# Deep Learning for Vision & Language

Segmentation, AutoEncoders, Variational AutoEncoders, Introduction to Diffusion Models
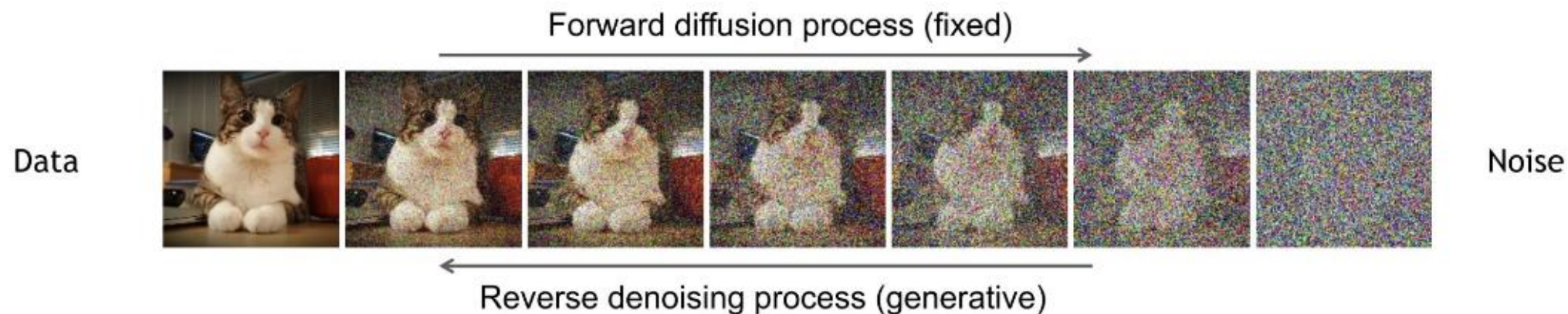
RICE UNIVERSITY

# Denoising Diffusion Probabilistic Models (DDPM)

**Forward diffusion:** Markov chain of diffusion steps to slowly add gaussian noise to data

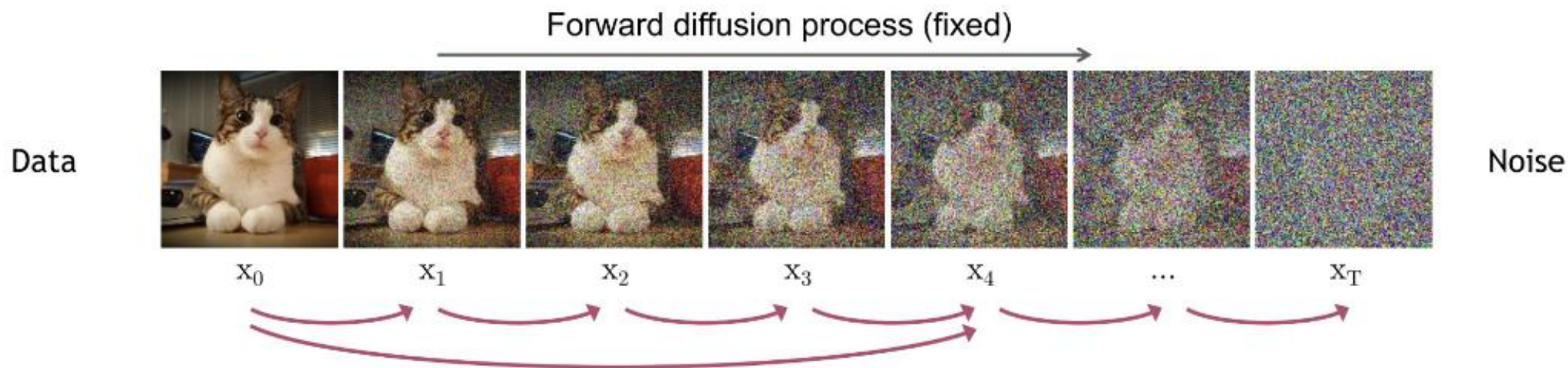**Reverse diffusion:** A model is trained to generate data from noise by iterative denoising



Forward diffusion process (fixed)

Data

Noise

Reverse denoising process (generative)

## Denoising Diffusion Probabilistic Models

**Jonathan Ho**
UC Berkeley
jonathanho@berkeley.edu

**Ajay Jain**
UC Berkeley
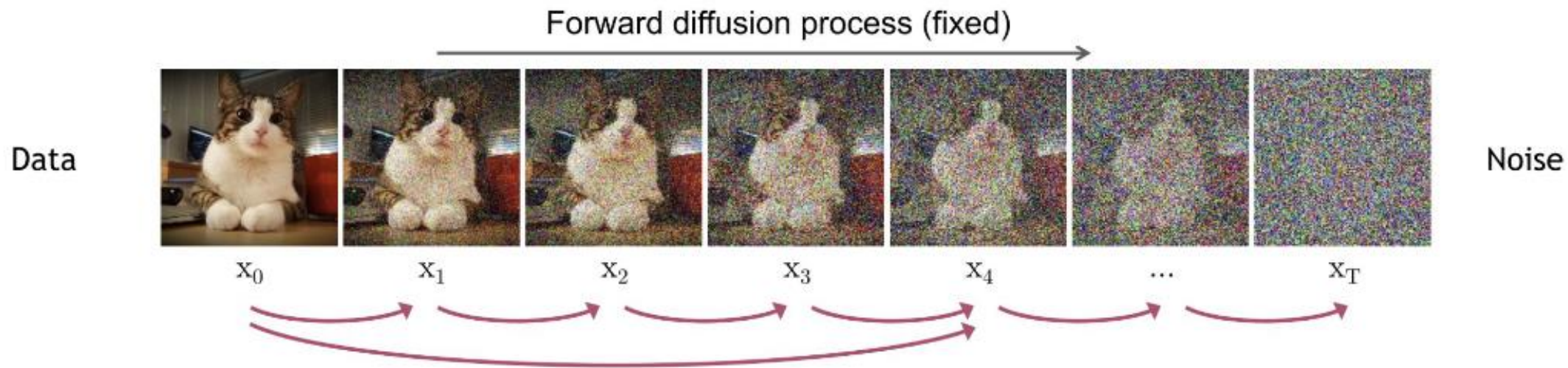ajayj@berkeley.edu

**Pieter Abbeel**
UC Berkeley
pabbeel@cs.berkeley.edu

# DDPM | Forward diffusion



We add a small amount of gaussian noise to a sample $x_0$ in **T** timesteps to produces noised samples, $\{x_1, x_2, \ldots, x_T\}$. The steps are controlled by the noise schedule as follows:

$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1-\beta_t}\mathbf{x}_{t-1}, \beta_t\mathbf{I}) \quad q(\mathbf{x}_{1:T}|\mathbf{x}_0) = \prod_{t=1}^{T} q(\mathbf{x}_t|\mathbf{x}_{t-1})$$

Forward diffusion process (fixed)

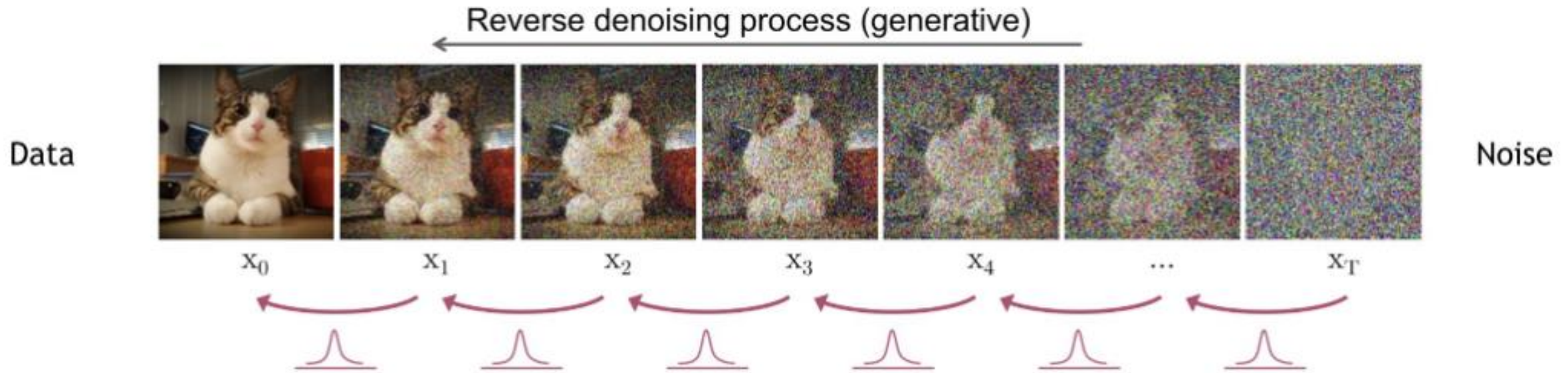Data     Noise

$x_0$  $x_1$  $x_2$  $x_3$  $x_4$  ...  $x_T$

$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1-\beta_t}\mathbf{x}_{t-1}, \beta_t \mathbf{I}) \quad q(\mathbf{x}_{1:T}|\mathbf{x}_0) = \prod_{t=1}^{T} q(\mathbf{x}_t|\mathbf{x}_{t-1})$$

Define $\bar{\alpha}_t = \prod_{s=1}^{t}(1-\beta_s)$  $\Rightarrow$  $q(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1-\bar{\alpha}_t)\mathbf{I}))$  (Diffusion Kernel)

For sampling:  $\mathbf{x}_t = \sqrt{\bar{\alpha}_t}\,\mathbf{x}_0 + \sqrt{(1-\bar{\alpha}_t)}\,\epsilon$   where  $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
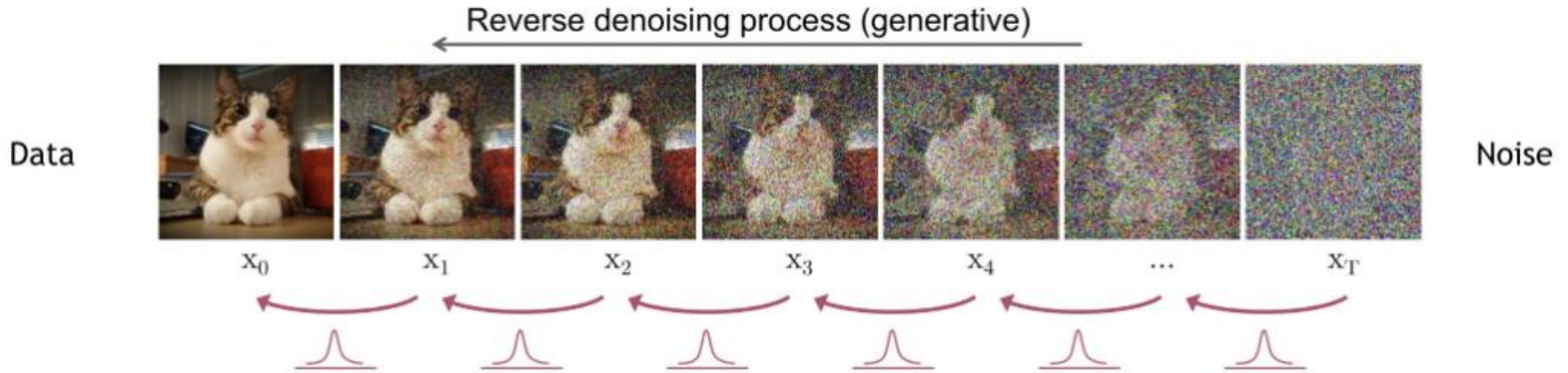
3

# DDPM | Reverse Diffusion



Reverse denoising process (generative)

Data $\quad x_0 \quad x_1 \quad x_2 \quad x_3 \quad x_4 \quad \dots \quad x_T \quad$ Noise

We learn a neural network model $(p_\theta)$ to approximate these conditional probabilities $q(x_{(t-1)} | x_t)$ in order to run the reverse diffusion process as follows:

$$p_\theta(\mathbf{x}_{0:T}) = p(\mathbf{x}_T) \prod_{t=1}^{T} p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) \quad p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t))$$
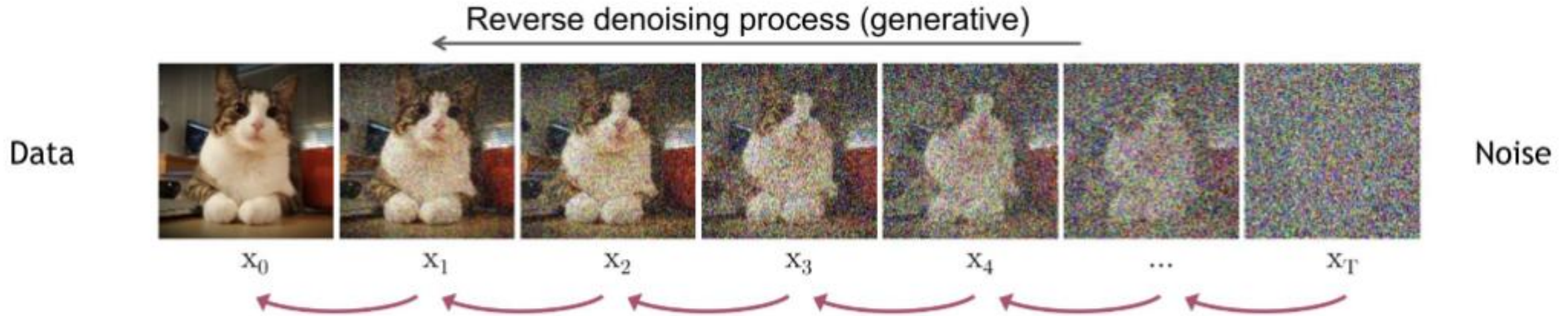
# DDPM | Reverse Diffusion



Reverse denoising process (generative)

Data | $x_0$ | $x_1$ | $x_2$ | $x_3$ | $x_4$ | ... | $x_T$ | Noise

We learn a neural network model ($p_\theta$) to approximate these conditional probabilities $q(x_{(t-1)} \mid x_t)$ in order to run the reverse diffusion process as follows:

$$p_\theta(\mathbf{x}_{0:T}) = p(\mathbf{x}_T) \prod_{t=1}^{T} p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) \quad p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t))$$

# How do we train?



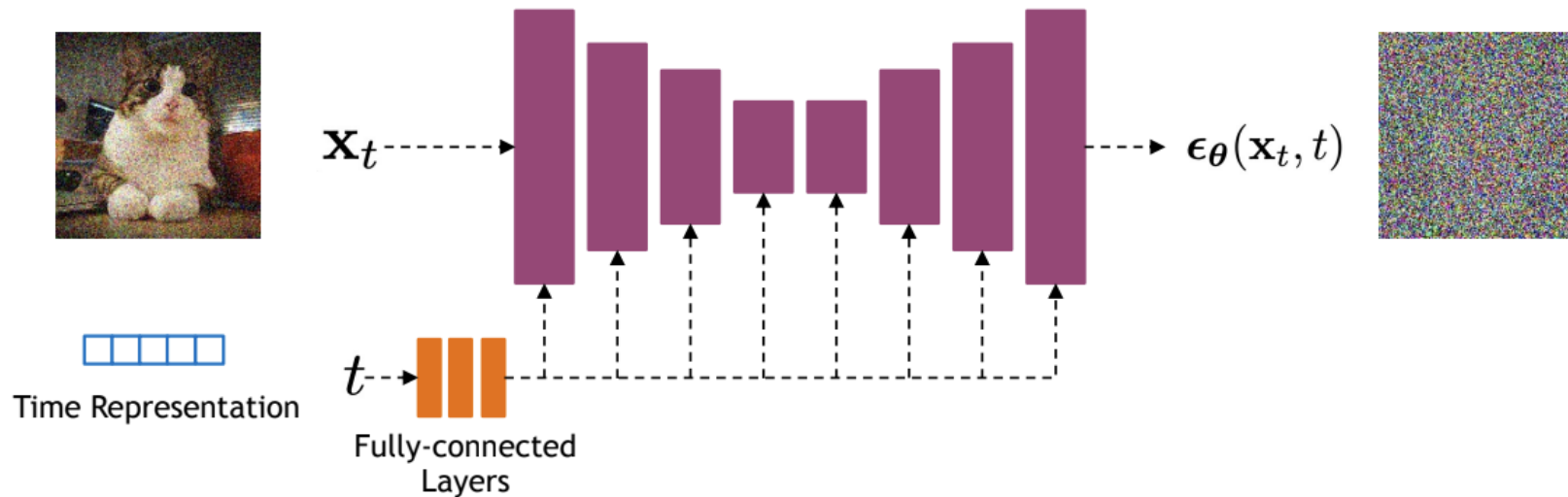Reverse denoising process (generative)

Data

Noise

$x_0$  $x_1$  $x_2$  $x_3$  $x_4$  ...  $x_T$

---

**Algorithm 1** Training

1: **repeat**
2:   $\mathbf{x}_0 \sim q(\mathbf{x}_0)$
3:   $t \sim \text{Uniform}(\{1,\dots,T\})$
4:   $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
5:   Take gradient descent step on
$$\nabla_\theta \left\| \epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1-\bar{\alpha}_t}\epsilon, t) \right\|^2$$
6: **until** converged

---

# Unet to model transition

Diffusion models often use U-Net architectures with ResNet blocks and self-attention layers to represent $\epsilon_\theta(\mathbf{x}_t, t)$
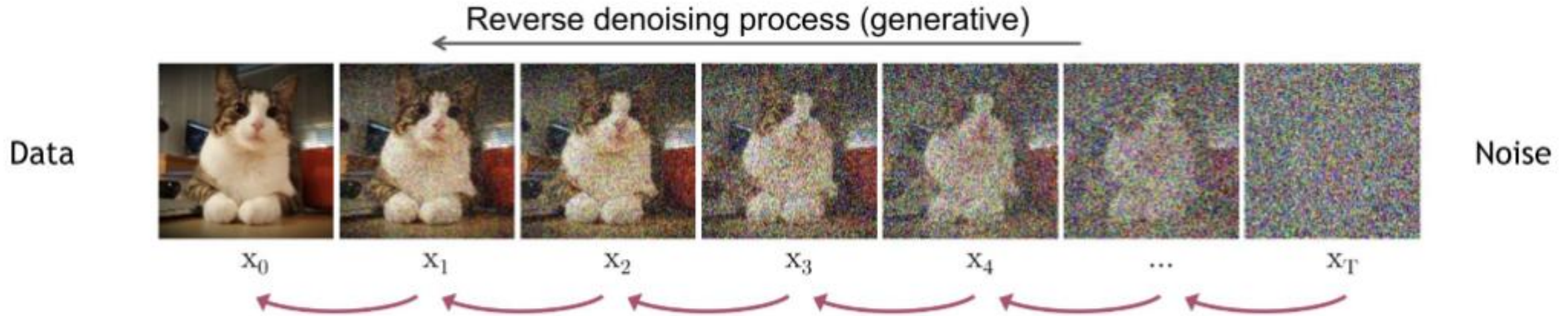


Time representation: sinusoidal positional embeddings or random Fourier features.

Time features are fed to the residual blocks using either simple spatial addition or using adaptive group normalization layers. (see Dharivwal and Nichol NeurIPS 2021)

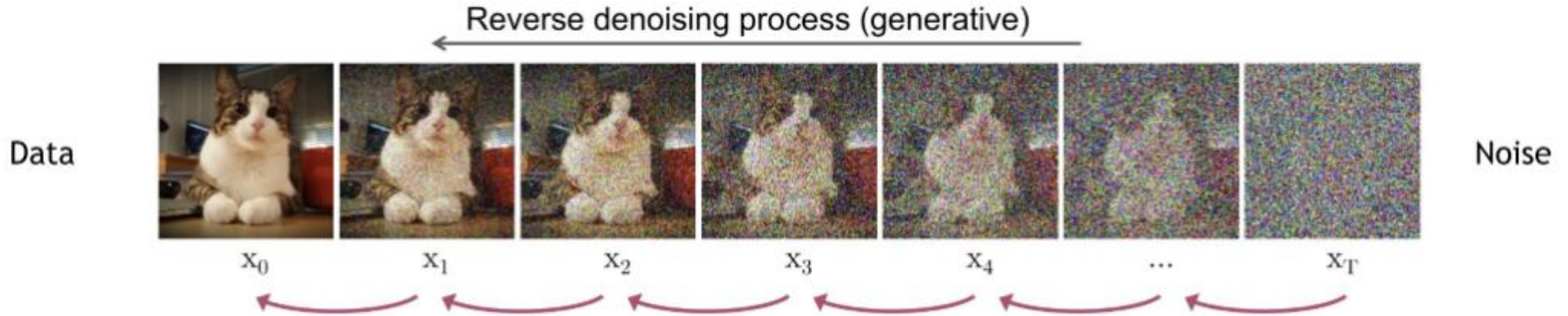https://cvpr2022-tutorial-diffusion-models.github.io/

# How do we train?



Reverse denoising process (generative)

Data ... Noise

$x_0$  $x_1$  $x_2$  $x_3$  $x_4$  ...  $x_T$

**Algorithm 2** Sampling

1: $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
2: **for** $t = T, \ldots, 1$ **do**
3: $\quad \mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ if $t > 1$, else $\mathbf{z} = \mathbf{0}$
4: $\quad \mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$
5: **end for**
6: **return** $\mathbf{x}_0$

# How do we train?



Reverse denoising process (generative)

Data        $x_0$     $x_1$     $x_2$     $x_3$     $x_4$     ...     $x_T$       Noise

**Algorithm 1** Training

1: **repeat**
2:   $\mathbf{x}_0 \sim q(\mathbf{x}_0)$
3:   $t \sim \text{Uniform}(\{1, \ldots, T\})$
4:   $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
5:   Take gradient descent step on
$$\nabla_\theta \left\| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta(\sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\epsilon}, t) \right\|^2$$
6: **until** converged

**Algorithm 2** Sampling

1: $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
2: **for** $t = T, \ldots, 1$ **do**
3:   $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ if $t > 1$, else $\mathbf{z} = \mathbf{0}$
4:   $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$
5: **end for**
6: **return** $\mathbf{x}_0$

# Imagen by Google



A cute corgi lives in a house made out of sushi.



A cute sloth holding a small treasure chest. A bright golden glow is coming from the chest.

# Imagen by Google

## 2.2 Diffusion models and classifier-free guidance

Here we give a brief introduction to diffusion models; a precise description is in Appendix A. Diffusion models [63, 28, 65] are a class of generative models that convert Gaussian noise into samples from a learned data distribution via an iterative denoising process. These models can be conditional, for example on class labels, text, or low-resolution images [e.g. 16, 29, 59, 58, 75, 41, 54]. A diffusion model $\hat{\mathbf{x}}_\theta$ is trained on a denoising objective of the form

$$\mathbb{E}_{\mathbf{x},\mathbf{c},\boldsymbol{\epsilon},t}\left[w_t\|\hat{\mathbf{x}}_\theta(\alpha_t\mathbf{x} + \sigma_t\boldsymbol{\epsilon}, \mathbf{c}) - \mathbf{x}\|_2^2\right] \tag{1}$$

where $(\mathbf{x}, \mathbf{c})$ are data-conditioning pairs, $t \sim \mathcal{U}([0, 1])$, $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, and $\alpha_t, \sigma_t, w_t$ are functions of $t$ that influence sample quality. Intuitively, $\hat{\mathbf{x}}_\theta$ is trained to denoise $\mathbf{z}_t := \alpha_t\mathbf{x} + \sigma_t\boldsymbol{\epsilon}$ into $\mathbf{x}$ using a squared error loss, weighted to emphasize certain values of $t$. Sampling such as the ancestral sampler [28] and DDIM [64] start from pure noise $\mathbf{z}_1 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and iteratively generate points $\mathbf{z}_{t_1}, \dots, \mathbf{z}_{t_T}$, where $1 = t_1 > \cdots > t_T = 0$, that gradually decrease in noise content. These points are functions of the $\mathbf{x}$-predictions $\hat{\mathbf{x}}_0^t := \hat{\mathbf{x}}_\theta(\mathbf{z}_t, \mathbf{c})$.

# Latent Diffusion Models (Stable Diffusion)

## High-Resolution Image Synthesis with Latent Diffusion Models

Robin Rombach[1] *     Andreas Blattmann[1] *     Dominik Lorenz[1]     Patrick Esser[R]     Björn Ommer[1]

[1]Ludwig Maximilian University of Munich & IWR, Heidelberg University, Germany     [R]Runway ML

https://github.com/CompVis/latent-diffusion

### Abstract

By decomposing the image formation process into a sequential application of denoising autoencoders, diffusion models (DMs) achieve state-of-the-art synthesis results on image data and beyond. Additionally, their formulation allows for a guiding mechanism to control the image generation process without retraining. However, since these models typically operate directly in pixel space, optimization of powerful DMs often consumes hundreds of GPU days and inference is expensive due to sequential evaluations. To enable DM training on limited computational

13 Apr 2022



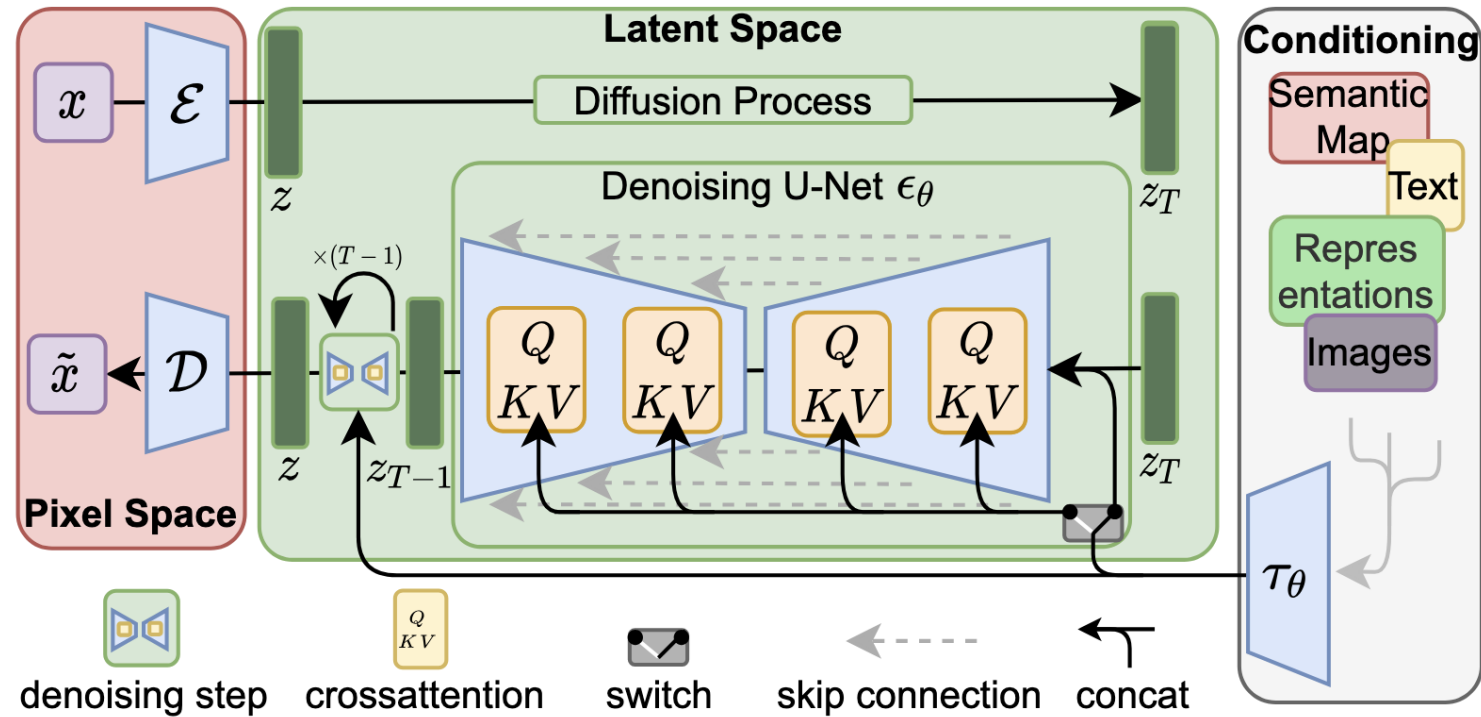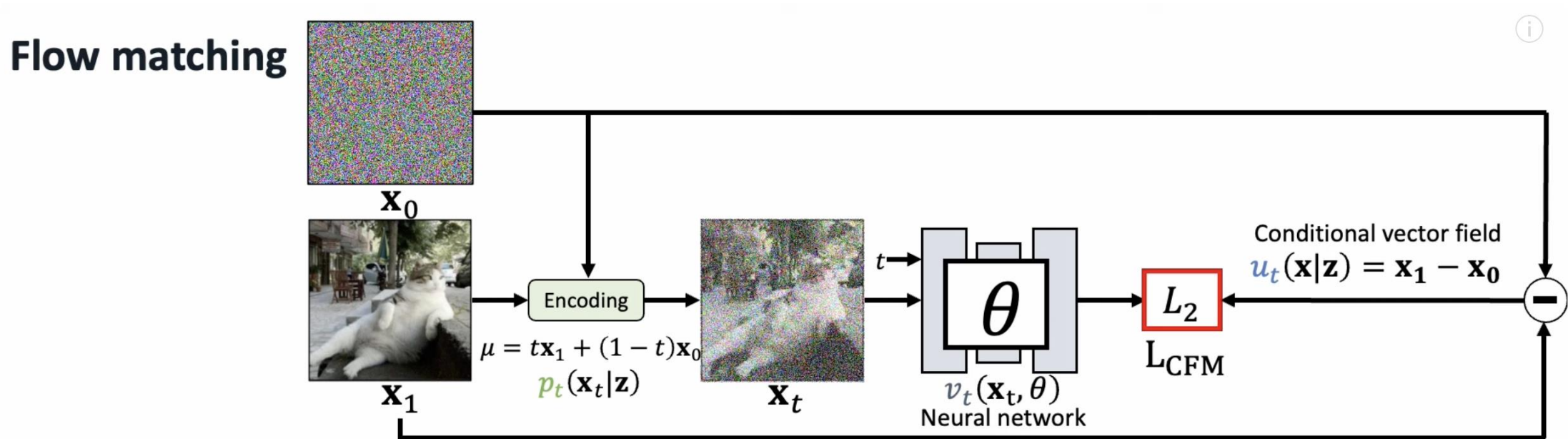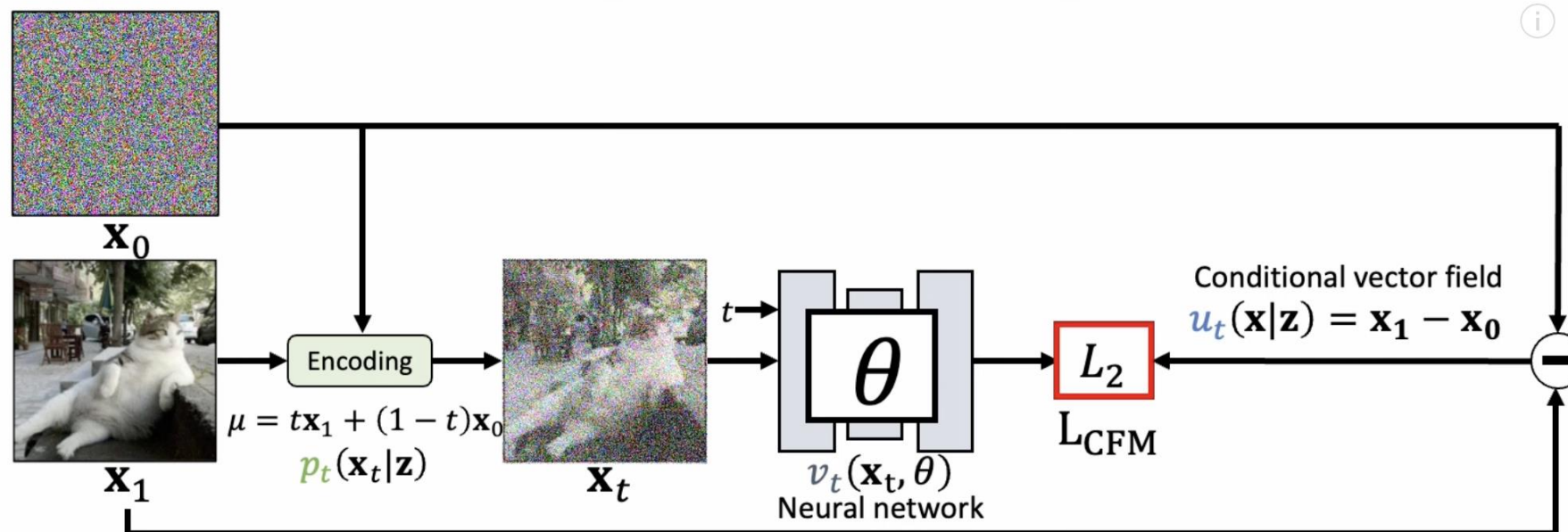| Input | ours ($f = 4$) PSNR: 27.4 R-FID: 0.58 | DALL-E ($f = 8$) PSNR: 22.8 R-FID: 32.01 | VQGAN ($f = 16$) PSNR: 19.9 R-FID: 4.98 |

# Latent Diffusion Models



Figure 3. We condition LDMs either via concatenation or by a more general cross-attention mechanism. See Sec. 3.3

# Another More Recent Improvement: Flow Matching
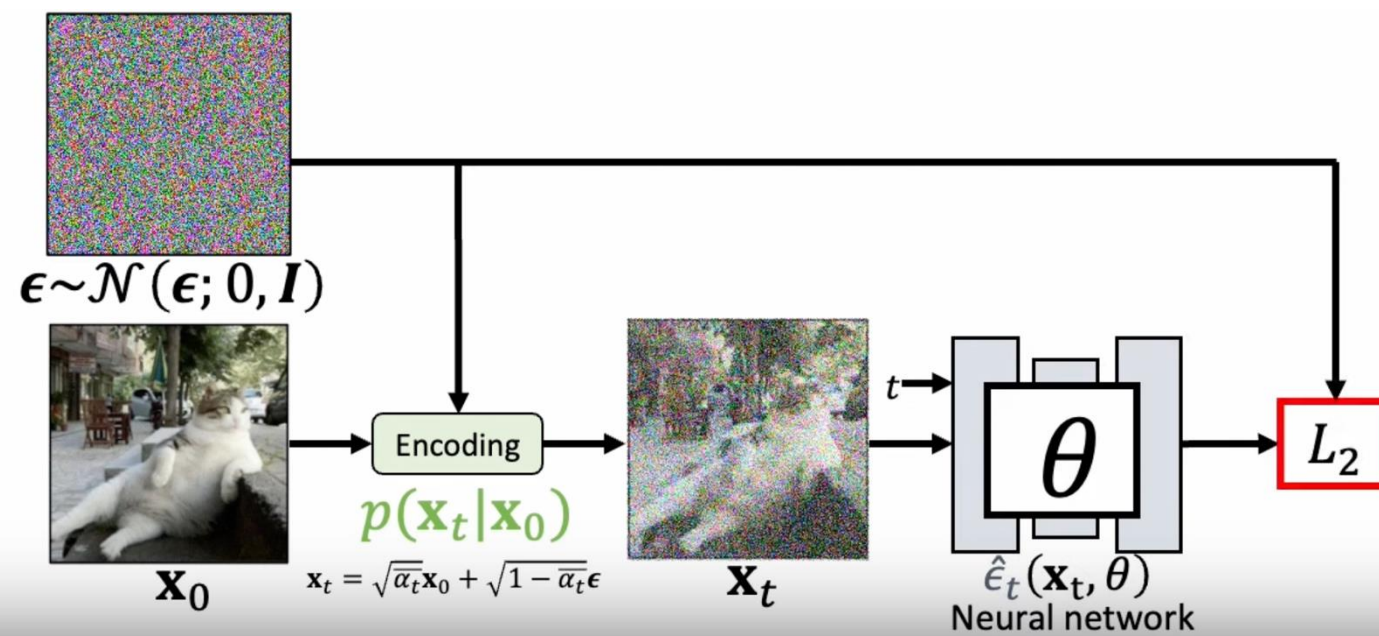


**Flow matching**

$$\mu = t\mathbf{x}_1 + (1-t)\mathbf{x}_0$$
$$p_t(\mathbf{x}_t|\mathbf{z})$$

Encoding

$\mathbf{x}_0$

$\mathbf{x}_1$

$\mathbf{x}_t$

$v_t(\mathbf{x}_t, \theta)$
Neural network

$L_2$

$L_{\text{CFM}}$

Conditional vector field
$$u_t(\mathbf{x}|\mathbf{z}) = \mathbf{x}_1 - \mathbf{x}_0$$

# Flow matching



$\mathbf{x_0}$

Encoding

$\mu = t\mathbf{x_1} + (1-t)\mathbf{x_0}$

$p_t(\mathbf{x_t}|\mathbf{z})$

$\mathbf{x_1}$

$\mathbf{x_t}$

$t \rightarrow$

$\theta$

$v_t(\mathbf{x_t}, \theta)$

Neural network

$L_2$

$\mathrm{L_{CFM}}$

Conditional vector field

$u_t(\mathbf{x}|\mathbf{z}) = \mathbf{x_1} - \mathbf{x_0}$

# Diffusion

$\epsilon \sim \mathcal{N}(\epsilon; 0, \boldsymbol{I})$

Encoding

$p(\mathbf{x_t}|\mathbf{x_0})$

$\mathbf{x_0}$

$\mathbf{x}_t = \sqrt{\overline{\alpha_t}}\mathbf{x_0} + \sqrt{1-\overline{\alpha_t}}\epsilon$

$\mathbf{x_t}$

$t \rightarrow$

$\theta$

$\hat{\epsilon}_t(\mathbf{x_t}, \theta)$

Neural network

$L_2$

15

# AV-Link: Temporally-Aligned Diffusion Features for Cross-Modal Audio-Video Generation

Moayed Haji-Ali[1,2,*]         Willi Menapace[2]         Aliaksandr Siarohin[2]         Ivan Skorokhodov[2]
Alper Canberk[2]         Kwot Sin Lee[2]         Vicente Ordonez[1]         Sergey Tulyakov[2]

[1]Rice University         [2]Snap Inc

Project Webpage: https://snap-research.github.io/AVLink

# AV-Link: Temporally-Aligned Diffusion Features for Cross-Modal Audio-Video Generation

Moaye



okhodov$^2$
kov$^2$

## 3.1. Background

We base our generative models on the Flow Matching framework [47, 50]. Flow Matching expresses generation of data $\mathbf{X}_1 \sim p_d$ as the progressive transformation of $\mathbf{X}_0$ following a path connecting samples from the two distributions. In its simplest formulation [50], the path is instantiated as a linear interpolation between the samples:

$$\mathbf{X}_t = t\mathbf{X}_1 + (1-t)\mathbf{X}_0, \tag{1}$$

and $\mathbf{X}_0 \sim p_n = \mathcal{N}(0, \mathbb{I})$ originate from a noise distribution.

We can move along the path following the velocity $v_t = \frac{d\mathbf{X}_t}{dt} = \mathbf{X}_1 - \mathbf{X}_0$ approximated by learnable $\mathcal{G}$ minimizing:

$$\mathcal{L} = \mathbb{E}_{t \sim p_t, \mathbf{X}_1 \sim p_d, \mathbf{X}_0 \sim p_n} \left\| \mathcal{G}(\mathbf{X}_t, t) - v_t \right\|_2^2, \tag{2}$$

where $p_t$ indicates a training distribution over time $t$, which we instantiate as a logit normal distribution [17].

## 3.1. Background

We base our generative models on the Flow Matching framework [47, 50]. Flow Matching expresses generation of data $\mathbf{X}_1 \sim p_d$ as the progressive transformation of $\mathbf{X}_0$ following a path connecting samples from the two distributions. In its simplest formulation [50], the path is instantiated as a linear interpolation between the samples:

$$\mathbf{X}_t = t\mathbf{X}_1 + (1-t)\mathbf{X}_0, \qquad (1)$$

and $\mathbf{X}_0 \sim p_n = \mathcal{N}(0, \mathbb{I})$ originate from a noise distribution.

We can move along the path following the velocity $v_t = \frac{d\mathbf{X}_t}{dt} = \mathbf{X}_1 - \mathbf{X}_0$ approximated by learnable $\mathcal{G}$ minimizing:

$$\mathcal{L} = \mathbb{E}_{t \sim p_t, \mathbf{X}_1 \sim p_d, \mathbf{X}_0 \sim p_n} \left\| \mathcal{G}(\mathbf{X}_t, t) - v_t \right\|_2^2, \qquad (2)$$

where $p_t$ indicates a training distribution over time $t$, which we instantiate as a logit normal distribution [17].

Assume we have the following:

$$\mathbf{X}_0 \quad \text{and} \quad \frac{d\mathbf{X}_t}{dt} = \mathcal{G}(\mathbf{X}_t, t)$$

How do we find $\mathbf{X}_t$ ?

## 3.1. Background

We base our generative models on the Flow Matching framework [47, 50]. Flow Matching expresses generation of data $\mathbf{X}_1 \sim p_d$ as the progressive transformation of $\mathbf{X}_0$ following a path connecting samples from the two distributions. In its simplest formulation [50], the path is instantiated as a linear interpolation between the samples:

$$\mathbf{X}_t = t\mathbf{X}_1 + (1-t)\mathbf{X}_0, \tag{1}$$

and $\mathbf{X}_0 \sim p_n = \mathcal{N}(0, \mathbb{I})$ originate from a noise distribution.

We can move along the path following the velocity $v_t = \frac{d\mathbf{X}_t}{dt} = \mathbf{X}_1 - \mathbf{X}_0$ approximated by learnable $\mathcal{G}$ minimizing:

$$\mathcal{L} = \mathbb{E}_{t \sim p_t, \mathbf{X}_1 \sim p_d, \mathbf{X}_0 \sim p_n} \left\| \mathcal{G}(\mathbf{X}_t, t) - v_t \right\|_2^2, \tag{2}$$

where $p_t$ indicates a training distribution over time $t$, which we instantiate as a logit normal distribution [17].

At inference time, an ODE solver such as first-order Euler can be employed to produce samples $\mathbf{X}_1$ starting from Gaussian noise $\mathbf{X}_0$ using the model's velocity estimates.

Assume we have the following:

$$\mathbf{X}_0 \quad \text{and} \quad \frac{d\mathbf{X}_t}{dt} = \mathcal{G}(\mathbf{X}_t, t)$$
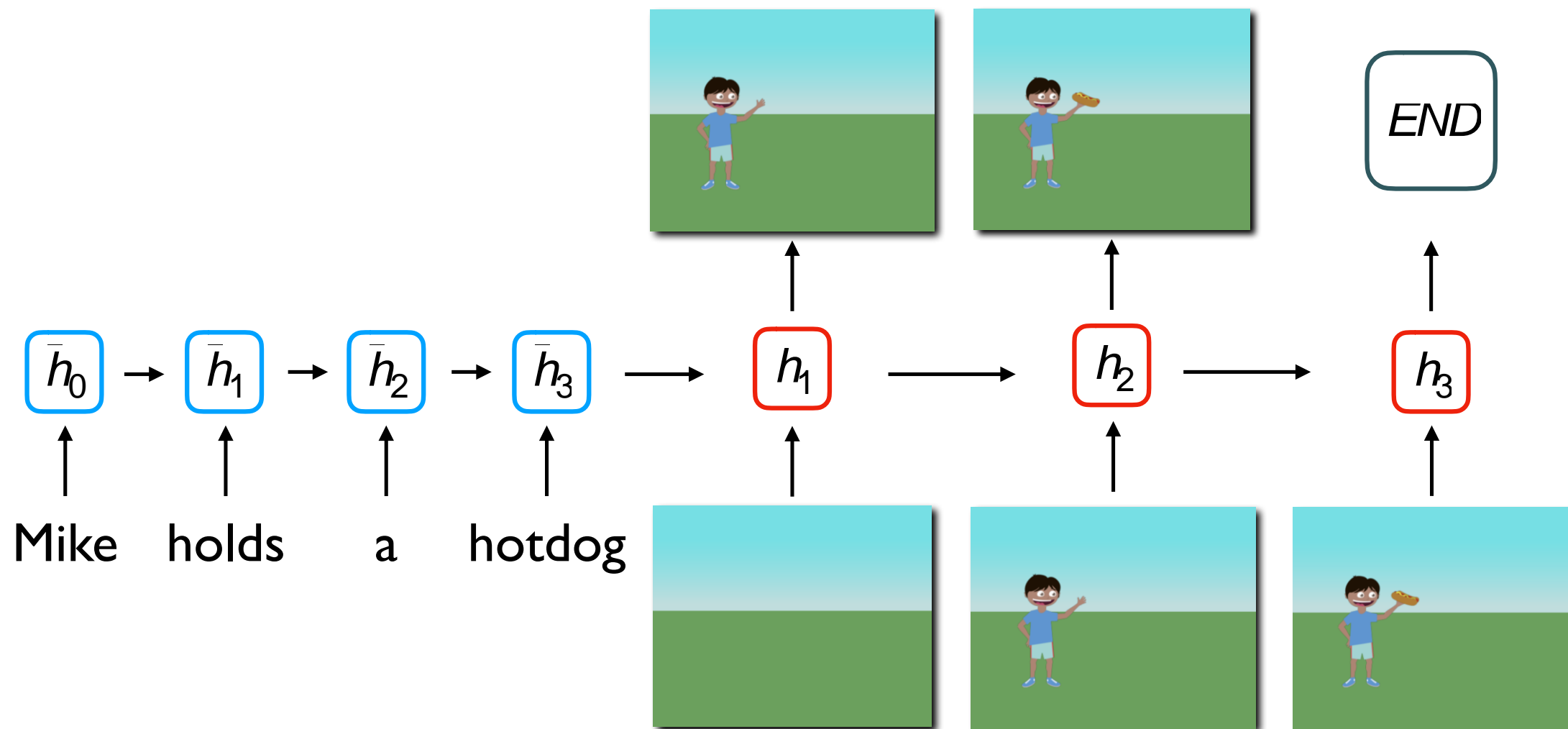
How do we find $\mathbf{X}_t$ ?

# Also check this material to know more

- Entire class on Flow Matching and Diffusion Models: https://diffusion.csail.mit.edu/

# Alternative Methods to Diffusion

- Auto-Regressive models (LLMs to Generate Images!)
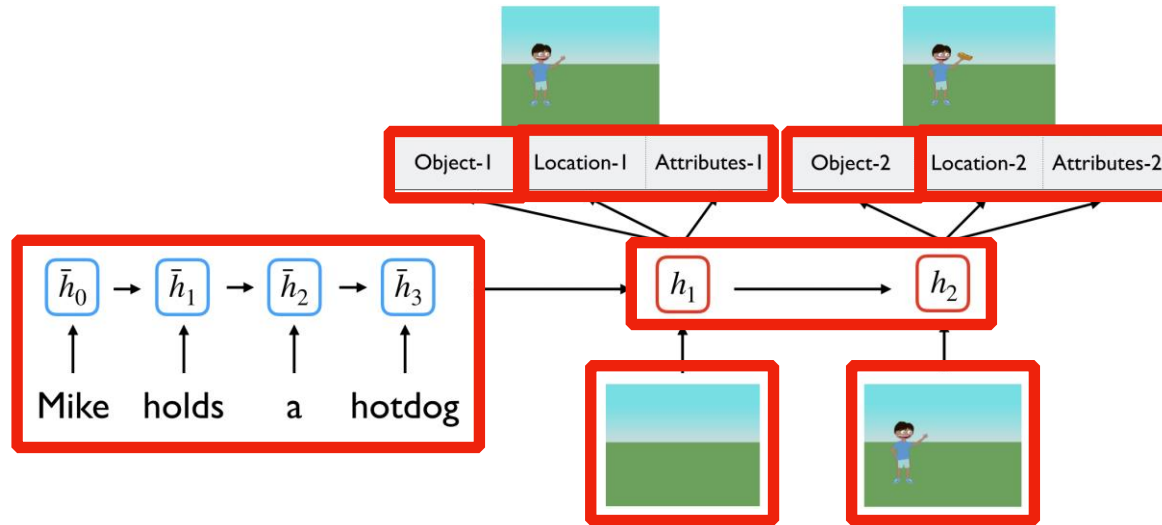
# Text to Scene as Machine Translation!

# The actual model



Object-1 | Location-1 | Attributes-1 | Object-2 | Location-2 | Attributes-2

$\bar{h}_0 \rightarrow \bar{h}_1 \rightarrow \bar{h}_2 \rightarrow \bar{h}_3$

Mike  holds  a  hotdog

$h_1 \longrightarrow h_2$

## Objective

objects     locations

$$L = -w_o \sum_t \log p(o_t) - w_l \sum_t \log p(l_t)$$
$$-\sum_k w_k \sum_t \log p(R_t^k) +$$
$$+w_a^O L_{attn}^O + w_a^A L_{attn}^A,$$

attributes

$h_i^E = \text{BiGRU}(x_i, h_{i-1}^E, h_{i+1}^E) \quad \Omega(B_i \quad h_t^D = \text{ConvGRU}(\Omega( \quad p(o_t) \propto \Theta^o([u_t^o; o_t \quad p(l_t, \{R_t^k\}) = \Theta^a([u_t^a; o_t; c_t^a])$

Encourage attention weights to fully use the input text.

$$L_{attn} = \sum_i [1 - \sum_t \alpha_{t,i}]^2$$

(A) Text Encoder    (B) Image Encoder    (C) Convolutional Recurrent Module    (D) Attention Modules    (E) Object Prediction    (F) Attribute Prediction

Concat   t   t-1   Concat

Input sentence    Canvas    Recurrent hidden state    Language context    Object OneHot    Location map    Attribute maps

uvavision / **Text2Scene**

Watch ▾  6   ★ Star  26   Fork  6

<> Code   ⓘ Issues 0   Pull requests 0   Projects 0   Wiki   Security   Insights   Settings

[CVPR'19] Text2Scene: Generating Compositional Scenes from Textual Descriptions   Edit

Manage topics

⊙ 4 commits   ⌥ 1 branch   ◇ 0 releases   ⚌ 1 contributor

Branch: master ▾   New pull request   |   Create new file   Upload files   Find File   Clone or download ▾

fwtan Update README.md   Latest commit 5681f67 4 days ago

| data | cleaning up the codes, alpha version | 19 days ago |
| examples | cleaning up the codes, alpha version | 19 days ago |
| experiments/scripts | cleaning up the codes, alpha version | 19 days ago |
| lib | cleaning up the codes, alpha version | 19 days ago |
| tools | cleaning up the codes, alpha version | 19 days ago |
| README.md | Update README.md | 4 days ago |

📖 README.md

## Text2Scene: Generating Compositional Scenes from Textual Descriptions

3:32 ⦆   ⚲vislang.ai   ⬆

Besides Mike and Jenny feel free to reference any of these other objects: bear, cat, dog, duck, owl, snake, hat, crown, pirate hat, viking hat, witch hat, glasses, pie, pizza, hot dog, ketchup, mustard, drink, bee, slide, sandbox, swing, tree, pine tree, apple tree, helicopter, balloon, sun, cloud, rocket, airplane, ball, football, basketball, baseball bat, shovel, tennis racket, kite, fire. Also feel free to describe Mike and Jenny with other attributes or action words such as sitting, running, jumping, kicking, standing, afraid, happy, scared, angry, etc.

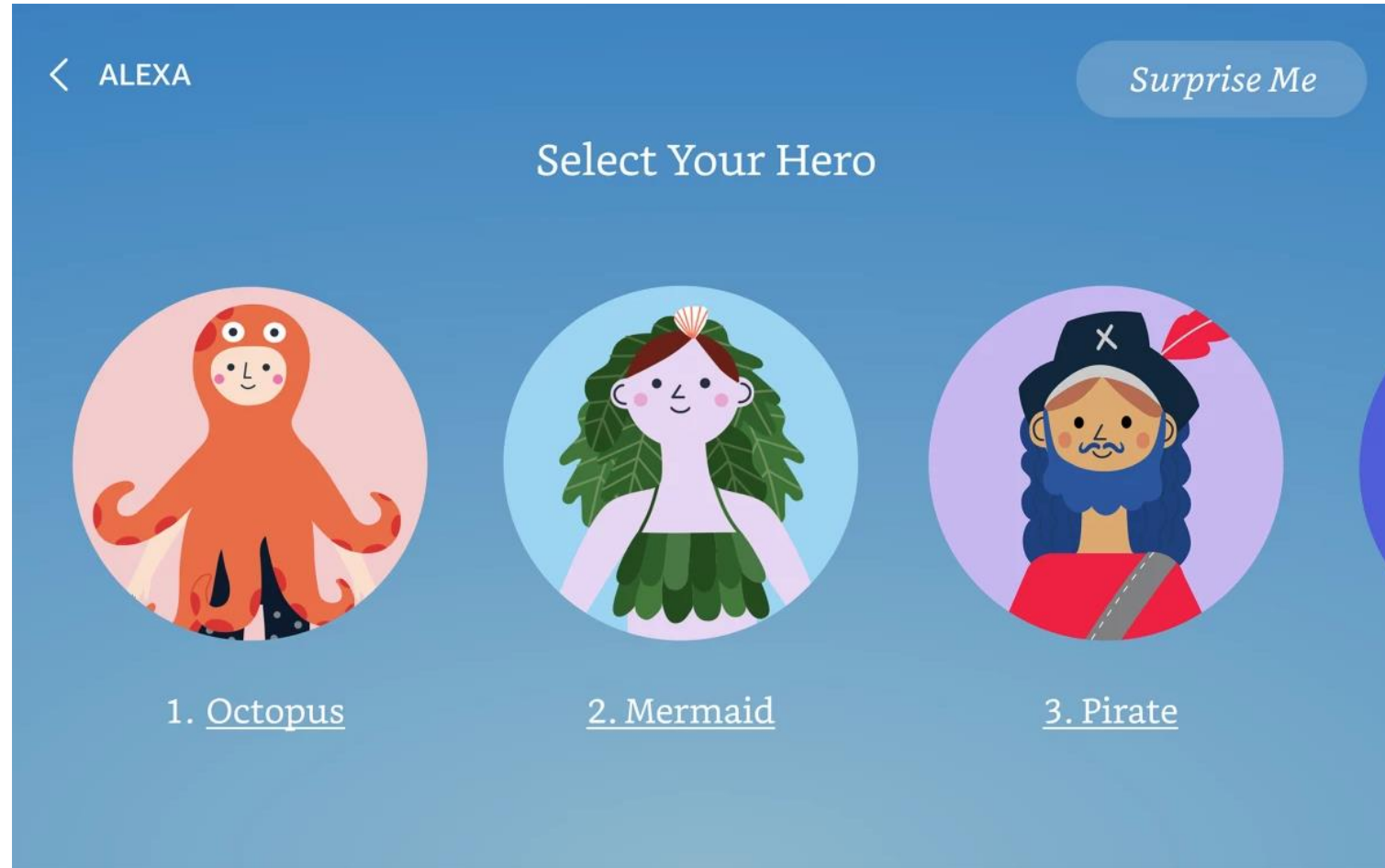#1  Mike is next to a tree
#2  Jenny is happy and kicks the b
#3  There is a fire

Generate Scene

## https://www.vislang.ai/text2scene

# Amazon Alexa AI

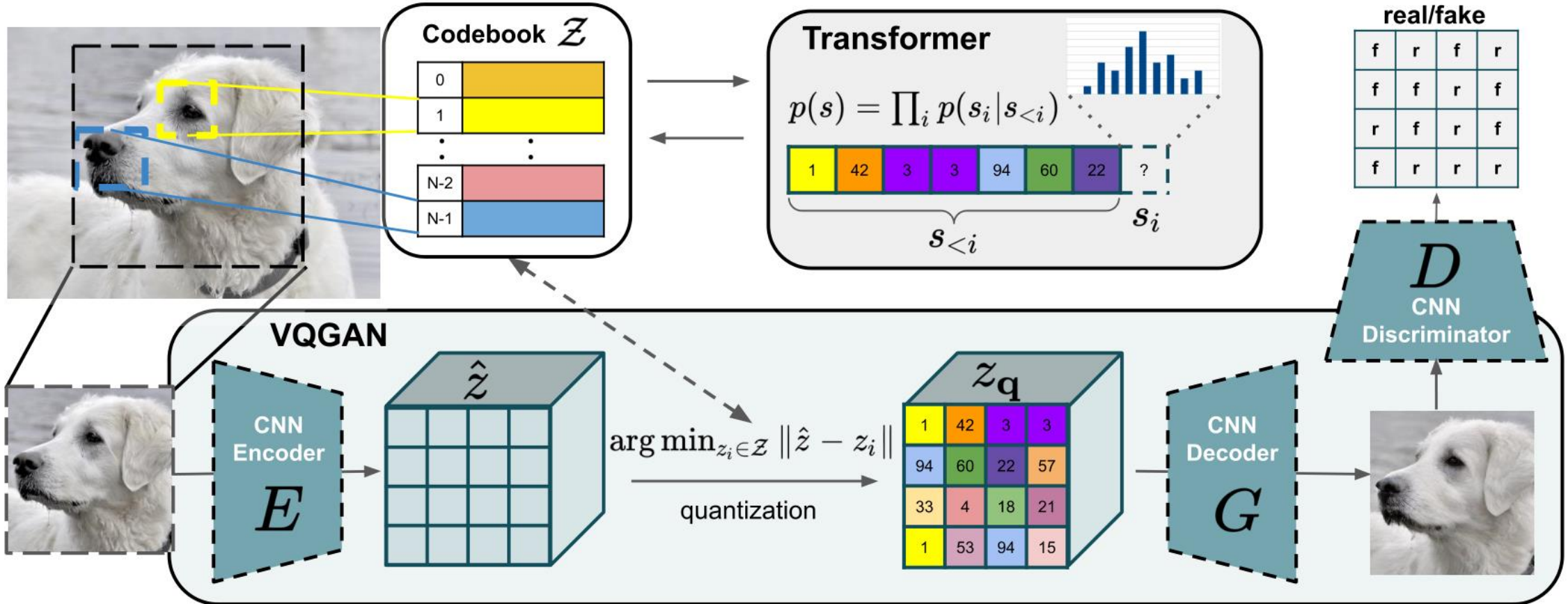# Amazon Alexa AI



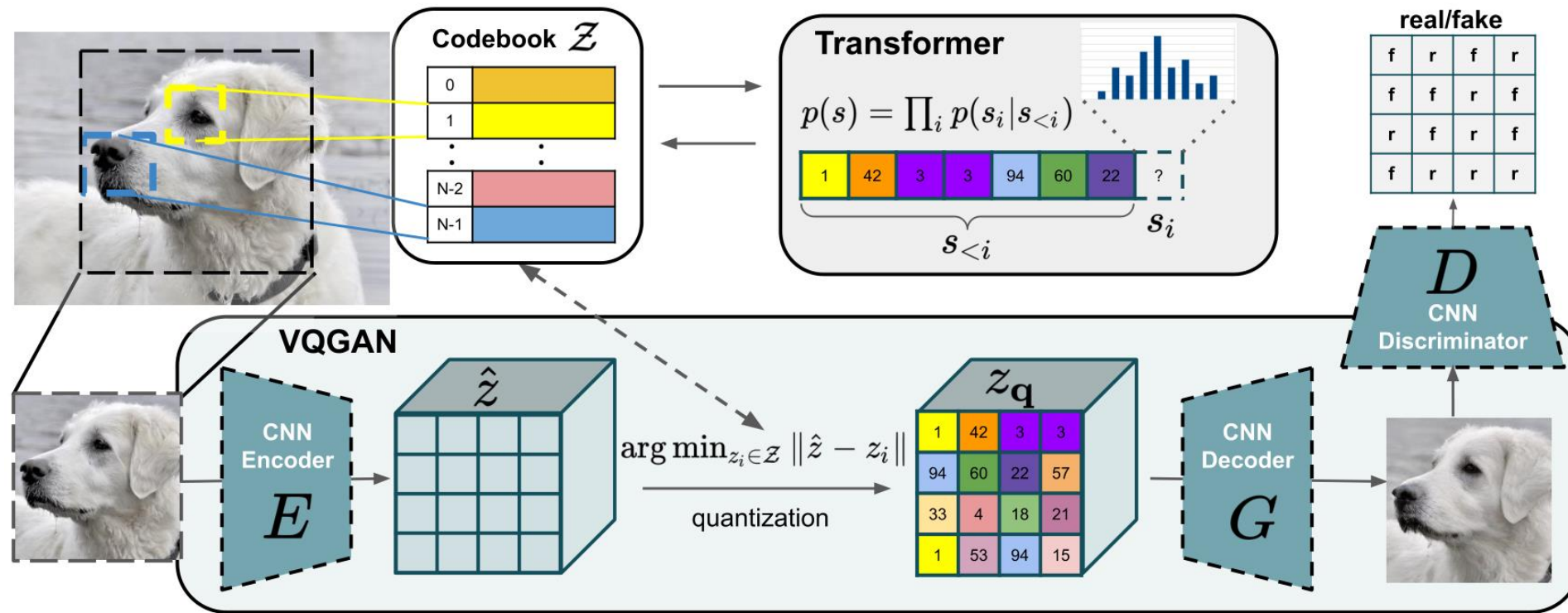Coraline the Mermaid and a scuba diver were having a contest.

https://www.amazon.science/blog/the-science-behind-alexas-new-interactive-story-creation-experience

# Vector Quantized - GAN

# Vector Quantized GAN (VQGAN)



$$Q^* = \arg\min_{E,G,\mathcal{Z}} \max_D \mathbb{E}_{x \sim p(x)} \Big[ \mathcal{L}_{\mathrm{VQ}}(E, G, \mathcal{Z})$$
$$+ \lambda \mathcal{L}_{\mathrm{GAN}}(\{E, G, \mathcal{Z}\}, D) \Big]$$

$$\mathcal{L}_{\mathrm{VQ}}(E, G, \mathcal{Z}) = \|x - \hat{x}\|^2 + \|\mathrm{sg}[E(x)] - z_{\mathbf{q}}\|_2^2$$
$$+ \|\mathrm{sg}[z_{\mathbf{q}}] - E(x)\|_2^2.$$

$$\mathcal{L}_{\mathrm{GAN}}(\{E, G, \mathcal{Z}\}, D) = [\log D(x) + \log(1 - D(\hat{x}))]$$

https://arxiv.org/abs/2012.09841        https://arxiv.org/pdf/2110.04627.pdf

# DALL-E (v1)

Zero-Shot Text-to-Image Generation

Aditya Ramesh[1]   Mikhail Pavlov[1]   Gabriel Goh[1]   Scott Gray[1]
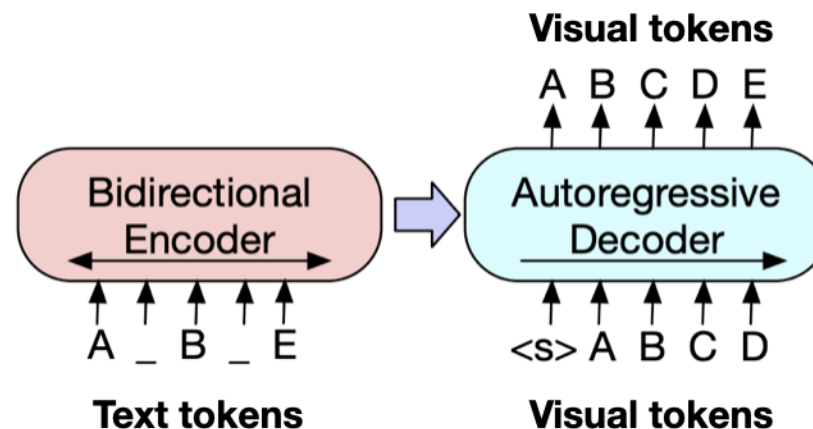Chelsea Voss[1]   Alec Radford[1]   Mark Chen[1]   Ilya Sutskever[1]

**Step 1:**

**Learn Discrete Dictionary of Visual Tokens**

**Step 2:**

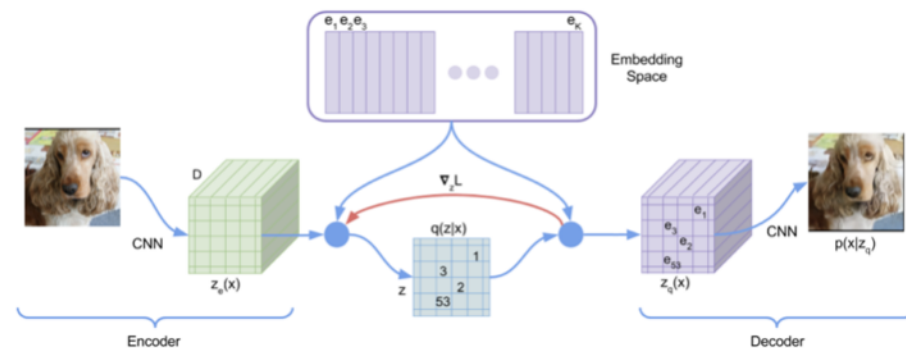**Build a scene as a composition of discrete visual tokens**



VQVAE — Oord, Vinyals, Kavukcuoglu, 2017
VQGAN — Esser, Rombach, Ommer, 2021
dVAE - DALL-E  — Ramesh et al 2021
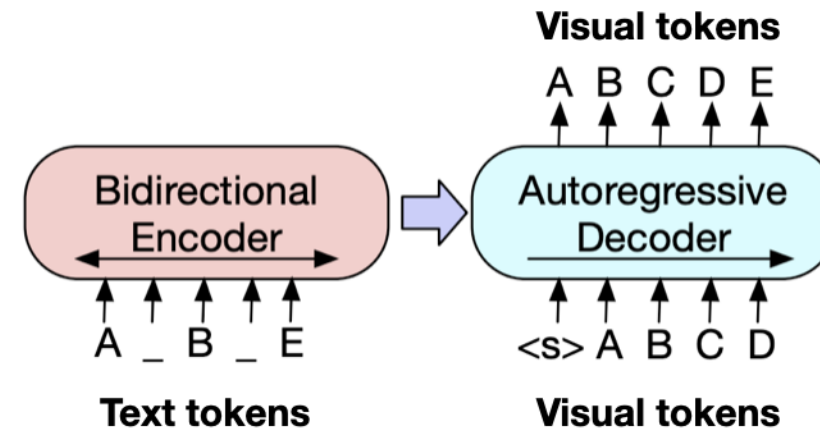


BART, GPT-3, etc

# DALL-E (v1)

**Step 1:**

**Learn Discrete Dictionary of Visual Tokens**

**Step 2:**

**Build a scene as a composition of discrete visual tokens**



VQVAE — Oord, Vinyals, Kavukcuoglu, 2017
VQGAN — Esser, Rombach, Ommer, 2021
dVAE - DALL-E — Ramesh et al 2021

BART, GPT-3, etc

an armchair in the shape of an avocado. . . .

# Questions