

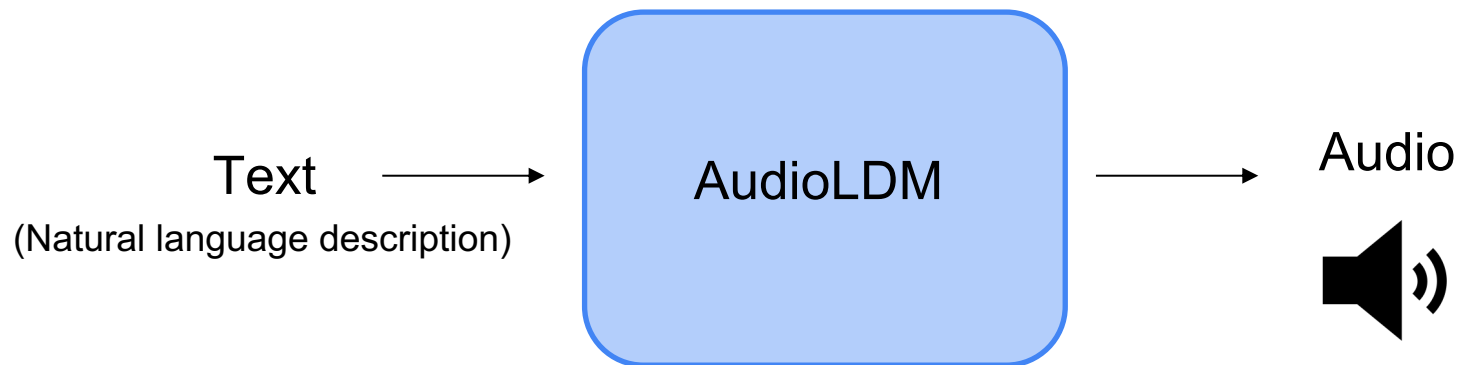
AudioLDM

Text-to-Audio Generation with
Latent Diffusion Models

By Haohe Liu et al.

Presentation by Yumi Kim
on 12. March 2024

Text-to-audio Generation



Demo

Prompt : sound of birds in the forest



Demo

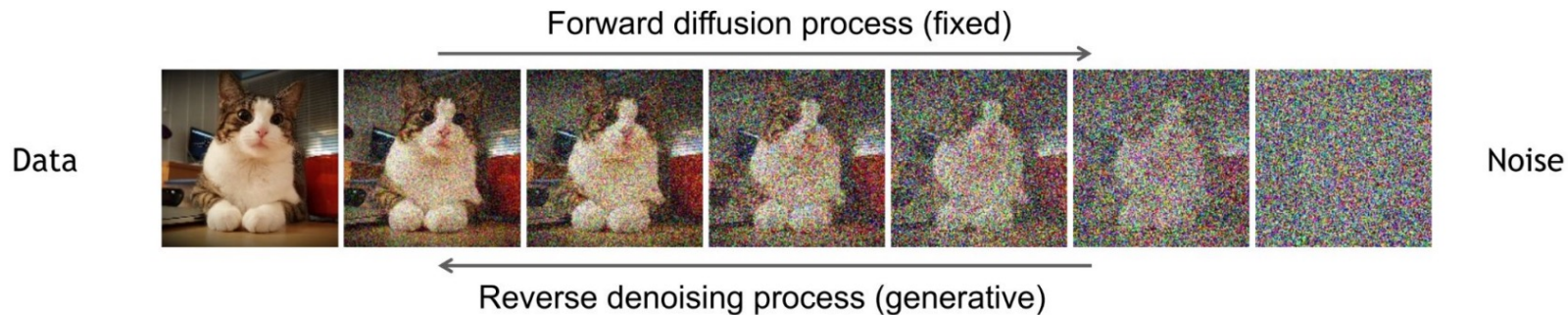
*Prompt: A woman in pointed high heels walking **on a cracked wooden floor***



*Prompt: A woman with pointed high heels walking **down a large corridor***

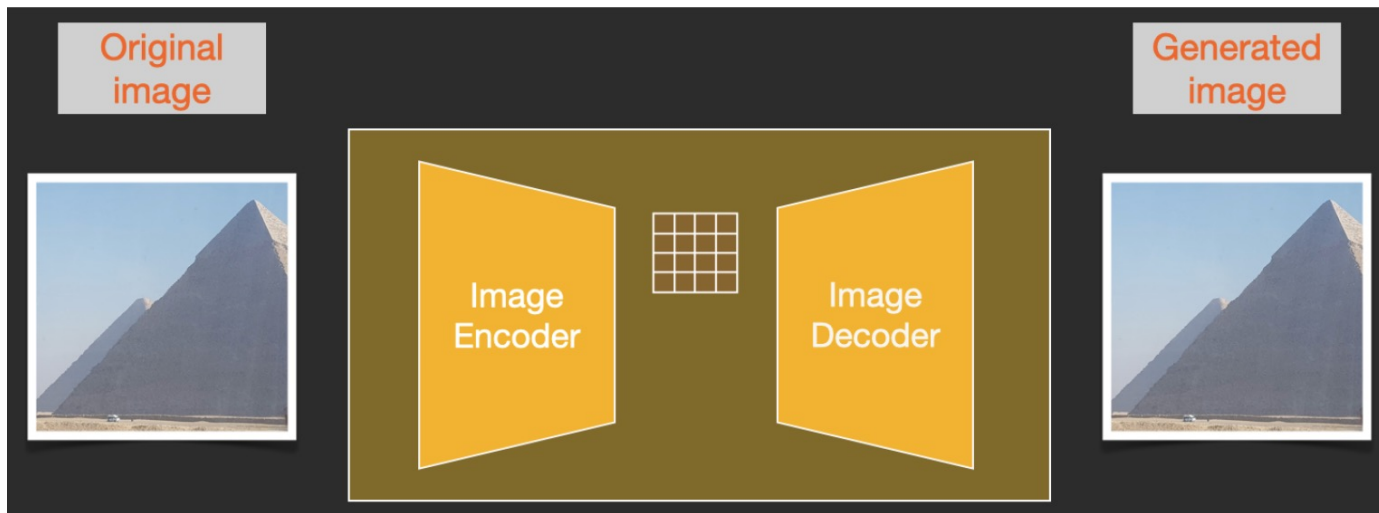


Diffusion model

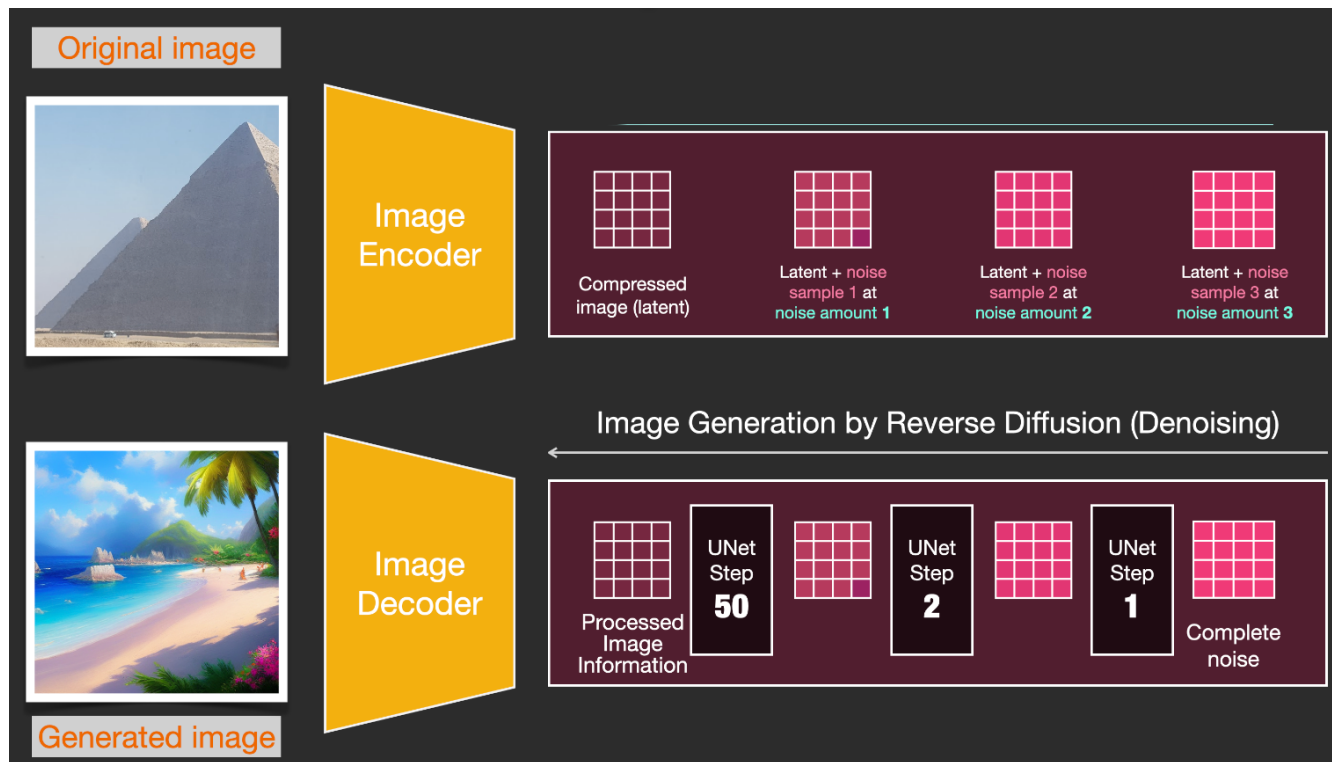


Latent Diffusion Model (LDM)

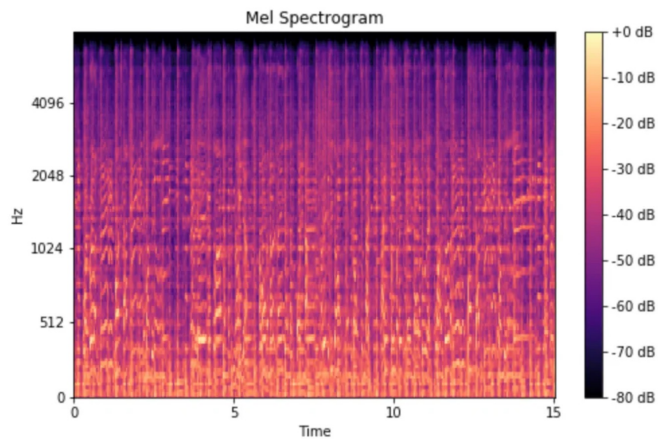
By Stable Diffusion(Rombach et al. 2022)



Latent Diffusion Model

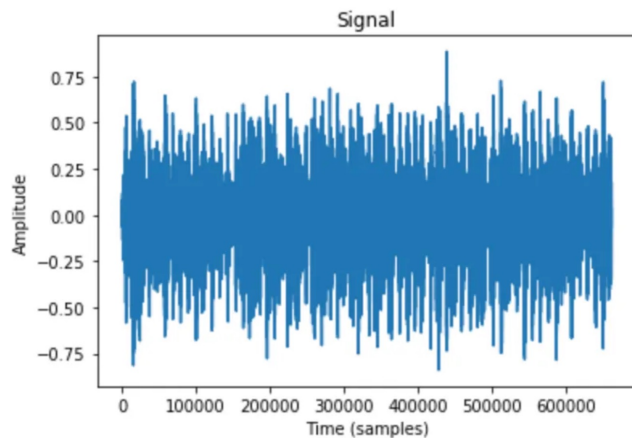


Mel Spectrogram



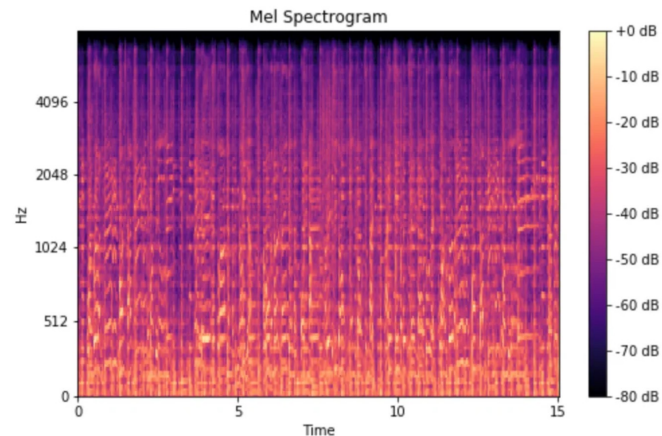
Mel Spectrogram

Mel Spectrogram



Waveform

Short FFT
→
Mel Scale

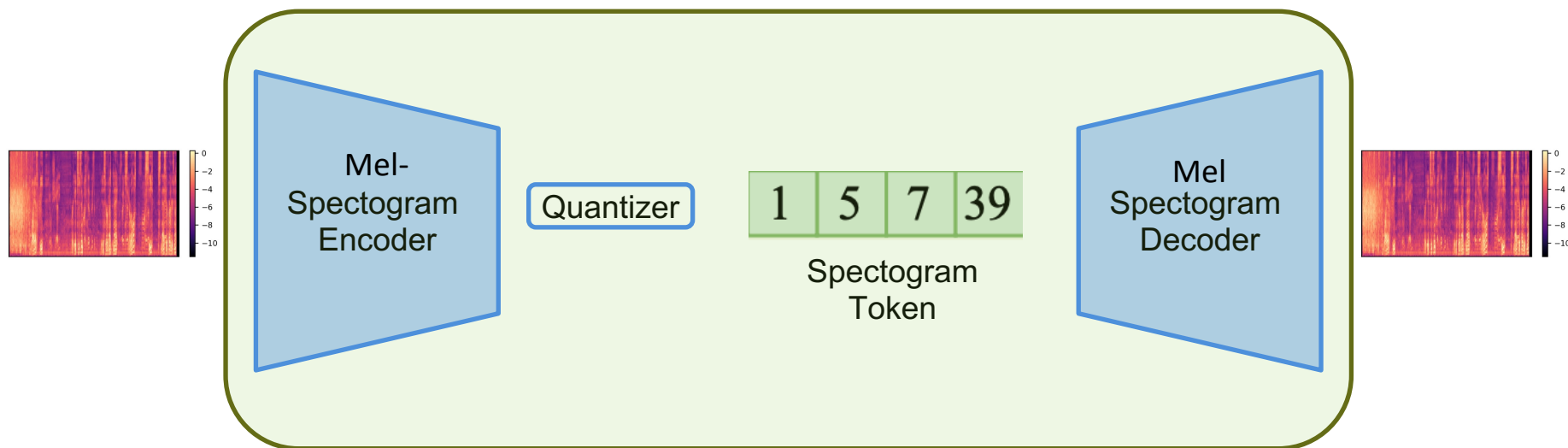


Mel Spectrogram

Related Work

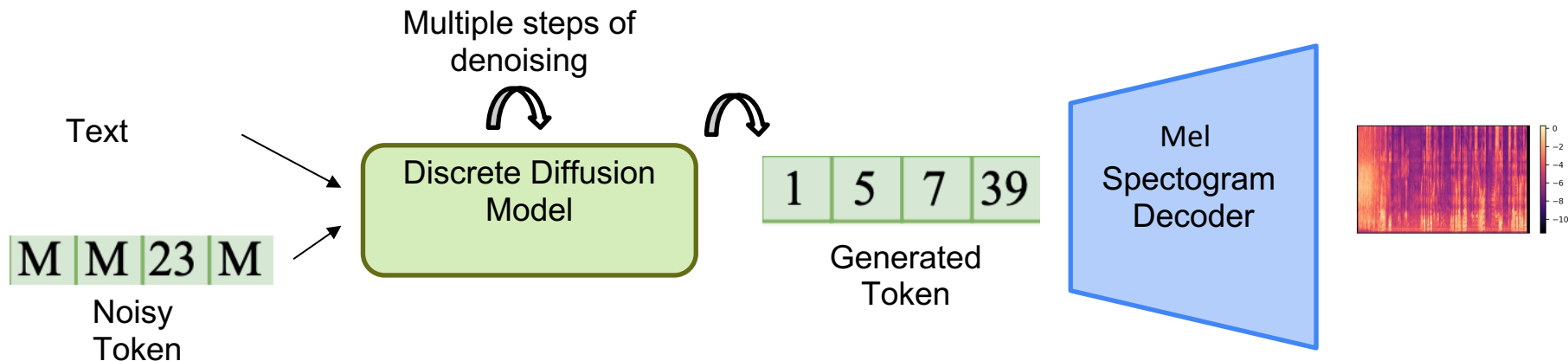
- text-to-audio generation system : DiffSound (*Yang et.al 2021*)

VQ- VAE



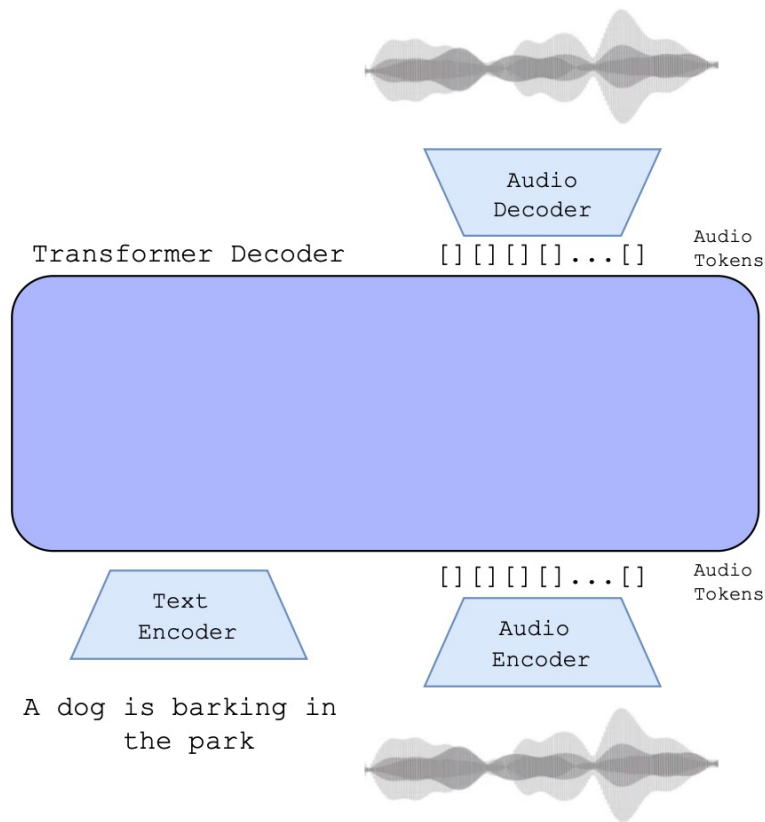
Related Work

- non-autoregressive token-decoder based on the discrete diffusion model



Related Work

- Text to audio Generation: AudioGen
 - **Discrete** audio tokens
 - **Autoregressive** Transformer-based decoder
 - generates Audio tokens based on text



Previous works

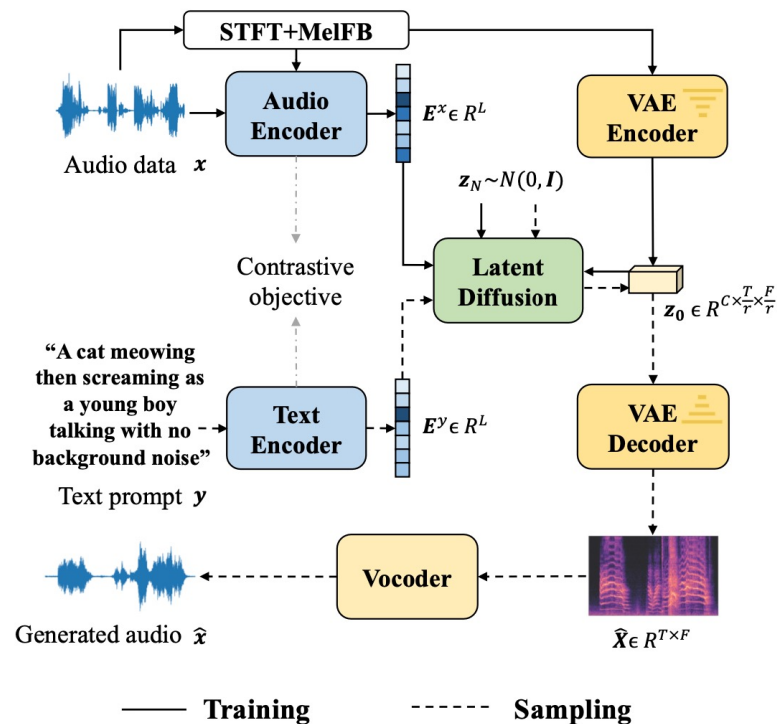
- Discrete tokens
- Paired Audio-text training data
 - Error prone text captions
 - Some captions don't describe the audio well
 - E.g. highly abstracted “Boats : Battleships-5.25 conveyor space”

AudioLDM

- No need for audio-text paired data for LDM training
- Continuous audio representation
- text-guided Audio manipulation in zero-shot fashion
 - Style Transfer and Inpainting
- Computationally efficient
 - trained on 1 GPU

Key components

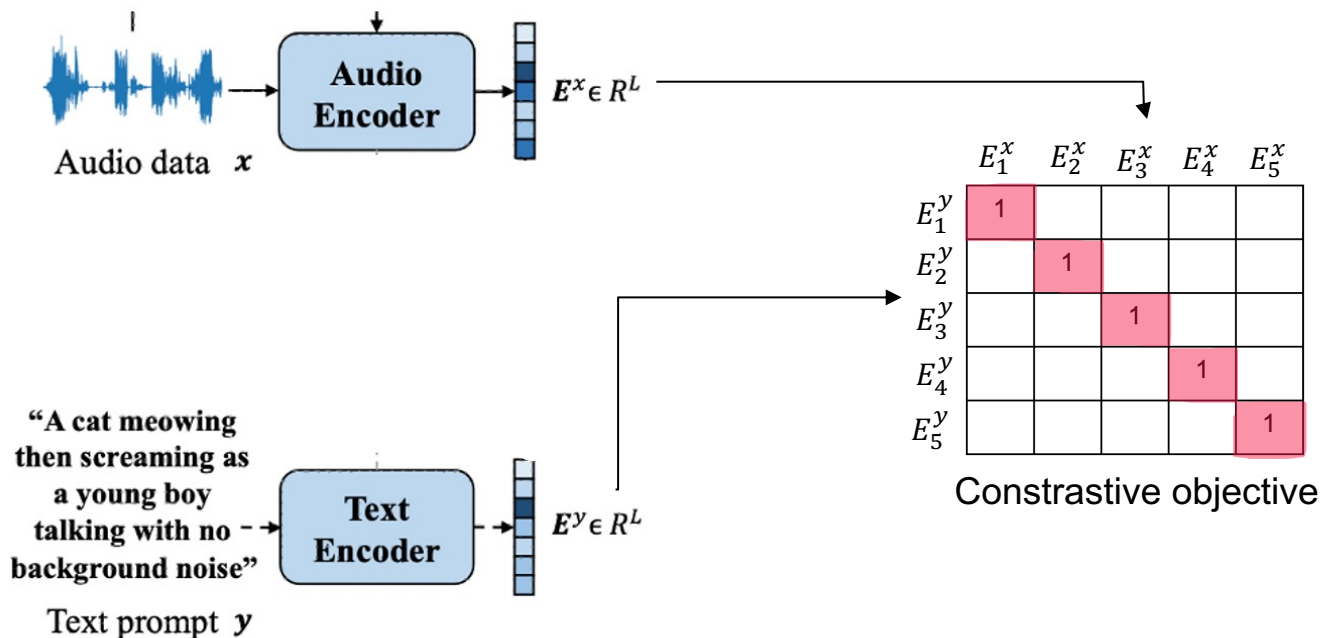
1. Contrastive Language-Audio Pretraining (CLAP) Encoders
2. Mel-spectrogram VAE
3. Latent Diffusion Models (LDM)
4. Mel-to-Waveform Vocoder



(a) Training and sampling process of AudioLDM

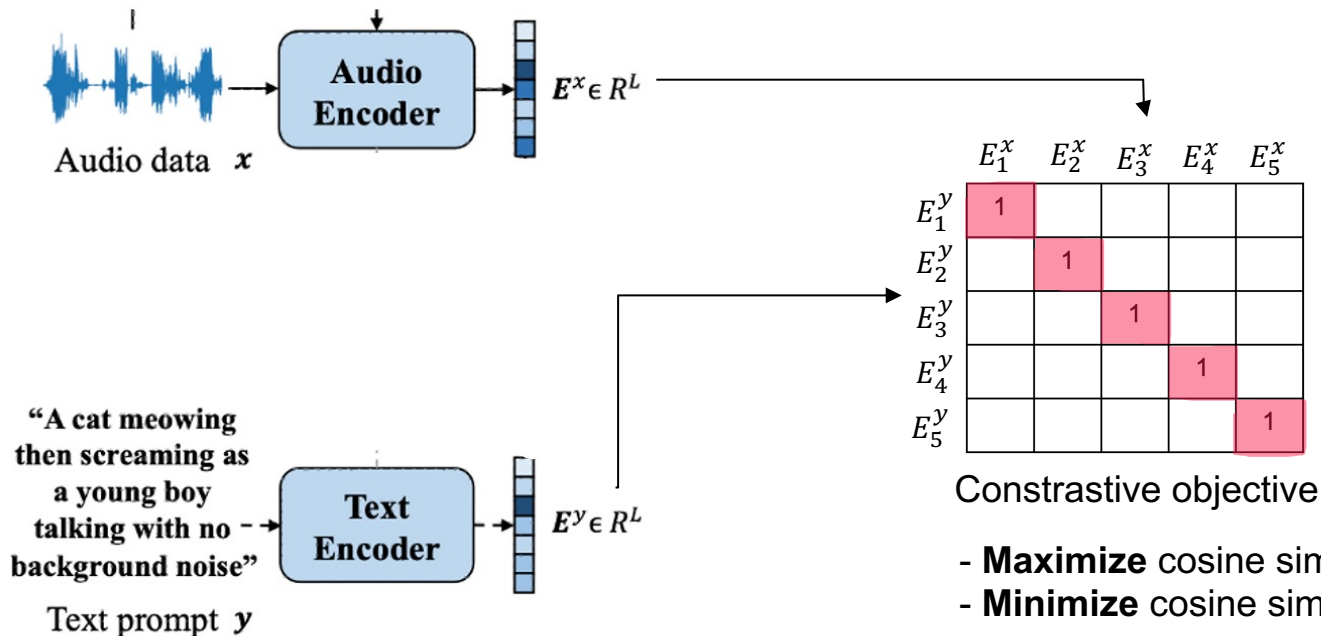
1. Contrastive Language-Audio Pretraining (CLAP)

By Wu et.al

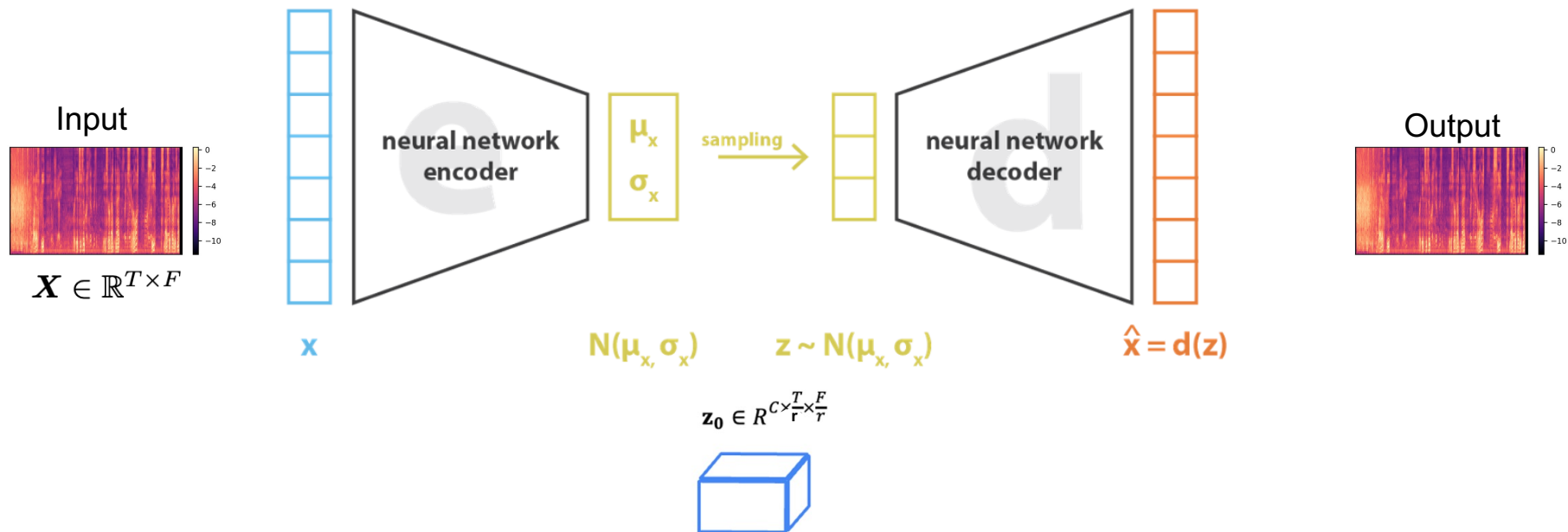


1. Contrastive Language-Audio Pretraining (CLAP)

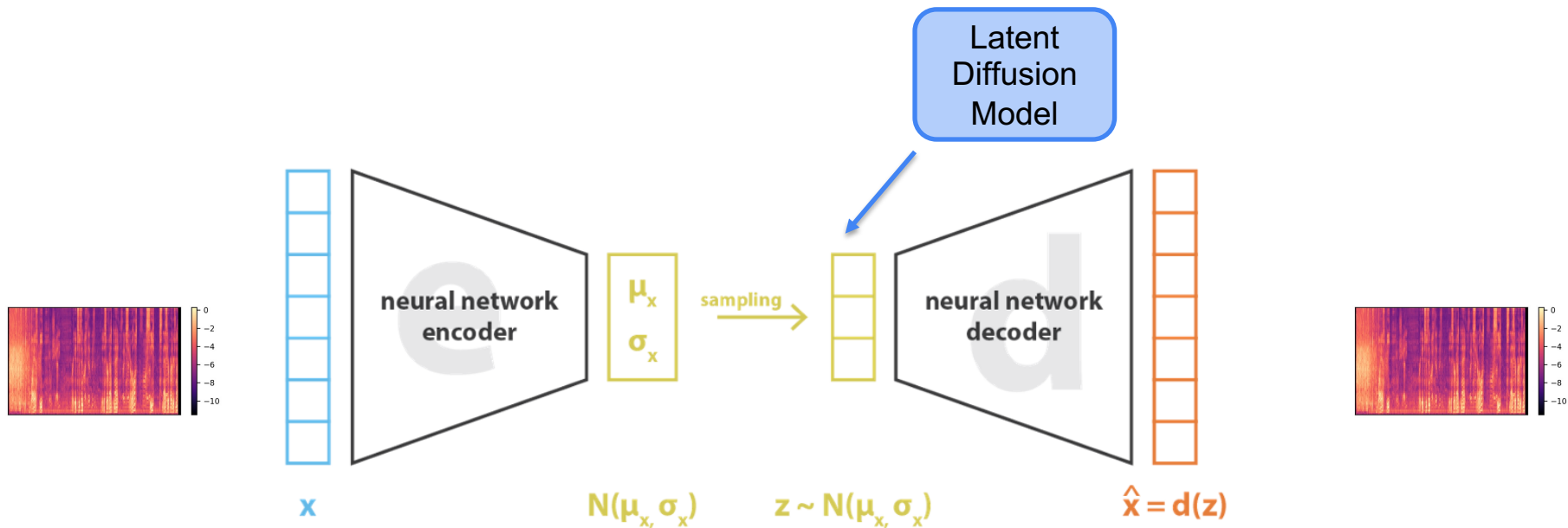
By Wu et.al



Step 2. Training Mel-spectrogram VAE

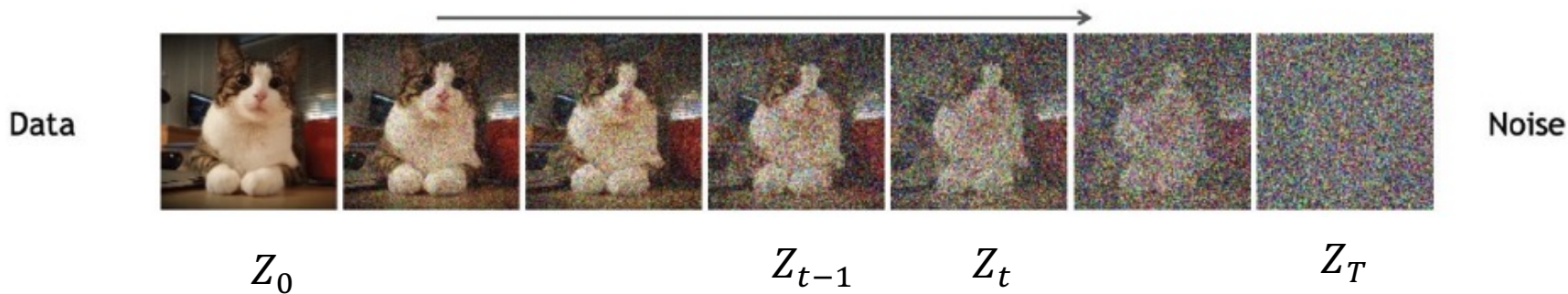


Step 2. Training Mel-spectrogram VAE



Diffusion Model : Forward process

Single Forward step : $q(\mathbf{z}_t | \mathbf{z}_{t-1}) = \mathcal{N}(\mathbf{z}_t; \sqrt{1 - \beta_t}\mathbf{z}_{t-1}, \beta_t\mathbf{I})$



- β_t : predefined noise schedule

Forward process

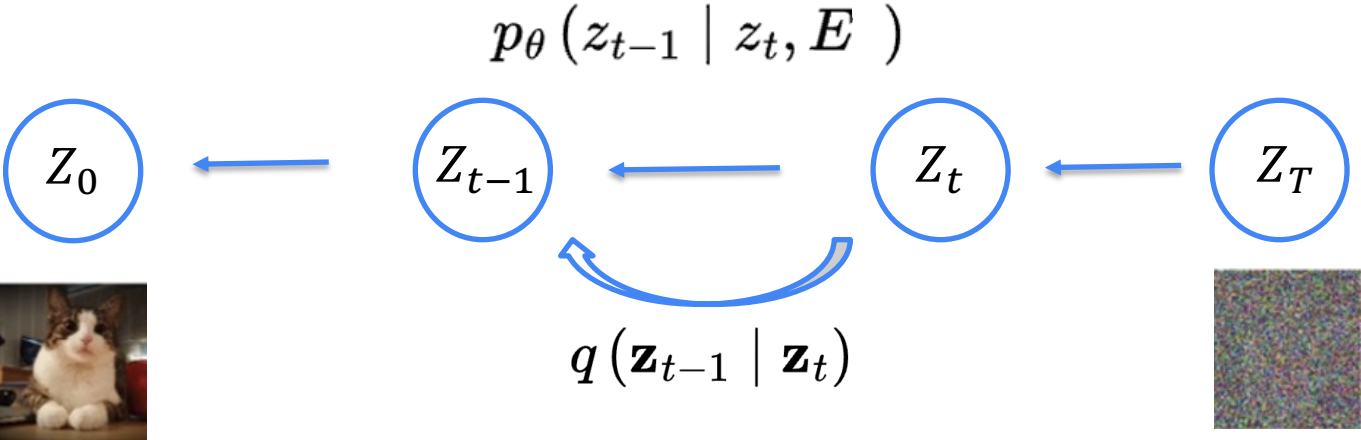
- Diffusion models define a Markov chain of diffusion steps

$$q(\mathbf{z}_{1:T} | \mathbf{z}_0) = \prod_{t=1}^T q(\mathbf{z}_t | \mathbf{z}_{t-1})$$

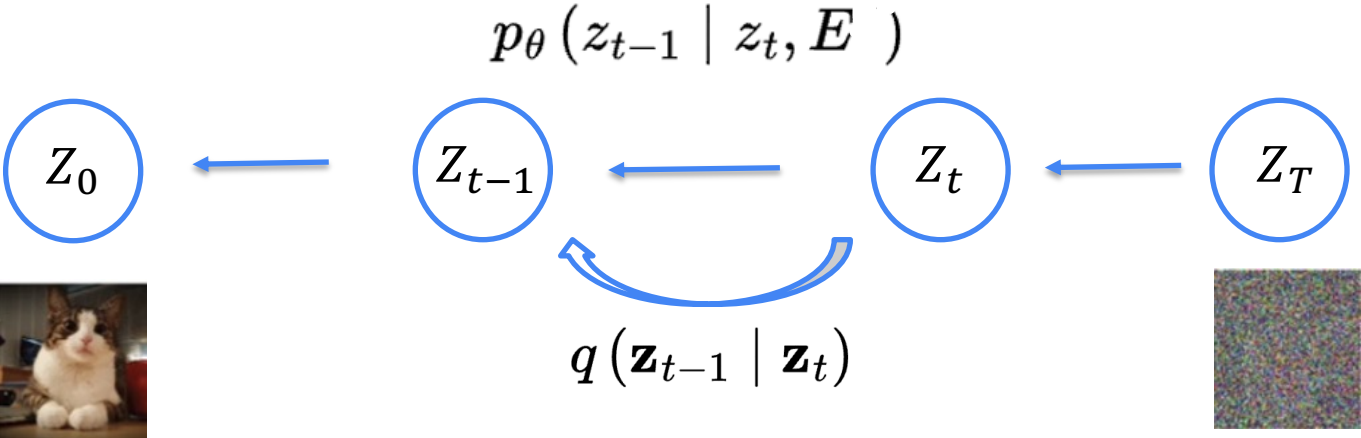
- Using this notation $\alpha_t = 1 - \beta_t$ we can write: $\mathbf{z}_t = \sqrt{\bar{\alpha}_t} \mathbf{z}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon$
 $\bar{\alpha}_t := \prod_{s=1}^t \alpha_s$

ϵ is from $\mathcal{N}(0, \mathbf{I})$

Reverse process



Reverse process



$$p_\theta(z_{0:T} | E^y) := p(z_T) \prod_{t=1}^T p_\theta(z_{t-1} | z_t, E^y)$$

Reverse process

Following Denoising Diffusion Probabilistic Models (DDPM) by Ho et al.

$$p_{\theta}(\mathbf{z}_{t-1} \mid \mathbf{z}_t, \mathbf{E}^y) = \mathcal{N}(\mathbf{z}_{t-1}; \boldsymbol{\mu}_{\theta}(\mathbf{z}_t, t, \mathbf{E}^y), \boldsymbol{\sigma}_n^2 \mathbf{I})$$

$$\boldsymbol{\mu}_{\theta}(\mathbf{z}_t, t, \mathbf{E}^y) = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{z}_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_t \right)$$

Reverse process

Following Denoising Diffusion Probabilistic Models (DDPM) by Ho et al.

$$p_{\theta}(\mathbf{z}_{t-1} \mid \mathbf{z}_t, \mathbf{E}^y) = \mathcal{N}(\mathbf{z}_{t-1}; \boldsymbol{\mu}_{\theta}(\mathbf{z}_t, t, \mathbf{E}^y), \boldsymbol{\sigma}_t^2 \mathbf{I})$$

$$\boldsymbol{\mu}_{\theta}(\mathbf{z}_t, t, \mathbf{E}^y) = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{z}_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_{\theta}(\mathbf{z}_t, t, \mathbf{E}^y) \right)$$

$$\boldsymbol{\sigma}_t^2 = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t \quad \text{No need to learn}$$

Reverse process

Following Denoising Diffusion Probabilistic Models (DDPM) by Ho et al.

$$p_{\theta}(\mathbf{z}_{t-1} \mid \mathbf{z}_t, \mathbf{E}^y) = \mathcal{N}(\mathbf{z}_{t-1}; \boldsymbol{\mu}_{\theta}(\mathbf{z}_t, t, \mathbf{E}^y), \boldsymbol{\sigma}_t^2 \mathbf{I})$$

$$\boldsymbol{\mu}_{\theta}(\mathbf{z}_t, t, \mathbf{E}^y) = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{z}_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_{\theta}(\mathbf{z}_t, t, \mathbf{E}^y) \right)$$

$$\boldsymbol{\sigma}_t^2 = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t \quad \text{No need to learn}$$

Loss function: $L_t(\theta) = \mathbb{E}_{\mathbf{z}_0, \boldsymbol{\epsilon}, t} \left\| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta}(\mathbf{z}_t, t, \mathbf{E}^x) \right\|_2^2$

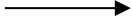
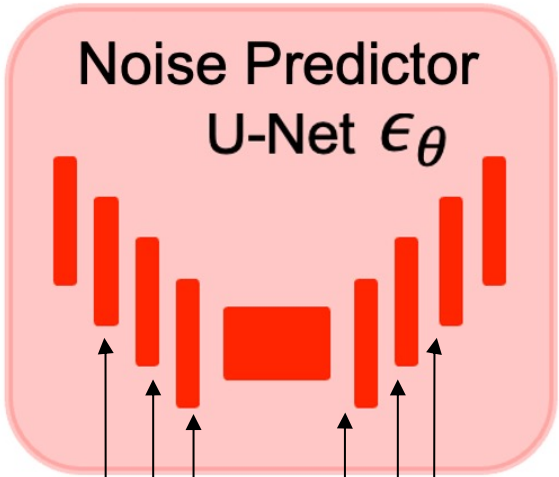
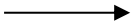
Noise Predictor U-Net



CLAP embedding

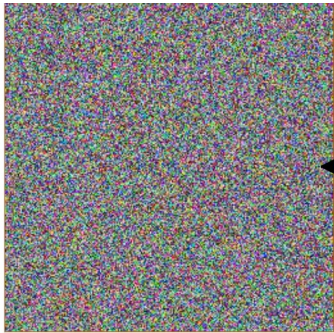


Step embedding for step t

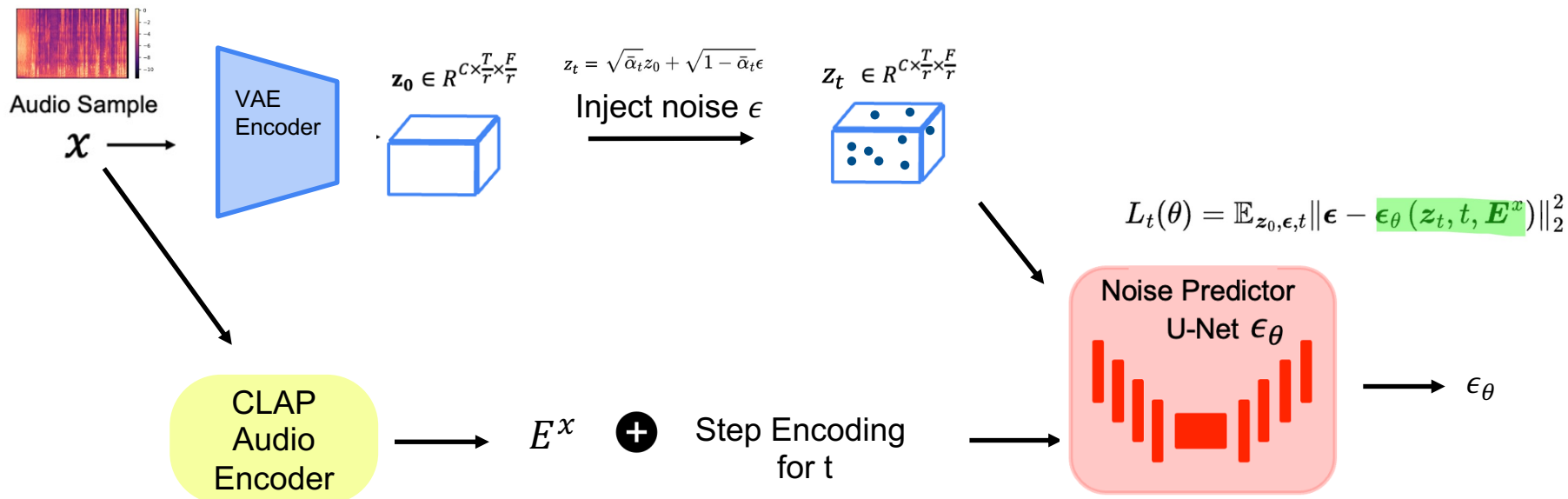


$$\epsilon_\theta(z_t, t, E)$$

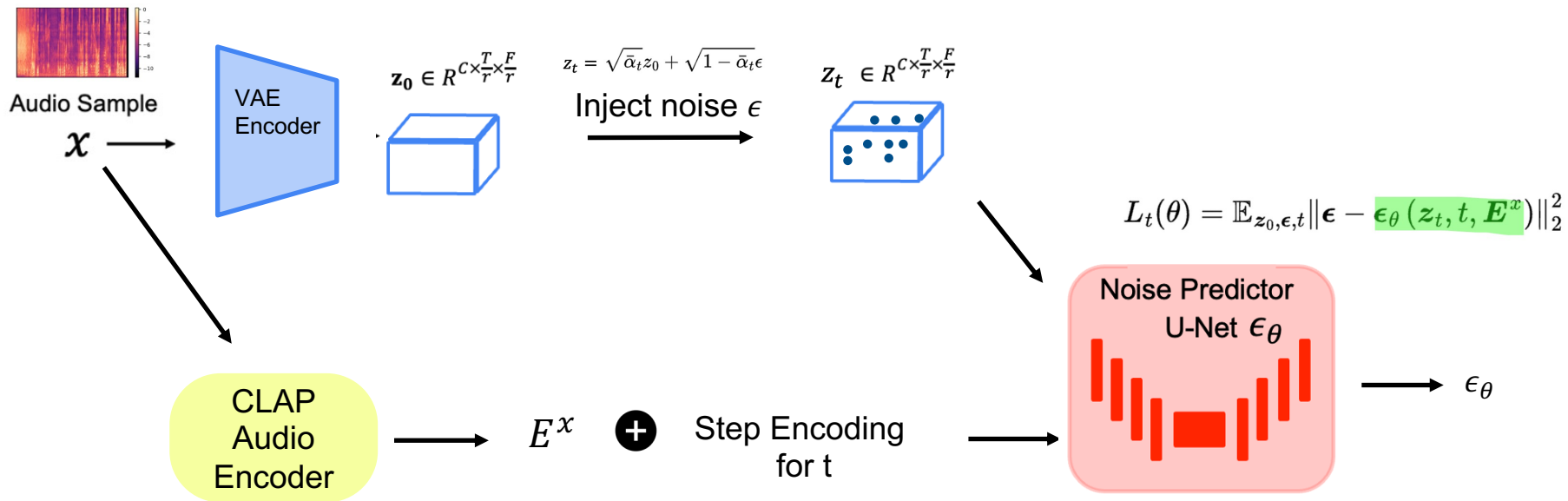
Predicted Noise



Step 3. Training Diffusion model

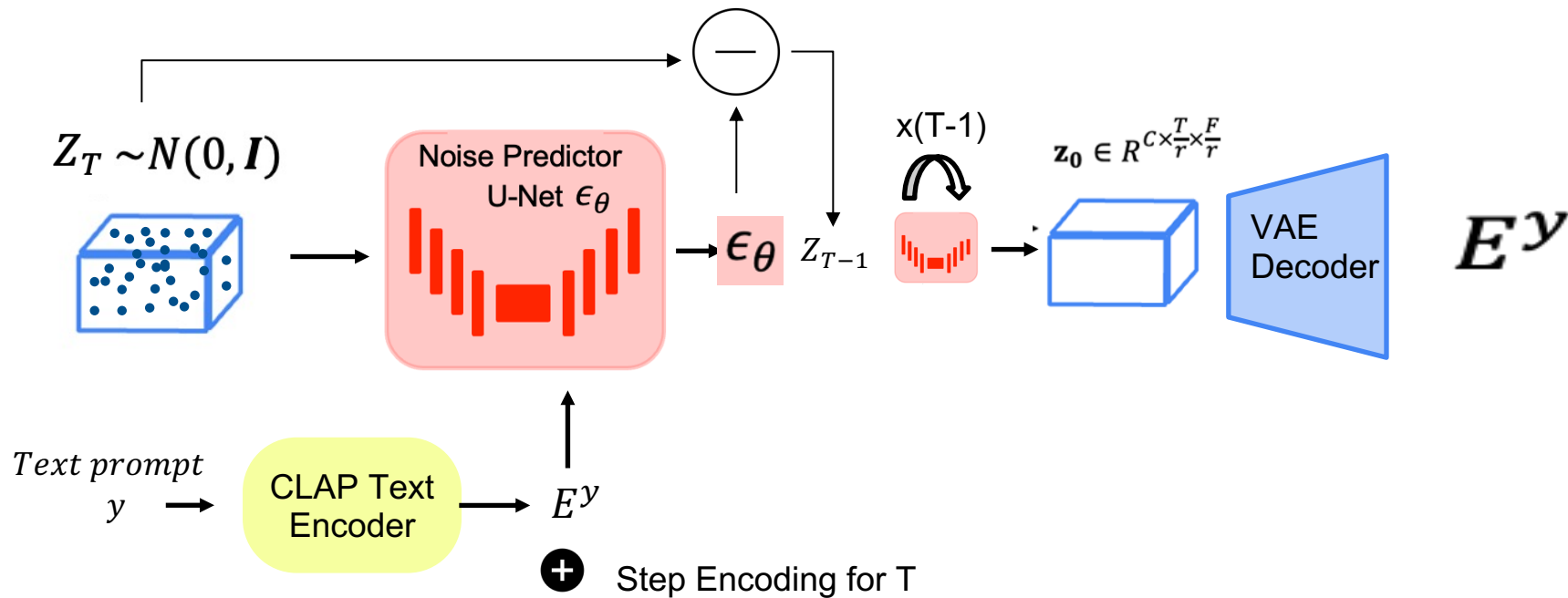


Step 3. Training Diffusion model



No Need for paired audio-text data

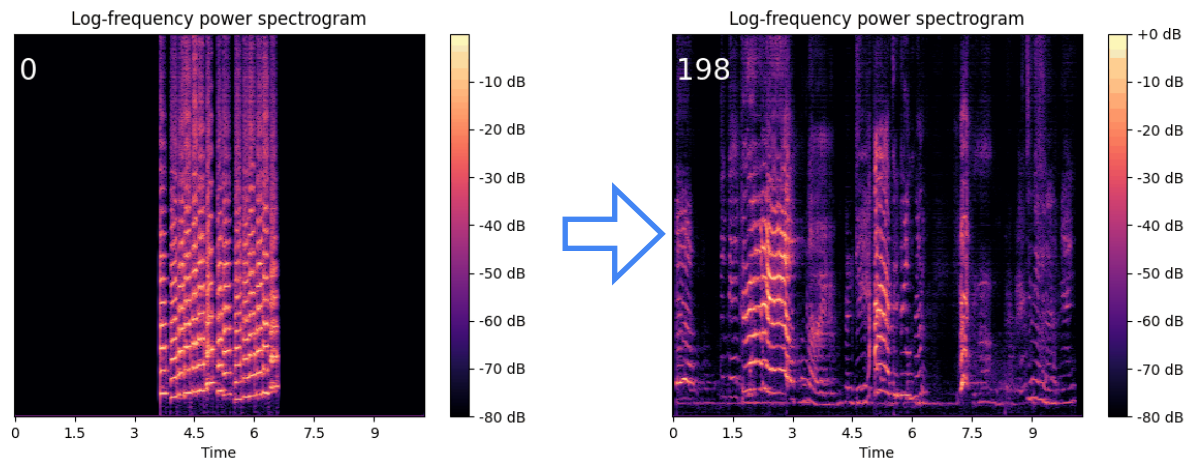
Generation



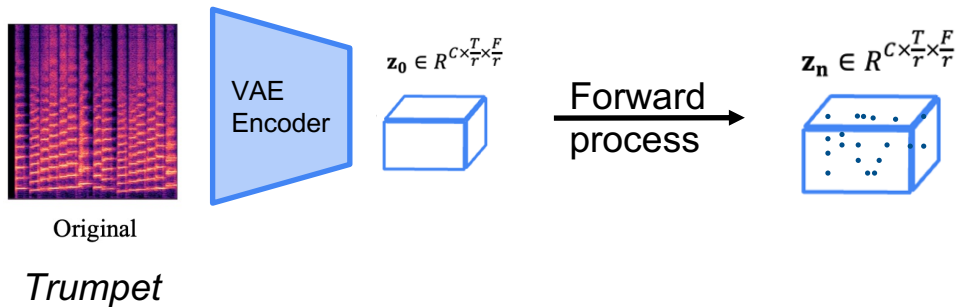
Text-guided audio Manipulation

- Style Transfer

From *trumpet* to *children singing*

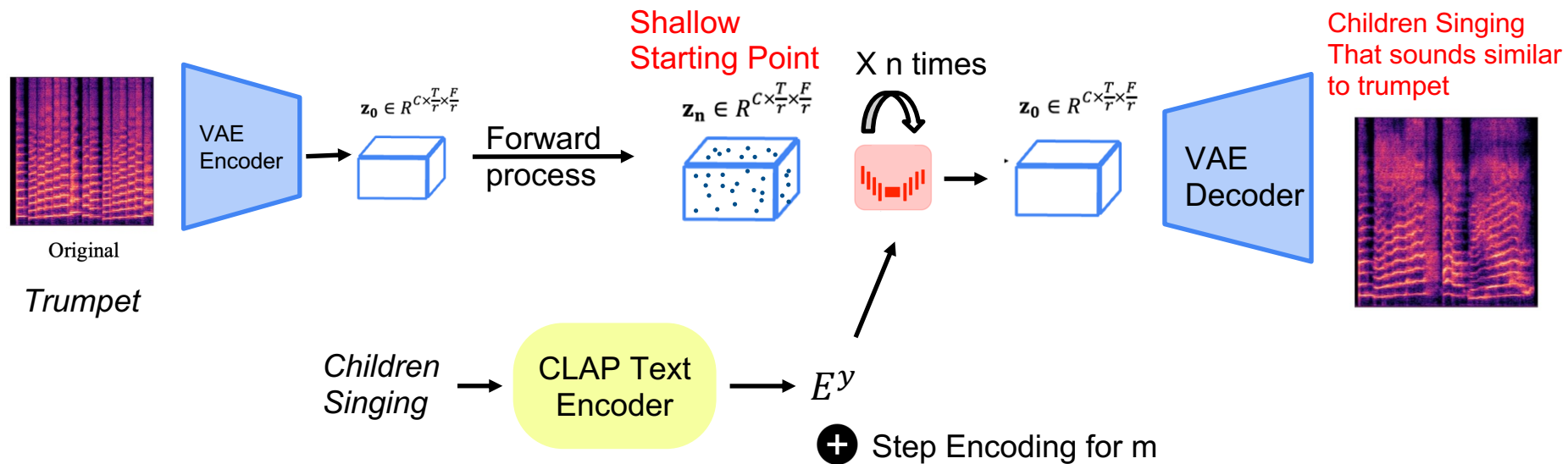


Style Transfer



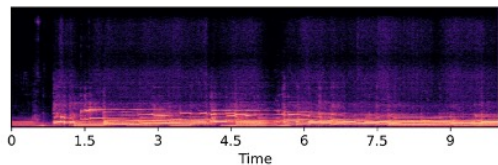
Only do $n < T$ steps and
not the full forward process!

Style Transfer

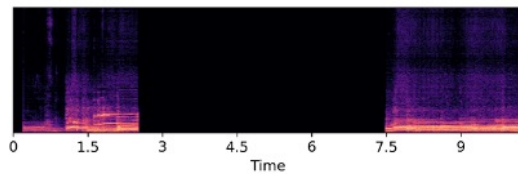


Inpainting

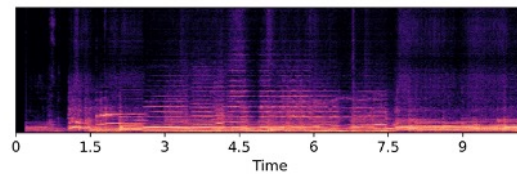
Ground-Truth



Unprocessed

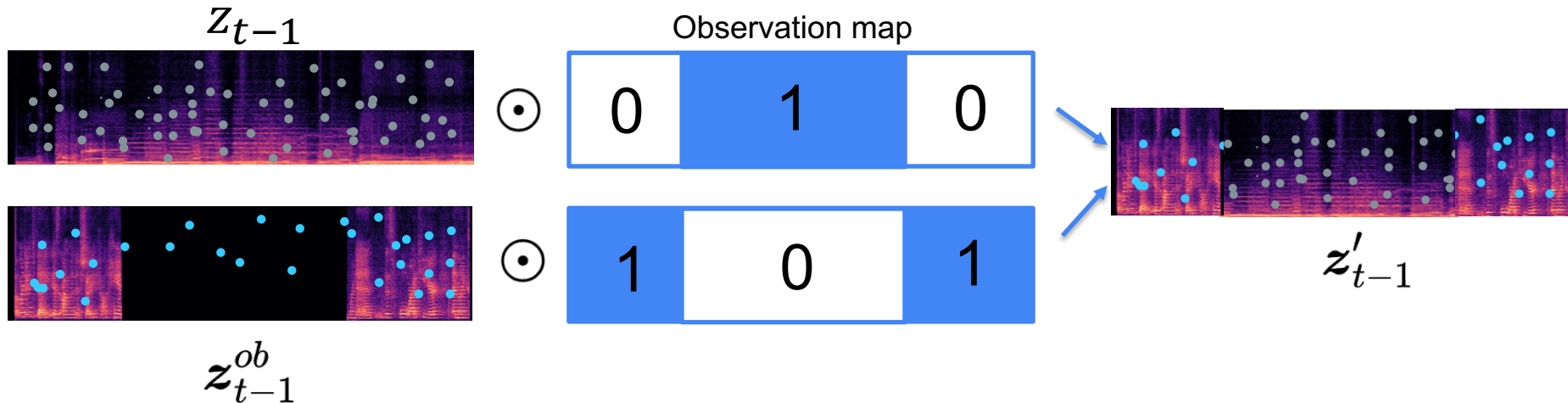


Inpainting result



At each denoising step

- replace generated z_{t-1} with $z'_{t-1} = (1 - m) \odot z_{t-1} + m \odot z_{t-1}^{ob}$



Results

Training Data

- AudioSet
- AudioCaps
- FreeSound
- BBC Sound Effect Library



=> ~9000 hours of audio samples with 16kHz sampling rate

Evaluation Metric

- Objective evaluation
 - FAD: Frechet Audio Distance
 - KL: Kullback-Leibler Divergence
 - IS: Inception Score

Evaluation Metric

- Objective evaluation
 - FAD: Frechet Audio Distance
 - KL: Kullback-Leibler Divergence
 - IS: Inception Score
- Subjective evaluation
 - overall quality (OVL)
 - relevance to the input text (REL)

Evaluation on AudioCaps Dataset

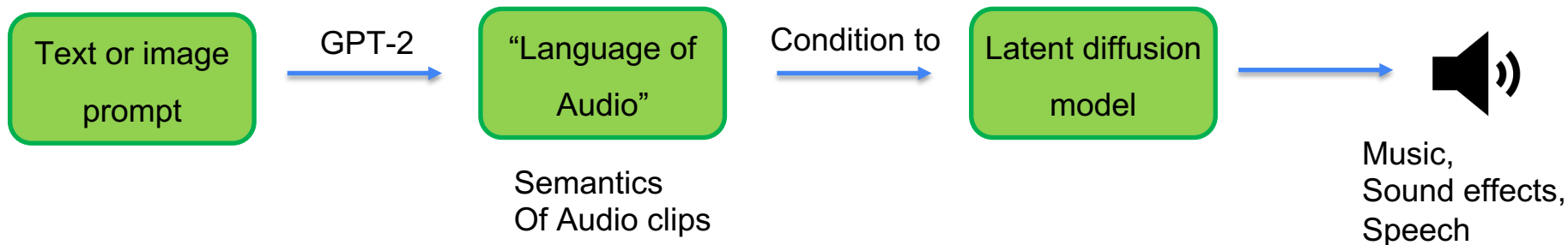
Model	Params	Duration (h)	IS \uparrow	KL \downarrow	FAD \downarrow	OVL \uparrow	REL \uparrow
Ground truth	-	-	-	-	-	83.61 \pm 1.1	80.11 \pm 1.2
DiffSound [†] (Yang et al., 2022)	400M	5420	4.01	2.52	7.75	45.00 \pm 2.6	43.83 \pm 2.3
AudioGen [†] (Kreuk et al., 2022)	285M	8067	-	2.09	3.13	-	-
AudioLDM-L-Full	739M	8886	8.13	1.59	1.96	65.91 \pm 1.0	65.97 \pm 1.6

Limitation

- Low sampling rate (16kHz) of data
- Limited generation quality
 - Incomprehensible speech generation

Future work : AudioLDM 2

- Higher sampling rate (48kHz)
- Improved Text-to-speech generation
- SOTA performance on Music and Sound effect generation
- Image Prompt



Thank you!