

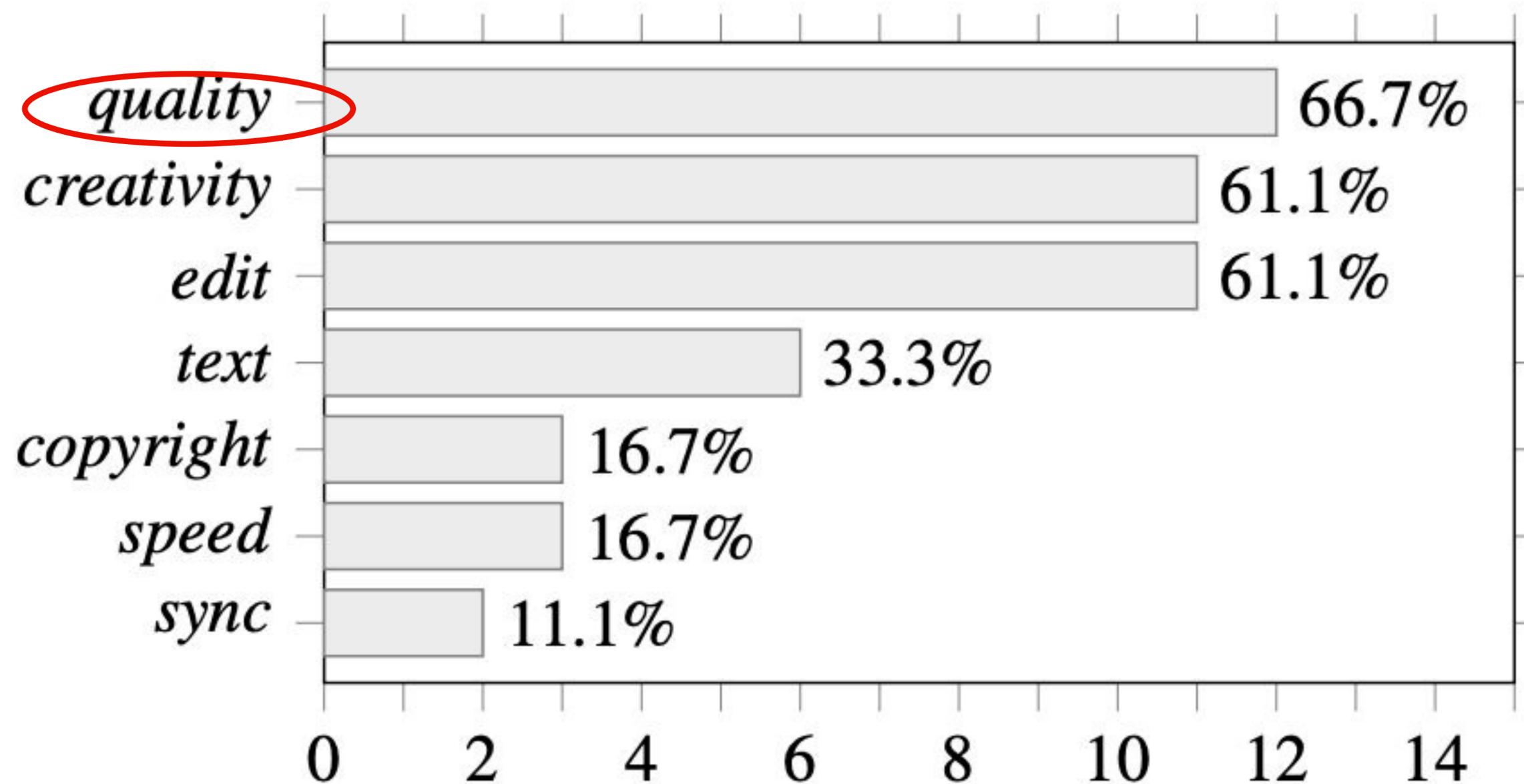
EDMSound: Spectrogram Based Diffusion Models for Efficient and High-Quality Audio Synthesis

Ge Zhu¹, Yutong Wen¹, Marc-André Carbonneau² and Zhiyao Duan¹

¹ ECE, University of Rochester

² Ubisoft La Forge

“What is the limitation(s) of the current text-conditioned audio generation as a product?”

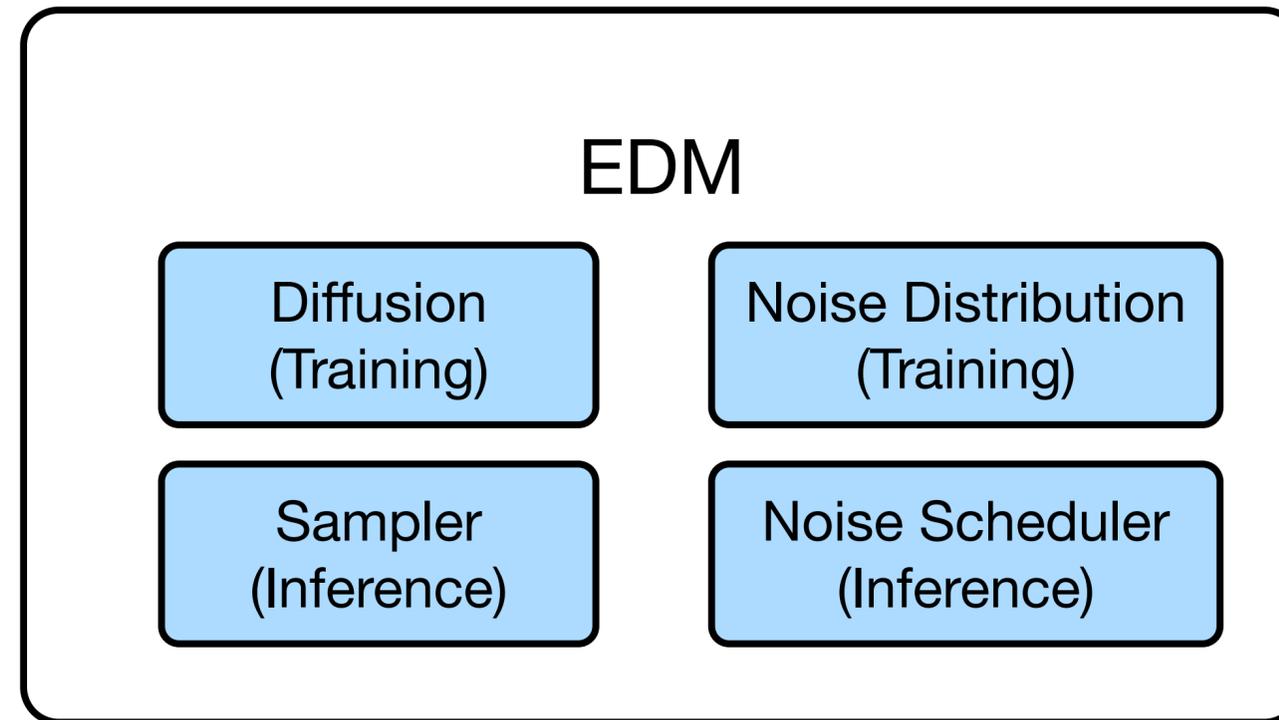


Overview

- EDMSound: Spectrogram based diffusion models
 - Training: Elucidated diffusion model (EDM) framework
 - Sampling: Exponential-integrator based deterministic solver
- Copy detection/memorization issue in diffusion models
 - Fine-tuned CLAP for audio

Elucidated Diffusion Model (EDM)

- Flexible diffusion training and inference with decoupled components

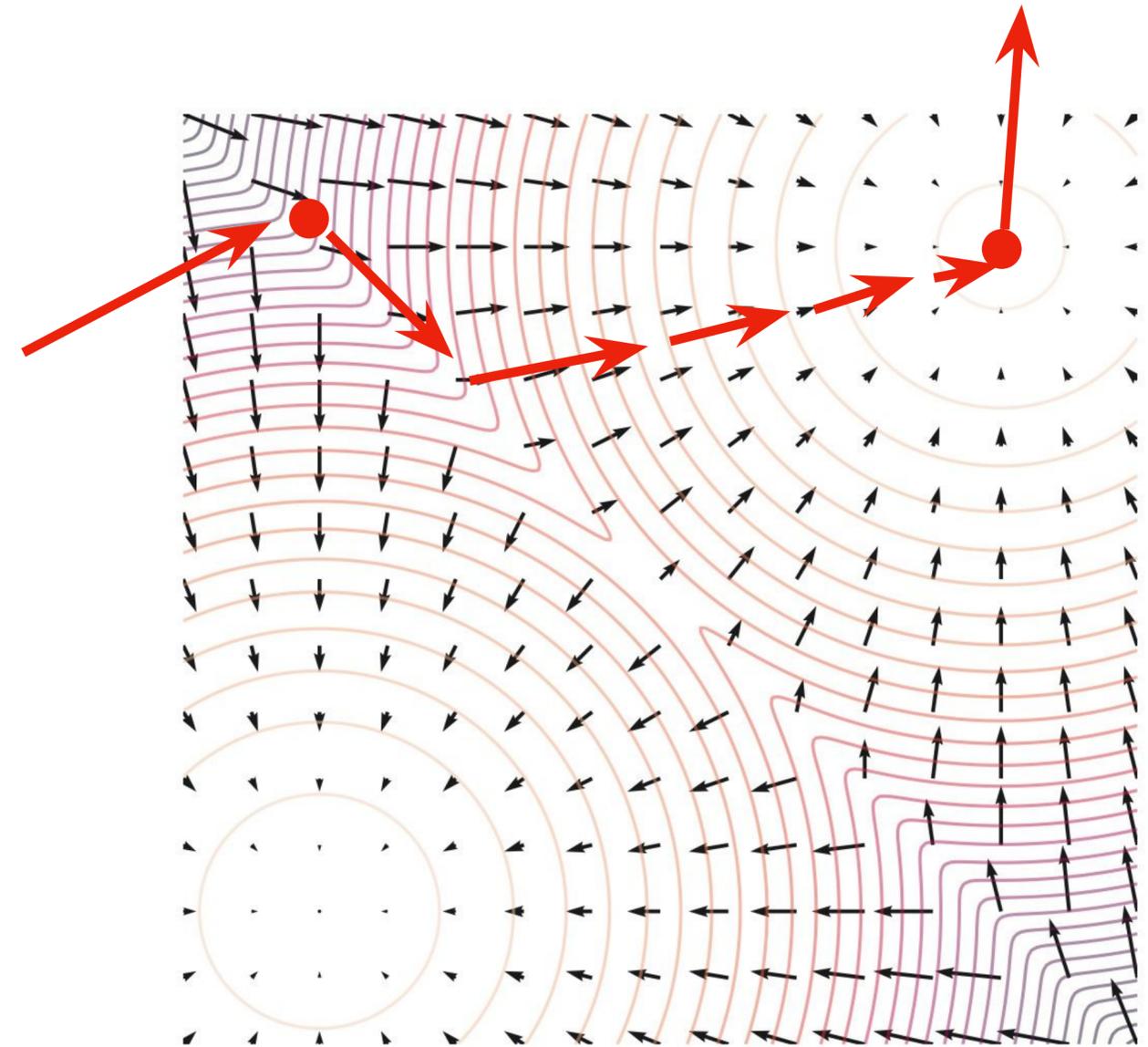
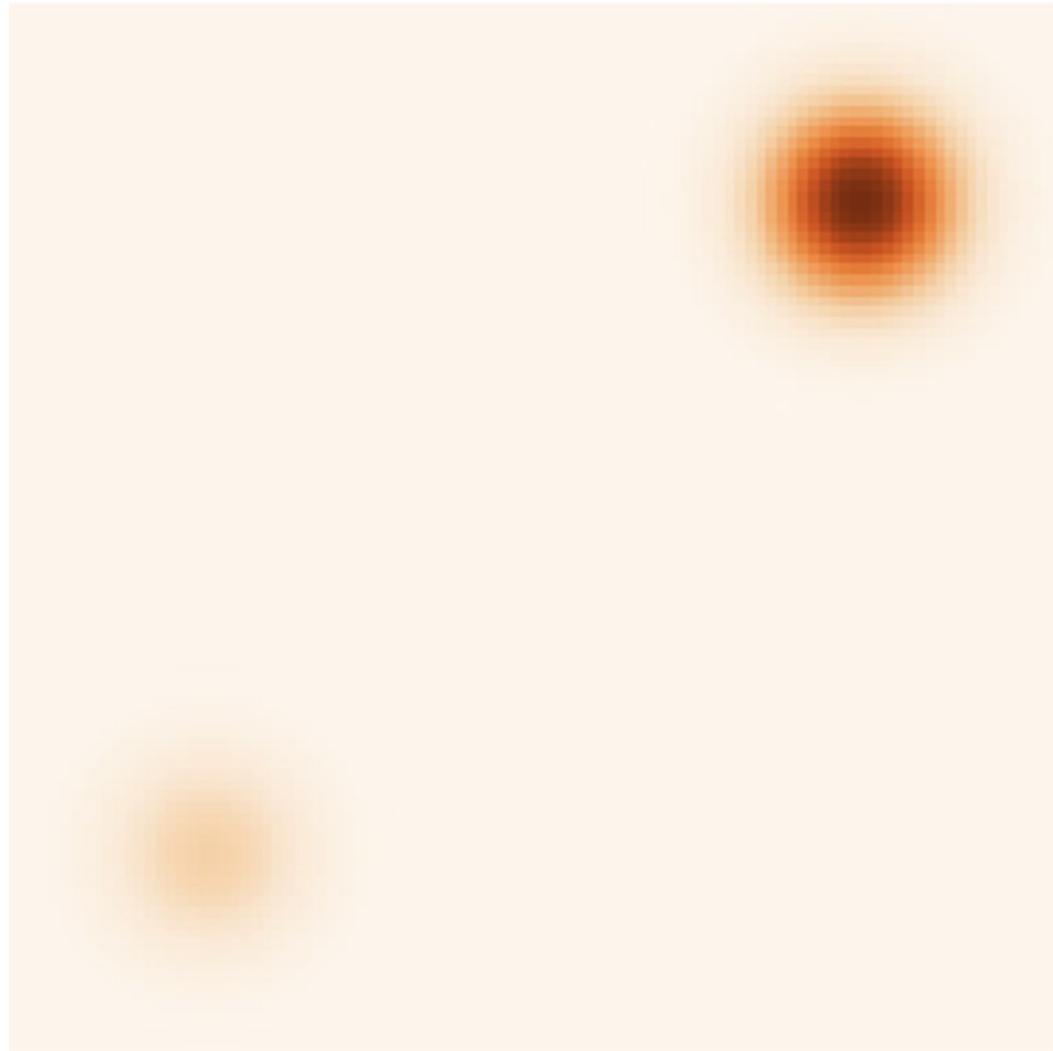


Perspectives on diffusion

1. *Diffusion models are autoencoders*
 2. *Diffusion models are deep latent variable models*
 3. *Diffusion models predict the score function*
 4. *Diffusion models solve reverse SDEs*
 5. *Diffusion models are flow-based models*
 6. *Diffusion models are recurrent neural networks*
 7. *Diffusion models are autoregressive models*
 8. *Diffusion models estimate expectations*
- ...

Diffusion model

$$\nabla \log p_t(x)$$



- Fitting (training) error: score function estimation
- Discretization (inference) error: limited steps or large step size

Especially important in
sampling with limited steps!

Elucidated Diffusion Model (EDM)

Fitting error:

- Pre-conditioning

Predicting x_0 or ϵ objectives are problematic

Network inputs depend on the noise level

- Loss reweighting

→ Both input and output magnitudes are fixed to unit variance

Diffusion Exponential Integrator Sampler

Ingredient 1: **Exponential Integrator** over Euler method.

Exact solution:

Original reverse sampling ODE:

$$\frac{d\hat{\mathbf{x}}}{dt} = \left[\mathbf{F}_t \hat{\mathbf{x}} - \frac{1}{2} \mathbf{G}_t \mathbf{G}_t^T \mathbf{s}_\theta(\hat{\mathbf{x}}, t) \right] \rightarrow \text{Semi-linear} \rightarrow \hat{\mathbf{x}}_t = \Psi(t, s) \hat{\mathbf{x}}_s + \int_s^t \Psi(t, \tau) \left[-\frac{1}{2} \mathbf{G}_\tau \mathbf{G}_\tau^T \mathbf{s}_\theta(\hat{\mathbf{x}}_\tau, \tau) \right] d\tau$$

Discretization:

$$\hat{\mathbf{x}}_{t-\Delta t} = \hat{\mathbf{x}}_t - \left[\mathbf{F}_t \hat{\mathbf{x}}_t - \frac{1}{2} \mathbf{G}_t \mathbf{G}_t^T \mathbf{s}_\theta(\hat{\mathbf{x}}_t, t) \right] \Delta t$$

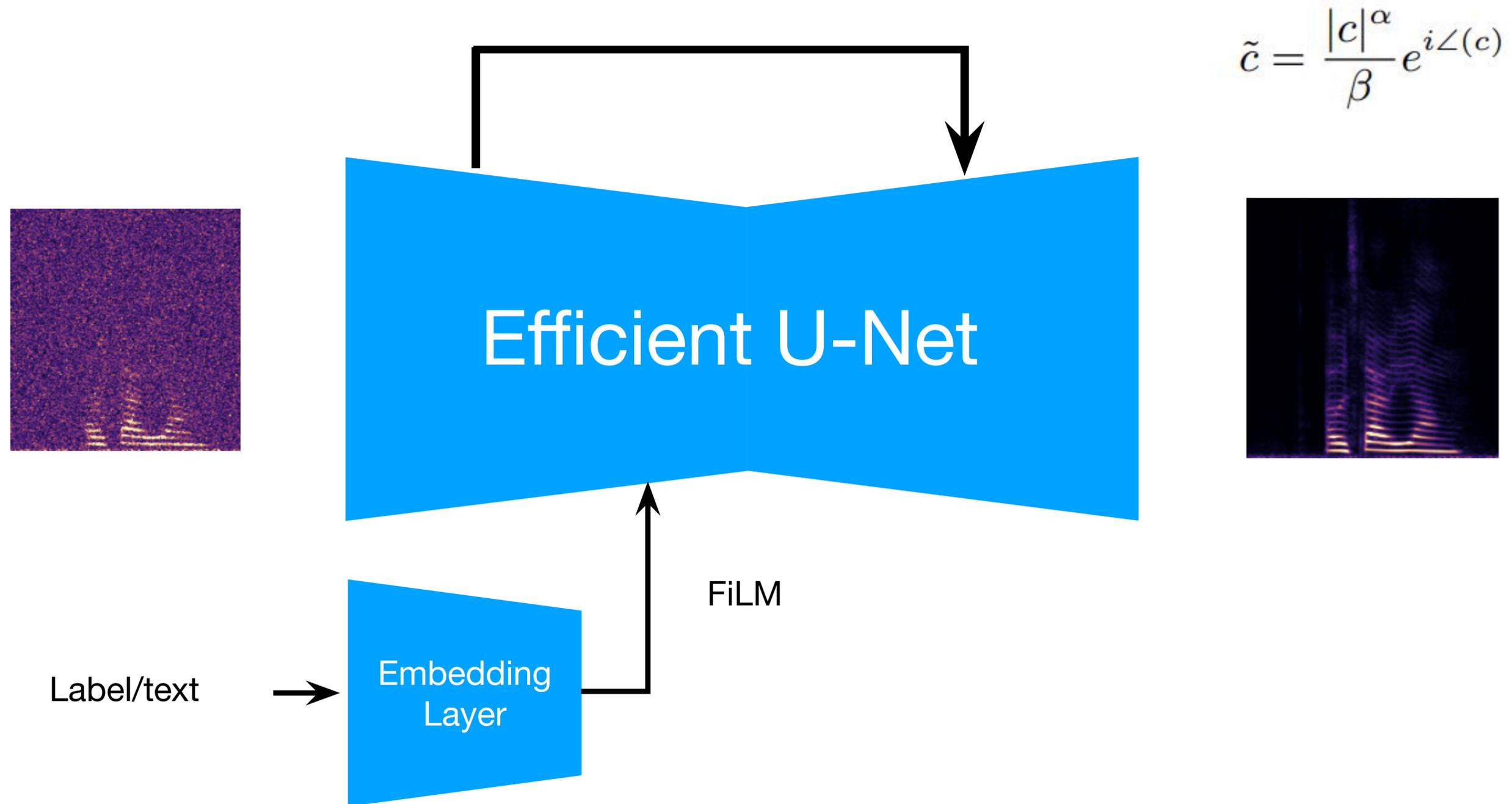
Discretization:

$$\hat{\mathbf{x}}_{t-\Delta t} = \Psi(t - \Delta t, t) \hat{\mathbf{x}}_t + \left[\int_t^{t-\Delta t} -\frac{1}{2} \Psi(t - \Delta t, \tau) \mathbf{G}_\tau \mathbf{G}_\tau^T d\tau \right] \mathbf{s}_\theta(\hat{\mathbf{x}}_t, t)$$

Ingredient 2: $\varepsilon_\theta(x, t)$ over $\mathbf{s}_\theta(x, t)$

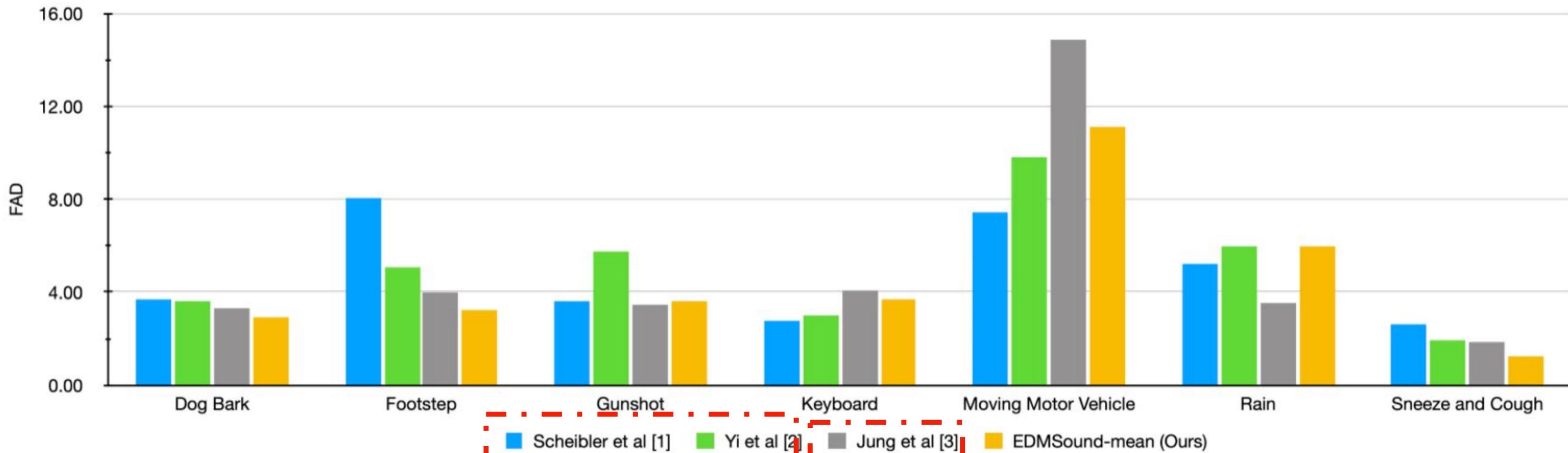
Ingredient 3: Polynomial extrapolation of ε_θ

Neural Network Architecture



Richter, Julius, et al. "Speech enhancement and dereverberation with diffusion-based generative models." IEEE/ACM Transactions on Audio, Speech, and Language Processing (2023).
Saharia, Chitwan, et al. "Photorealistic text-to-image diffusion models with deep language understanding." Advances in Neural Information Processing Systems 35 (2022): 36479-36494.

Fréchet Audio Distance on DCASE2023 Challenge Task7

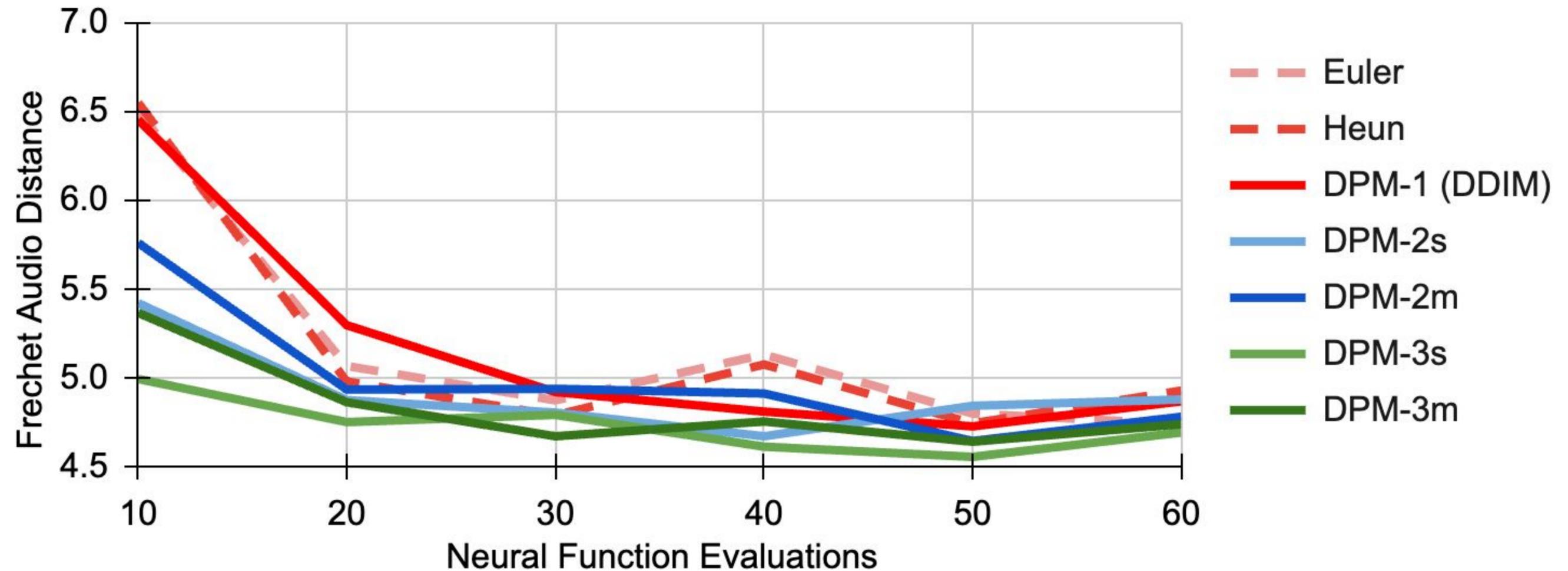


Open dataset
latent diffusion+HiFiGAN

Closed dataset
GAN+HiFiGAN

End-to-end Diffusion

Comparison of different samplers on DCASE2023 Task 7

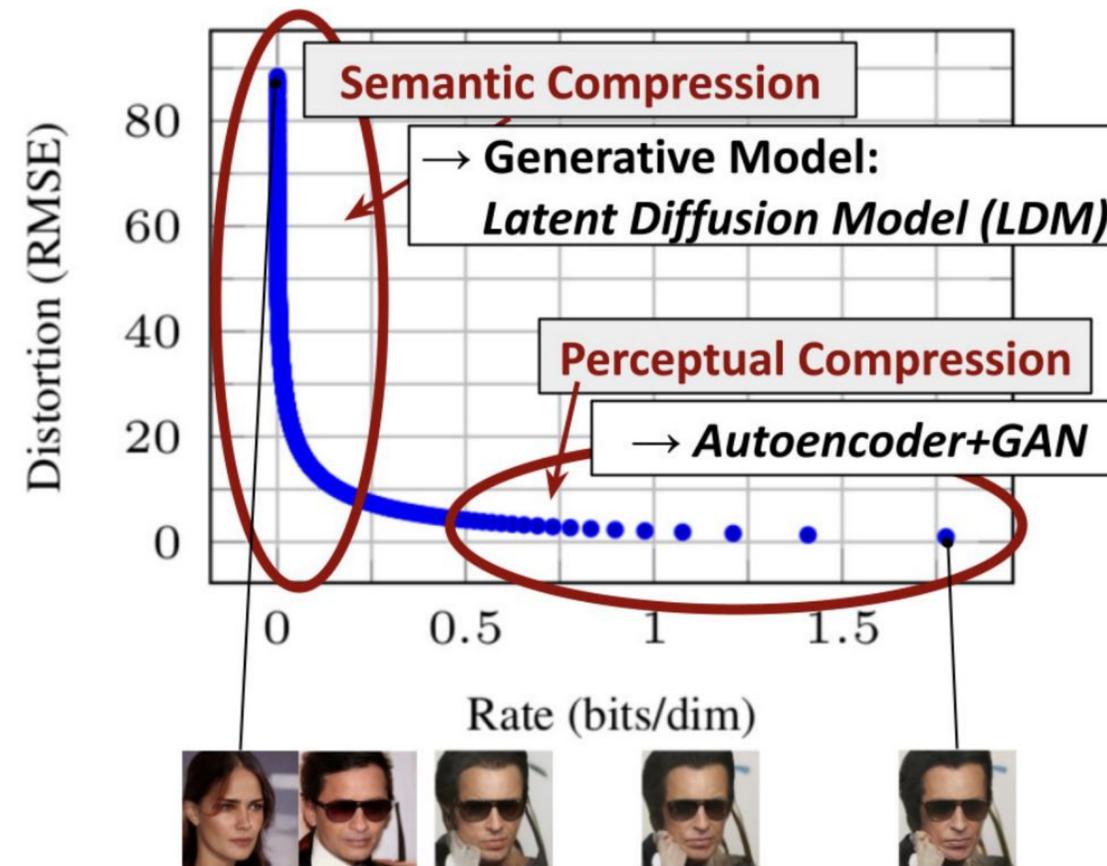


Lu, Cheng, et al. "Dpm-solver: A fast ode solver for diffusion probabilistic model sampling in around 10 steps." *Advances in Neural Information Processing Systems* 35 (2022): 5775-5787.

Lu, Cheng, et al. "Dpm-solver++: Fast solver for guided sampling of diffusion probabilistic models." *arXiv preprint arXiv:2211.01095* (2022).

Limitations on generation

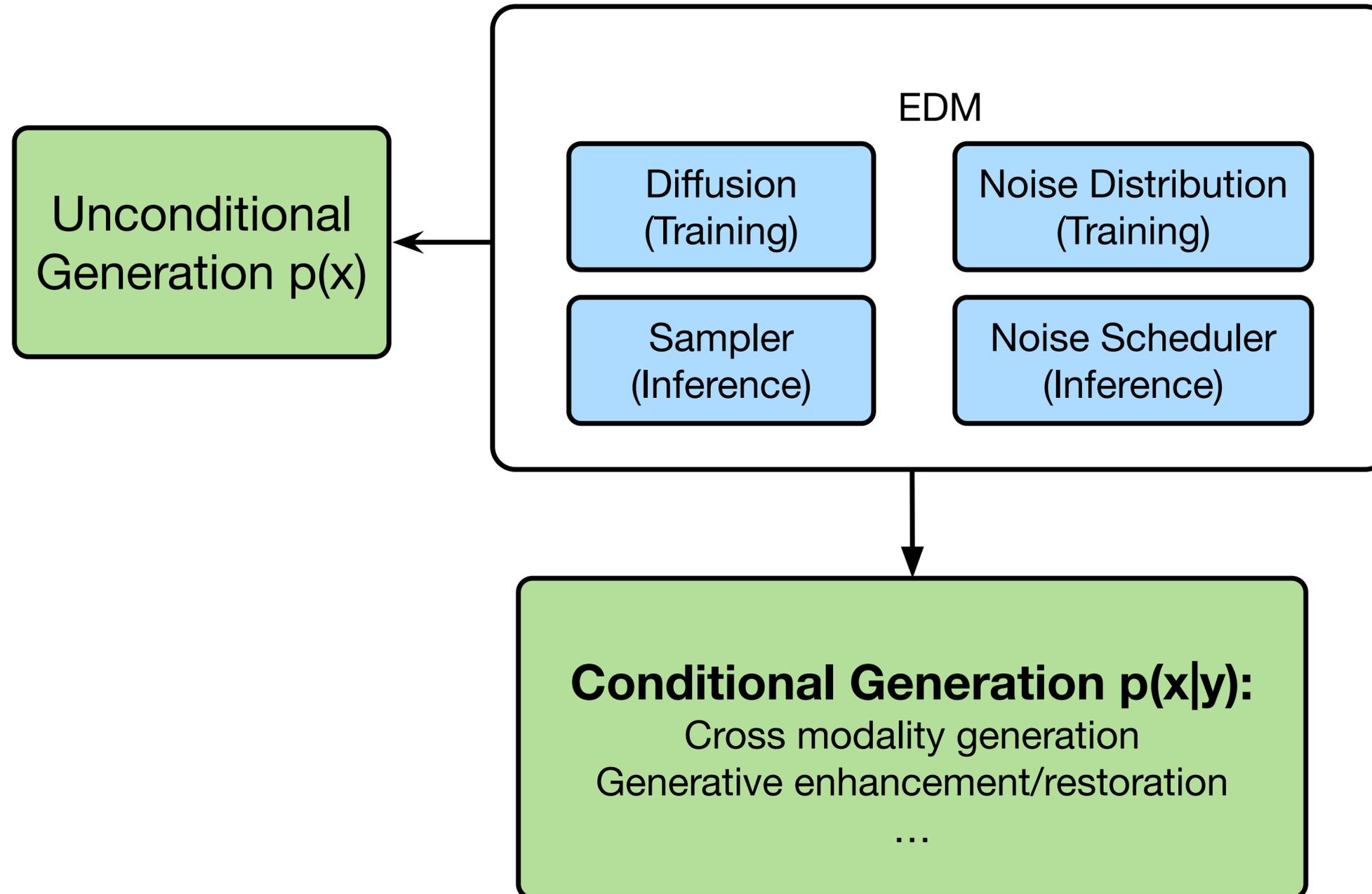
- Efficient samplers may not be as good as training-based method*
- Long context with rich semantics



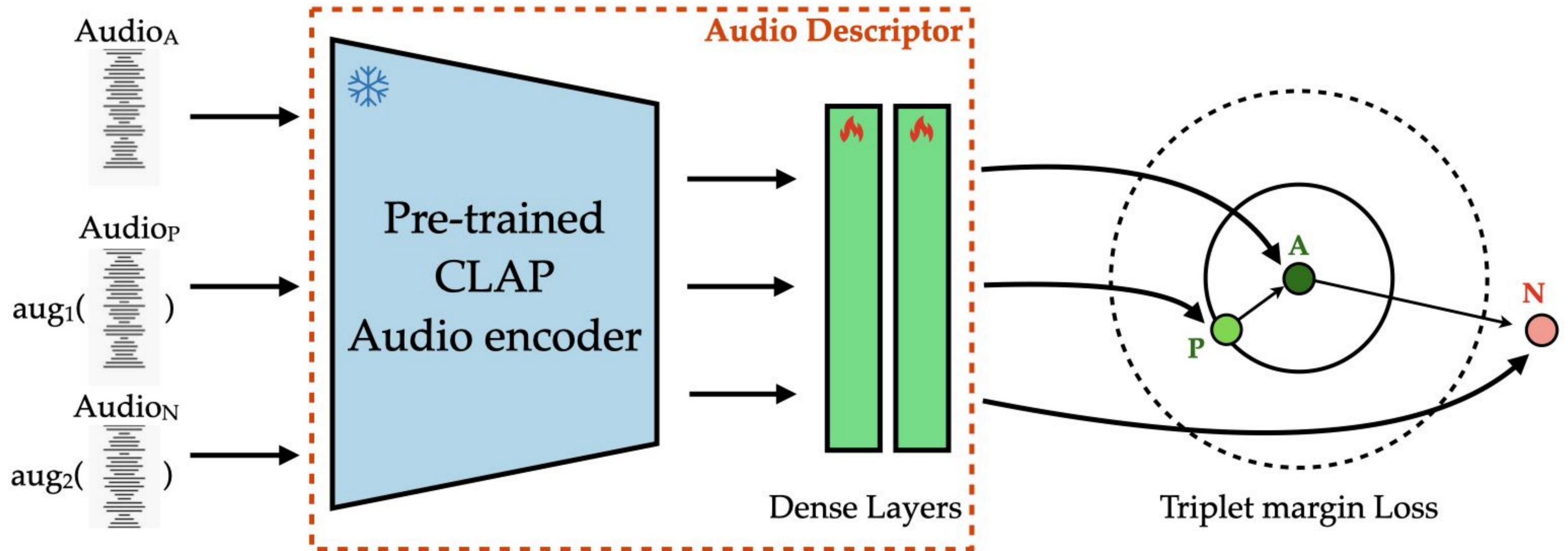
*Zheng, Kaiwen, et al. "DPM-Solver-v3: Improved Diffusion ODE Solver with Empirical Model Statistics." 37 Conference on Neural Information Processing Systems (2023).

Rombach, Robin, et al. "High-resolution image synthesis with latent diffusion models." *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 2022.

Reproducibility

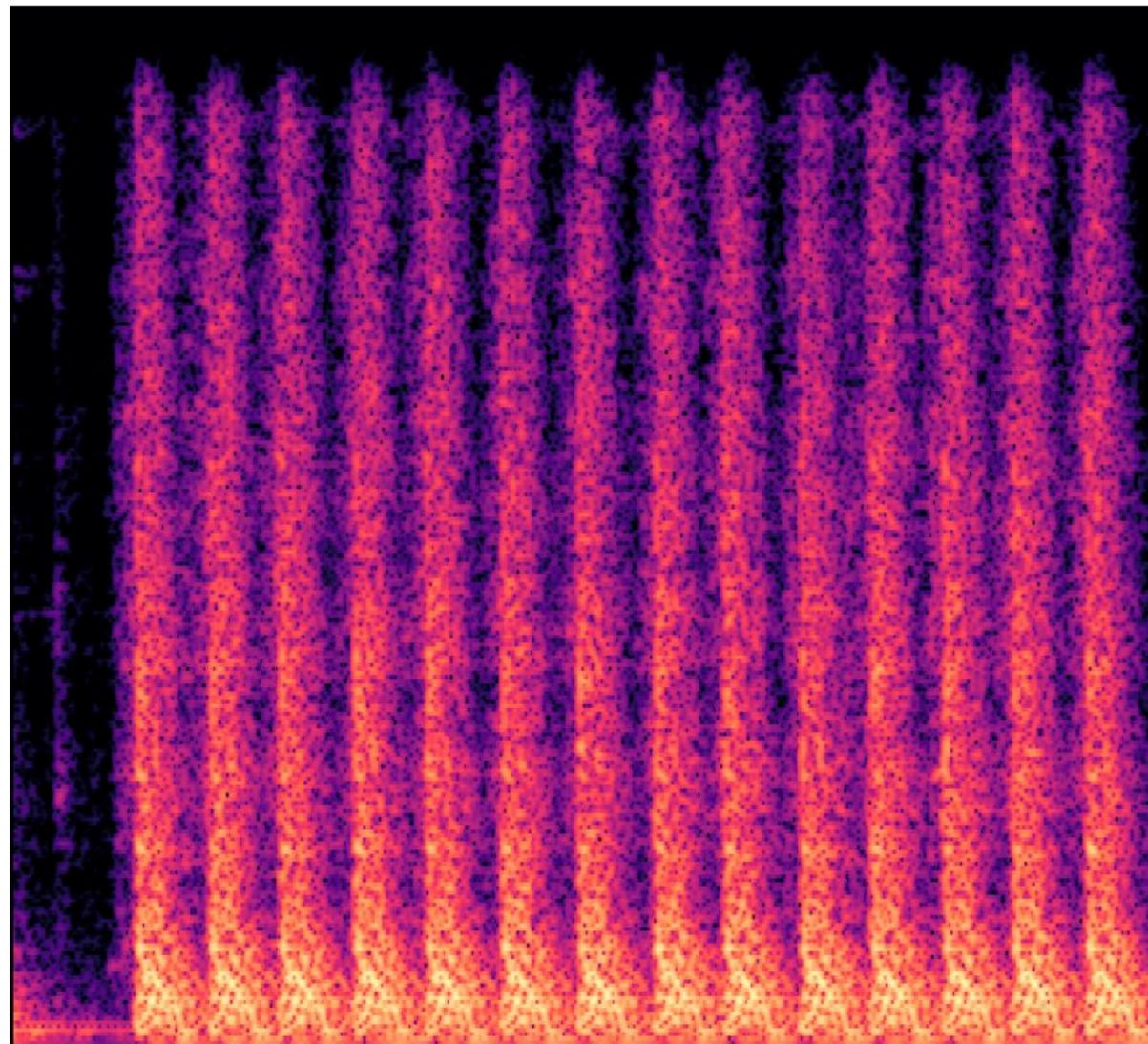


Copy detection (Memorization) in audio diffusion models

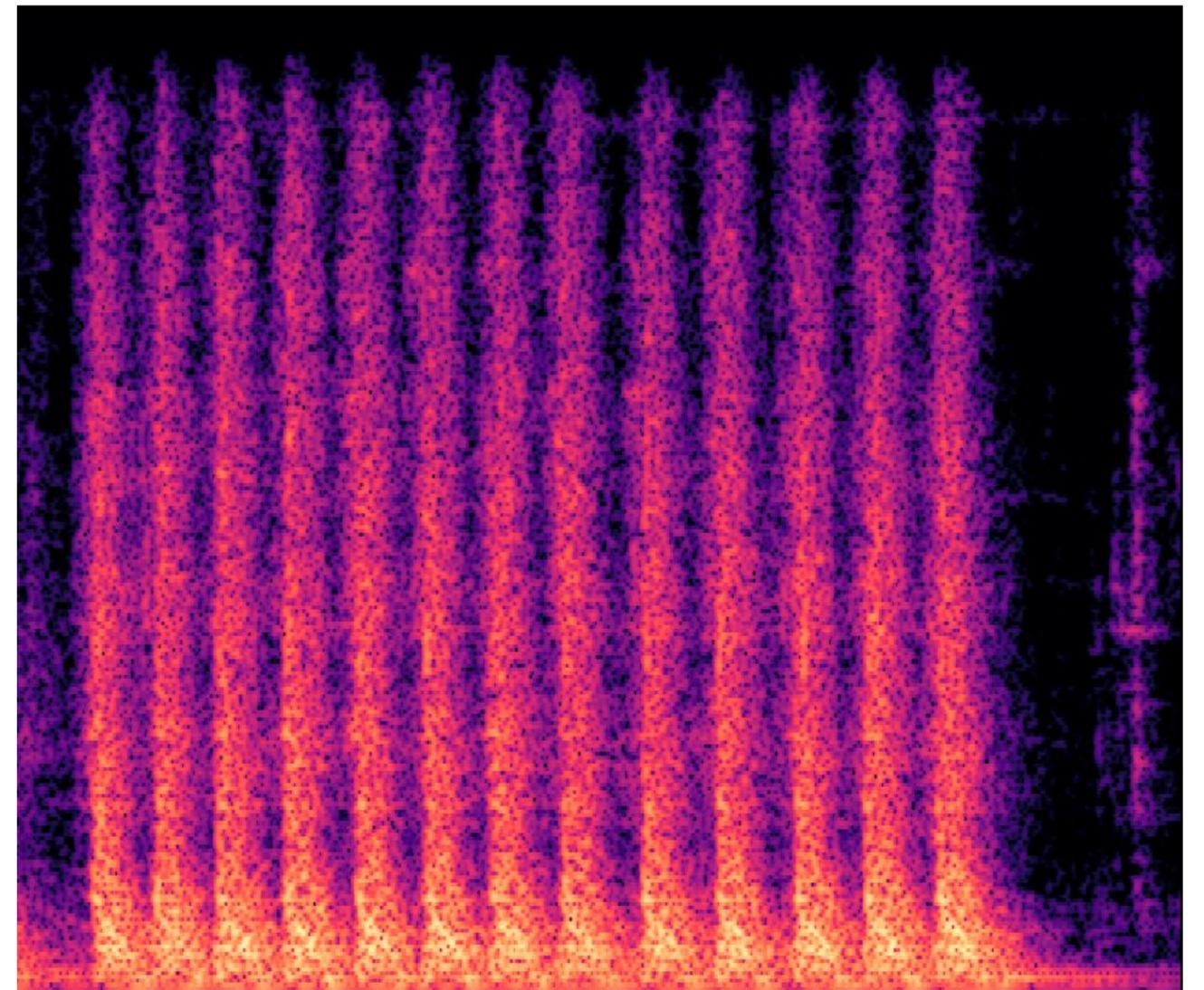


Most similar samples generated by EDMSound

<https://agentcooper2002.github.io/EDMSound/>



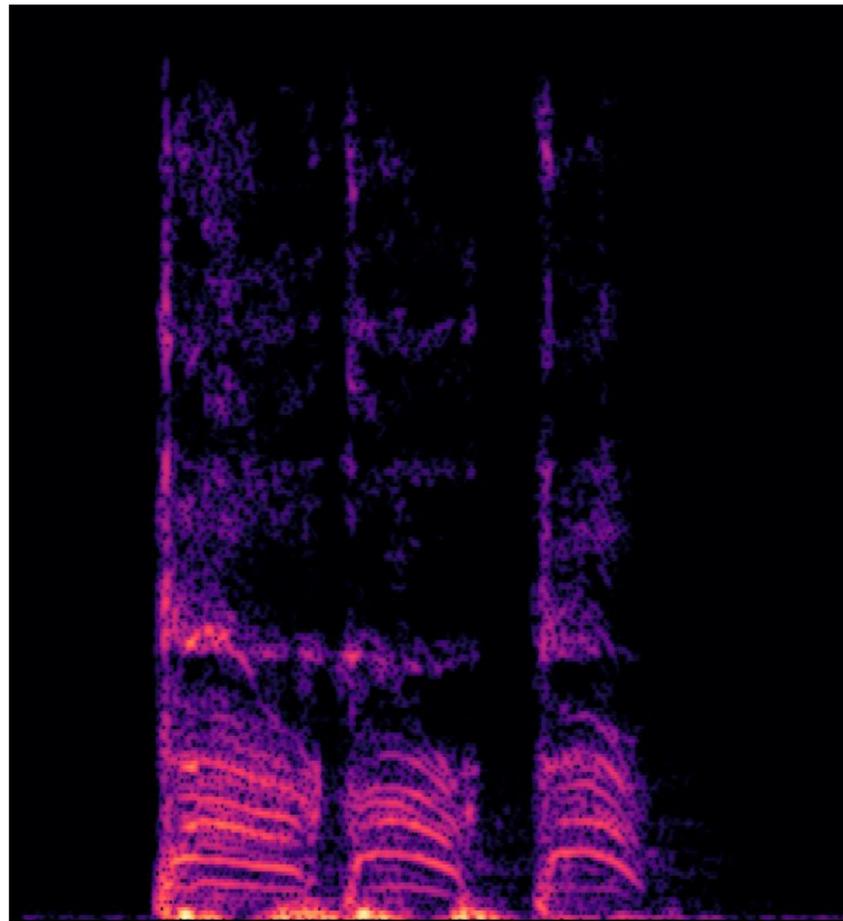
Sample from training set



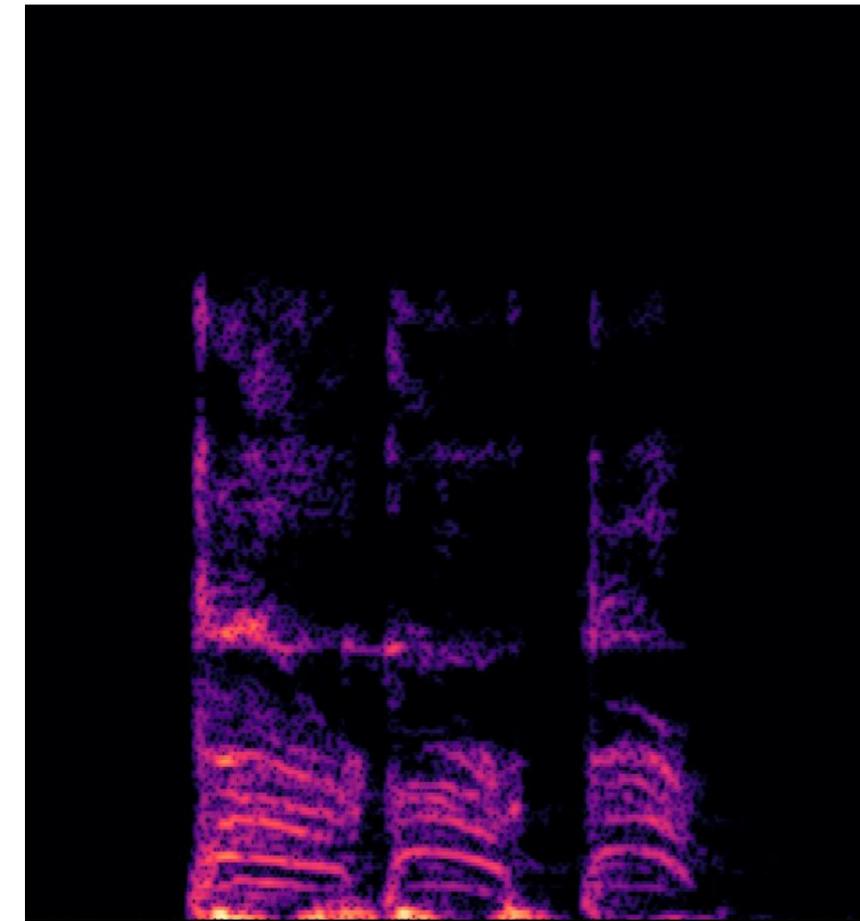
Generated Sample

Most similar samples generated by fine-tuned TANGO

<https://agentcooper2002.github.io/EDMSound/>



Sample from training set



Generated Sample

Ghosal, Deepanway, et al. "Text-to-Audio Generation using Instruction-Tuned LLM and Latent Diffusion Model." arXiv preprint arXiv:2304.13731 (2023).

R. Scheibler, et al. "Class-conditioned latent diffusion model for dcase 2023 foley sound synthesis challenge." Technical report, Tech. Rep., June, 2023.

Overview

- EDMSound: Spectrogram based diffusion models
 - Training: Elucidated diffusion model (EDM) framework
 - Sampling: Exponential-integrator based deterministic solver
- Copy detection/memorization issue in diffusion models
 - Fine-tuned CLAP for audio