,h}) s_\theta(X_t, \bar{X}, g_t(Y), t) \right)$$

$$\nabla \log p_{t|0}(X_t | X_0) \sim \mathcal{N}(0, I)$$

$$\hat{X}_{t-h} \leftarrow X'_{t-h} + \hat{\sigma}_{t,h} \nabla \log p_{t|0}(X_t | X_0)$$

return \hat{X}_0

Speaker Conditional Analysis

Input types for speaker conditioning $g_t(Y)$ compared in terms of speaker similarity.

| | <i>Diff-LibriTTS</i> | | | <i>Diff-VCTK</i> | | |
|---------------|----------------------|--------------|--------------|------------------|--------------|--------------|
| | <i>d-only</i> | <i>wodyn</i> | <i>whole</i> | <i>d-only</i> | <i>wodyn</i> | <i>whole</i> |
| Most similar | 27.0% | 38.0% | 34.1% | 27.2% | 46.7% | 23.6% |
| Least similar | 28.9% | 29.3% | 38.5% | 25.3% | 23.9% | 48.6% |

- d-only: Y = target Mel-spectrogram Y_0
- wodyn: Y = **Noisy** target Mel-spectrogram Y_t
- whole: $Y = \{Y_t, Y_{0.5/15}, Y_{1.5/15}, \dots, Y_{14.5/15}\}$, channel = 16

Any-to-Any Voice Conversion

| | VCTK test (9 speakers, 54 pairs) | | Whole test (25 speakers, 350 pairs) | |
|------------------------|-----------------------------------|-----------------------------------|-------------------------------------|-----------------------------------|
| | Naturalness | Similarity | Naturalness | Similarity |
| <i>AGAIN-VC</i> | 1.98 ± 0.05 | 1.97 ± 0.08 | 1.87 ± 0.03 | 1.75 ± 0.04 |
| <i>FragmentVC</i> | 2.20 ± 0.06 | 2.45 ± 0.09 | 1.91 ± 0.03 | 1.93 ± 0.04 |
| <i>VQMIVC</i> | 2.89 ± 0.06 | 2.60 ± 0.10 | 2.48 ± 0.04 | 1.95 ± 0.04 |
| <i>Diff-VCTK-ML-6</i> | 3.73 ± 0.06 | 3.47 ± 0.09 | 3.39 ± 0.04 | 2.69 ± 0.05 |
| <i>Diff-VCTK-ML-30</i> | 3.73 ± 0.06 | 3.57 ± 0.09 | 3.44 ± 0.04 | 2.71 ± 0.05 |
| <i>Ground truth</i> | 4.55 ± 0.05 | 4.52 ± 0.07 | 4.55 ± 0.05 | 4.52 ± 0.07 |

Conv Auto Encoder
Attention-based
Vector Quantization

Train on VCTK, **100** speakers

All subjective human evaluation was carried out on Amazon Mechanical Turk.

Any-to-Any Voice Conversion

| | VCTK test (9 speakers, 54 pairs) | | Whole test (25 speakers, 350 pairs) | |
|------------------------|-----------------------------------|-----------------------------------|-------------------------------------|-----------------------------------|
| | Naturalness | Similarity | Naturalness | Similarity |
| <i>AGAIN-VC</i> | 1.98 ± 0.05 | 1.97 ± 0.08 | 1.87 ± 0.03 | 1.75 ± 0.04 |
| <i>FragmentVC</i> | 2.20 ± 0.06 | 2.45 ± 0.09 | 1.91 ± 0.03 | 1.93 ± 0.04 |
| <i>VQMIVC</i> | 2.89 ± 0.06 | 2.60 ± 0.10 | 2.48 ± 0.04 | 1.95 ± 0.04 |
| <i>Diff-VCTK-ML-6</i> | 3.73 ± 0.06 | 3.47 ± 0.09 | 3.39 ± 0.04 | 2.69 ± 0.05 |
| <i>Diff-VCTK-ML-30</i> | 3.73 ± 0.06 | 3.57 ± 0.09 | 3.44 ± 0.04 | 2.71 ± 0.05 |
| <i>Ground truth</i> | 4.55 ± 0.05 | 4.52 ± 0.07 | 4.55 ± 0.05 | 4.52 ± 0.07 |

Conv Auto Encoder
Attention-based
Vector Quantization

Real-Time Factor on GPU (unknow model)

- 6 step: around **0.1**
- 30 step: around 0.5

Any-to-Any Voice Conversion

Train on LibriTTS

approximately **1100** speakers.

| | VCTK test (9 speakers, 54 pairs) | | Whole test (25 speakers, 350 pairs) | |
|----------------------------|-----------------------------------|-----------------------------------|-------------------------------------|-----------------------------------|
| | Naturalness | Similarity | Naturalness | Similarity |
| <i>Diff-LibriTTS-EM-6</i> | 1.68 ± 0.06 | 1.53 ± 0.07 | 1.57 ± 0.02 | 1.47 ± 0.03 |
| <i>Diff-LibriTTS-PF-6</i> | 3.11 ± 0.07 | 2.58 ± 0.11 | 2.99 ± 0.03 | 2.50 ± 0.04 |
| <i>Diff-LibriTTS-ML-6</i> | 3.84 ± 0.08 | 3.08 ± 0.11 | 3.80 ± 0.03 | 3.27 ± 0.05 |
| <i>Diff-LibriTTS-ML-30</i> | 3.96 ± 0.08 | 3.23 ± 0.11 | 4.02 ± 0.03 | 3.39 ± 0.05 |
| <i>BNE-PPG-VC</i> | 3.95 ± 0.08 | 3.27 ± 0.12 | 3.83 ± 0.03 | 3.03 ± 0.05 |

BEN-PPG-VC: combining a bottleneck feature extractor obtained from a phoneme recognizer with a seq2seq-based synthesis module.

Any-to-Any Voice Conversion

The proposed maximum likelihood (ML) sampling scheme over other sampling methods for a small number of inference steps.

| | VCTK test (9 speakers, 54 pairs) | | Whole test (25 speakers, 350 pairs) | |
|----------------------------|-----------------------------------|-----------------------------------|-------------------------------------|-----------------------------------|
| | Naturalness | Similarity | Naturalness | Similarity |
| <i>Diff-LibriTTS-EM-6</i> | 1.68 ± 0.06 | 1.53 ± 0.07 | 1.57 ± 0.02 | 1.47 ± 0.03 |
| <i>Diff-LibriTTS-PF-6</i> | 3.11 ± 0.07 | 2.58 ± 0.11 | 2.99 ± 0.03 | 2.50 ± 0.04 |
| <i>Diff-LibriTTS-ML-6</i> | 3.84 ± 0.08 | 3.08 ± 0.11 | 3.80 ± 0.03 | 3.27 ± 0.05 |
| <i>Diff-LibriTTS-ML-30</i> | 3.96 ± 0.08 | 3.23 ± 0.11 | 4.02 ± 0.03 | 3.39 ± 0.05 |
| <i>BNE-PPG-VC</i> | 3.95 ± 0.08 | 3.27 ± 0.12 | 3.83 ± 0.03 | 3.03 ± 0.05 |

BEN-PPG-VC: combining a bottleneck feature extractor obtained from a phoneme recognizer with a seq2seq-based synthesis module.

Maximum Likelihood Sampling

Euler-Maruyama



Probability Flow



Maximum Likelihood



CIFAR-10 images randomly sampled from VP DPM by running 10 reverse diffusion steps.

Conclusion

- Average Voice Encoder
a new disentanglement method.
- Diffusion-based Decoder
achieve good results both in terms of similarity and naturalness.
- Novel Sampling Scheme
High-quality results in just a few steps.