# Last Class

- Conditional GANs
- AutoEncoder Models (AEs, VAEs)

## Today:

- Text to image Models
- Sequence-to-sequence based text-to-image models
- Detour: Style Transfer – Input Feature Optimization.
- Reverse Diffusion Models

# Conditional GANs / Text-conditioned
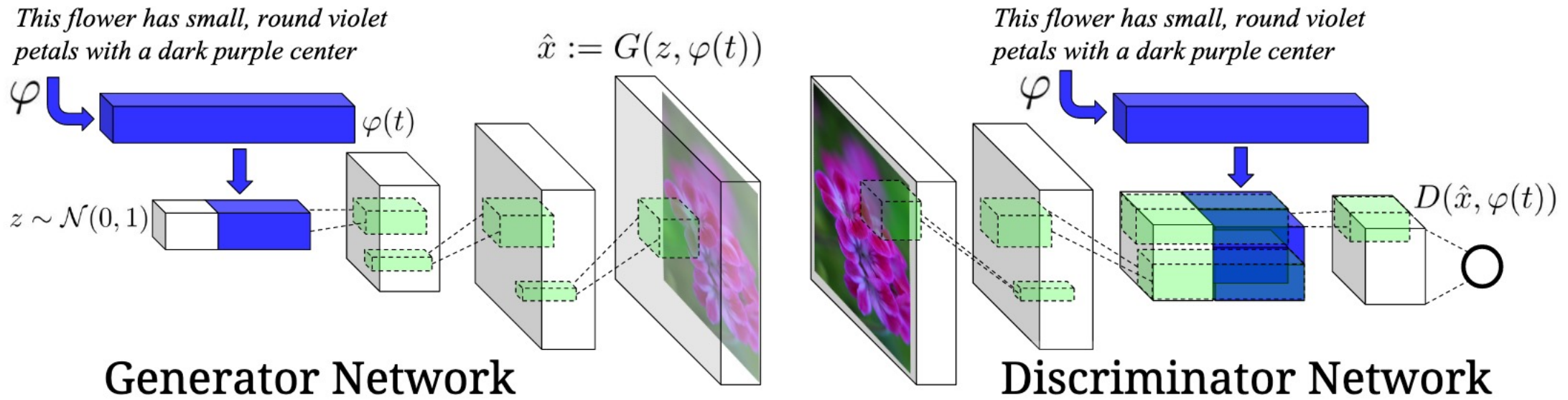
**Generative Adversarial Text to Image Synthesis**

**Scott Reed, Zeynep Akata, Xinchen Yan, Lajanugen Logeswaran**
**Bernt Schiele, Honglak Lee**

REEDSCOT[1], AKATA[2], XCYAN[1], LLAJAN[1]
SCHIELE[2], HONGLAK[1]

[1] University of Michigan, Ann Arbor, MI, USA (UMICH.EDU)
[2] Max Planck Institute for Informatics, Saarbrücken, Germany (MPI-INF.MPG.DE)

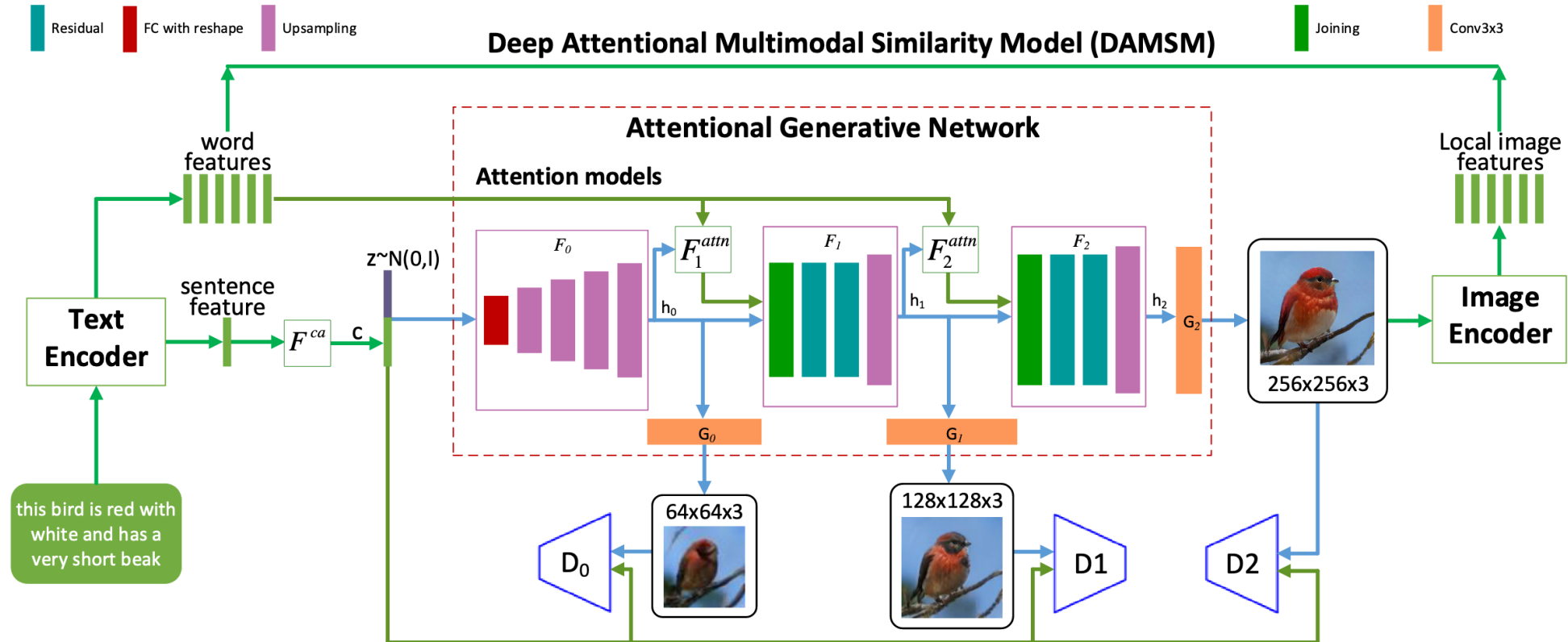# Conditional GANs / Text-conditioned

# Conditional GANs / Text-conditioned



this small bird has a pink breast and crown, and black primaries and secondaries.
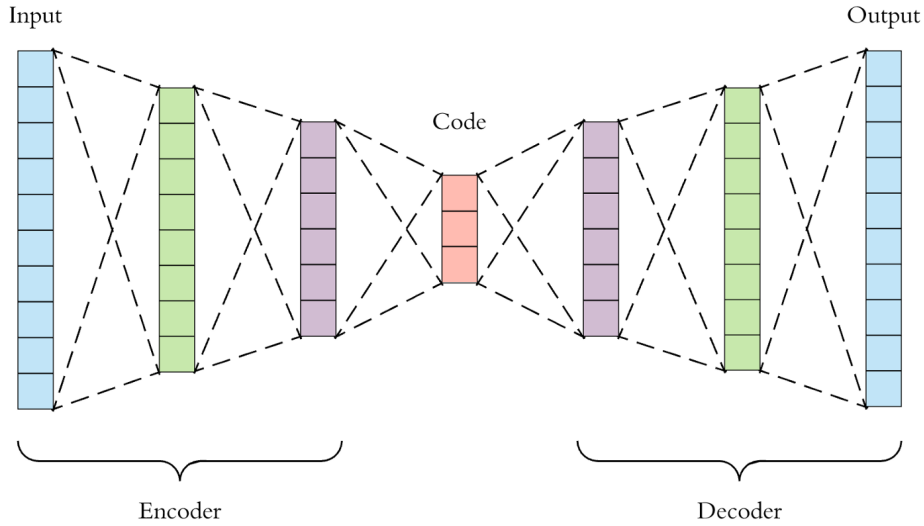
# Conditional GANs / Text-conditioned

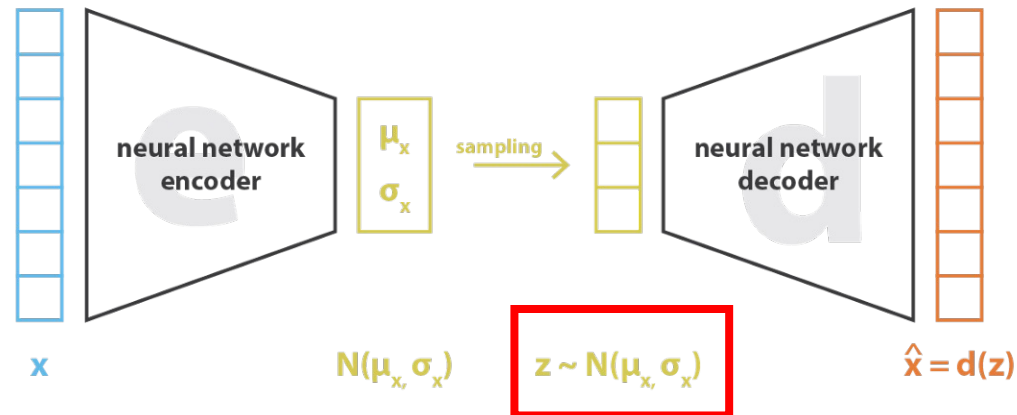# Conditional GANs / Text-conditioned



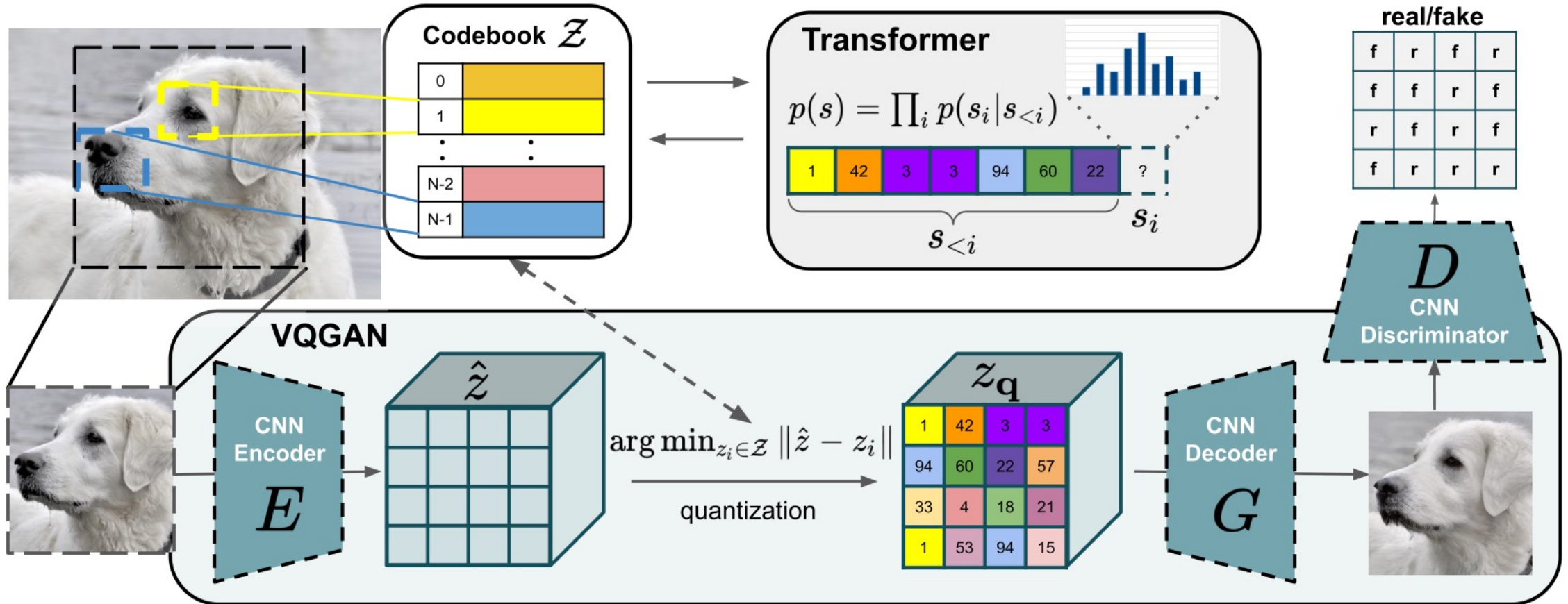this bird is red with white and has a very short beak

# AutoEncoder Models

# Variational AutoEncoder (VAE)



$$\text{loss} = \| x - \hat{x} \|^2 + KL[\, N(\mu_x, \sigma_x), N(0, I) \,]$$

# 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          10

# 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

# DALL-E (v1)

**Step 1:**

**Learn Discrete Dictionary of Visual Tokens**



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

**Step 2:**

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



BART, GPT-3, etc

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

# Text to Scene as Machine Translation!



$\bar{h}_0 \rightarrow \bar{h}_1 \rightarrow \bar{h}_2 \rightarrow \bar{h}_3 \rightarrow h_1 \rightarrow h_2 \rightarrow h_3$

$END$

Mike　holds　a　hotdog

Text2Scene: Generating Compositional Scenes from Textual Descriptions
Fuwen Tan, Song Feng, Vicente Ordonez. Intl. Conference on Computer Vision and Pattern Recognition. **CVPR 2019.**
Long Beach, California. June 2019.(~Oral presentation + Best Paper Finalist -- top 1% of submissions)

# 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

$h_1 \longrightarrow h_2$

objects

Objective

locations

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

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

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

Concat | Concat

t

t-1

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

attributes

Encourage attention weights to fully use the input text.

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

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

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

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

# 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

# More on the Idea of Feature Space Optimization

Gatys et. al. Image Style Transfer Using Convolutional Neural Networks. CVPR 2016

$$E_L = \sum \left(G^L - A^L\right)^2 \qquad \mathcal{L}_{total} = \alpha\mathcal{L}_{content} + \beta\mathcal{L}_{style}$$

$$G_{ij}^L = \sum_k F_{ik}^L F_{jk}^L.$$

$$A^L \quad G^L \qquad F^L$$

$$\frac{\partial E_L}{\partial F^L} \qquad \frac{\partial E_L}{\partial F^{L-1}} \qquad \mathcal{L}_{content} = \sum \left(F^l - P^l\right)^2$$

$$F^{L-1}$$

$$\frac{\partial \mathcal{L}_{total}}{\partial \vec{x}} \qquad \text{Gradient descent}$$

$$\mathcal{L}_{style} = \sum_l w_l E_l$$

$$\vec{a} = \qquad \vec{x} = \qquad \vec{p} =$$

$$\vec{x} := \vec{x} - \lambda \frac{\partial \mathcal{L}_{total}}{\partial \vec{x}}$$

20

# Idea 1: Image Reconstruction from Features

# Idea 1: Image Reconstruction from Features

$$\mathcal{L}_{content} = \sum \left( F^l - P^l \right)^2$$



$$\mathcal{L}_{total} = \alpha \mathcal{L}_{content}$$

$$\mathcal{L}_{content} = \sum \left( F^l - P^l \right)^2$$

$$\frac{\partial E_L}{\partial F^{L-1}}$$

$$\frac{\partial \mathcal{L}_{total}}{\partial \vec{x}} \quad \text{Gradient descent}$$

$$\vec{x} := \vec{x} - \lambda \frac{\partial \mathcal{L}_{total}}{\partial \vec{x}}$$

$\vec{x} =$

$\vec{p} =$

# Idea 1: Image Reconstruction from Features



$$\mathcal{L}_{total} = \alpha \mathcal{L}_{content}$$

$$F^L$$

$$\frac{\partial E_L}{\partial F^{L-1}}$$

$$\mathcal{L}_{content} = \sum \left(F^l - P^l\right)^2$$

$$F^{L-1}$$

conv5_3 $^4_2$
$_1$

pool4

conv4_3 $^4_2$
$_1$

pool3

conv3_3 $^4_2$
$_1$

pool2

conv2_1 $^2$

pool1

conv1_1 $^2$

input

$$\frac{\partial \mathcal{L}_{total}}{\partial \vec{x}}$$

Gradient descent

$$\vec{x} =$$

$$\vec{p} =$$

$$\vec{x} := \vec{x} - \lambda \frac{\partial \mathcal{L}_{total}}{\partial \vec{x}}$$

Gatys et. al. Image Style Transfer Using Convolutional Neural Networks. CVPR 2016

$$\mathcal{L}_{content} = \sum \left(F^l - P^l\right)^2$$

Idea 2: Backpropagation of Style

Gatys et. al. Image Style Transfer Using Convolutional Neural Networks. CVPR 2016



Idea 2: Backpropagation of Style

$$G_{ij}^{l} = \sum_{k} F_{ik}^{l} F_{jk}^{l}.$$

$$E_L = \sum \left(G^L - A^L\right)^2$$

Gatys et. al. Image Style Transfer Using Convolutional Neural Networks. CVPR 2016



Idea 2: Backpropagation of Style

$$E_L = \sum (G^L - A^L)^2 \qquad \mathcal{L}_{total} = \beta \mathcal{L}_{style}$$

$$G_{ij}^L = \sum_k F_{ik}^L F_{jk}^L.$$

$$\mathcal{L}_{style} = \sum_l w_l E_l$$

$$\vec{x} := \vec{x} - \lambda \frac{\partial \mathcal{L}_{total}}{\partial \vec{x}}$$

**Style Reconstructions**

Style Representations

Input image

Content Representations

Convolutional Neural Network

$$G_{ij}^l = \sum_k F_{ik}^l F_{jk}^l.$$

$$E_L = \sum \left(G^L - A^L\right)^2$$

Gatys et. al. Image Style Transfer Using Convolutional Neural Networks. CVPR 2016

$$E_L = \sum (G^L - A^L)^2 \qquad \mathcal{L}_{total} = \alpha\mathcal{L}_{content} + \beta\mathcal{L}_{style}$$

$$G^L_{ij} = \sum_k F^L_{ik}F^L_{jk}.$$

$$\frac{\partial E_L}{\partial F^L} \qquad \frac{\partial E_L}{\partial F^{L-1}} \qquad \mathcal{L}_{content} = \sum \left(F^l - P^l\right)^2$$

$$\mathcal{L}_{style} = \sum_l w_l E_l$$

$$\frac{\partial \mathcal{L}_{total}}{\partial \vec{x}} \qquad \text{Gradient descent}$$

$$\vec{x} =$$

$$\vec{a} = \qquad \vec{p} =$$

$$\vec{x} := \vec{x} - \lambda\frac{\partial \mathcal{L}_{total}}{\partial \vec{x}}$$

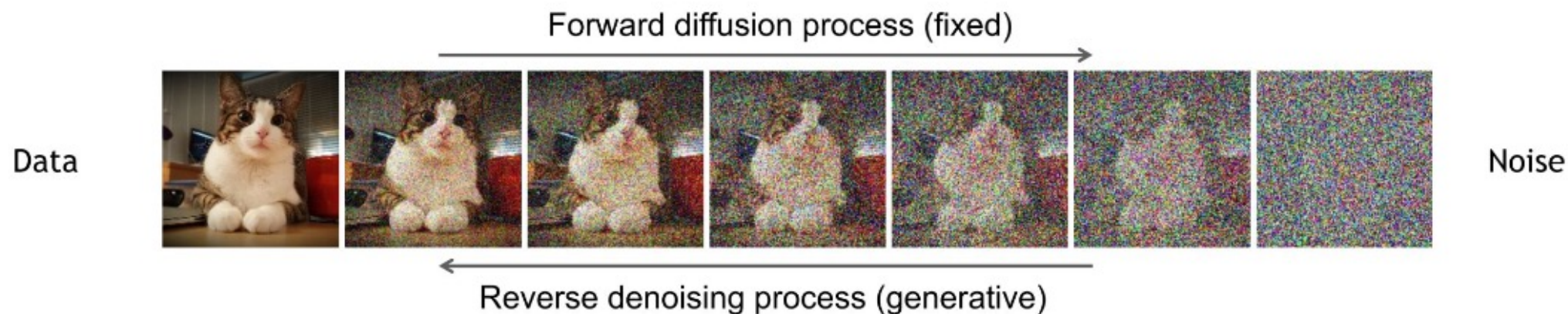# Denoising Diffusion Probabilistic Models (DDPM)

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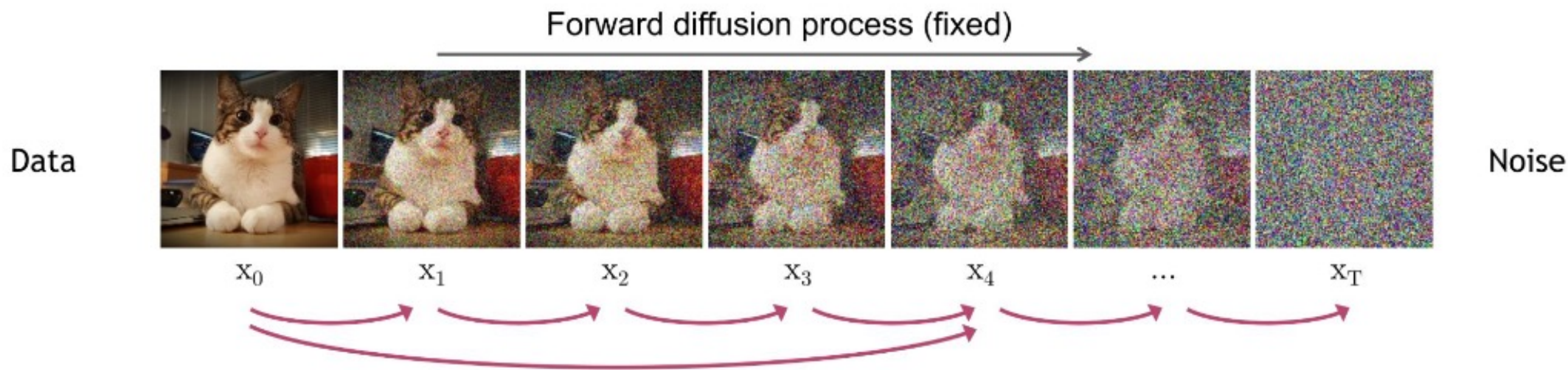
**Jonathan Ho**
UC Berkeley
jonathanho@berkeley.edu

**Ajay Jain**
UC Berkeley
ajayj@berkeley.edu
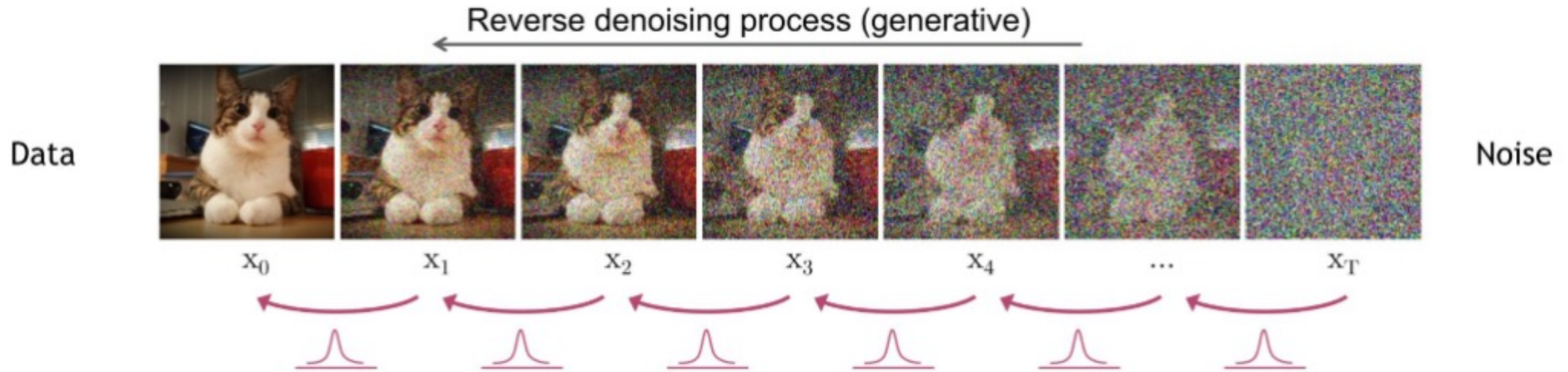
**Pieter Abbeel**
UC Berkeley
pabbeel@cs.berkeley.edu

31

# DDPM | Forward diffusion



Forward diffusion process (fixed)

Data                                                                                                 Noise

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

We add a small amount of gaussian noise to a sample $\mathbf{x}_0$ in **T** timesteps to produces noised samples, $\{\mathbf{x}_1, \mathbf{x}_2, ... , \mathbf{x}_T\}$. The steps are controlled by the noise schedule as follows:

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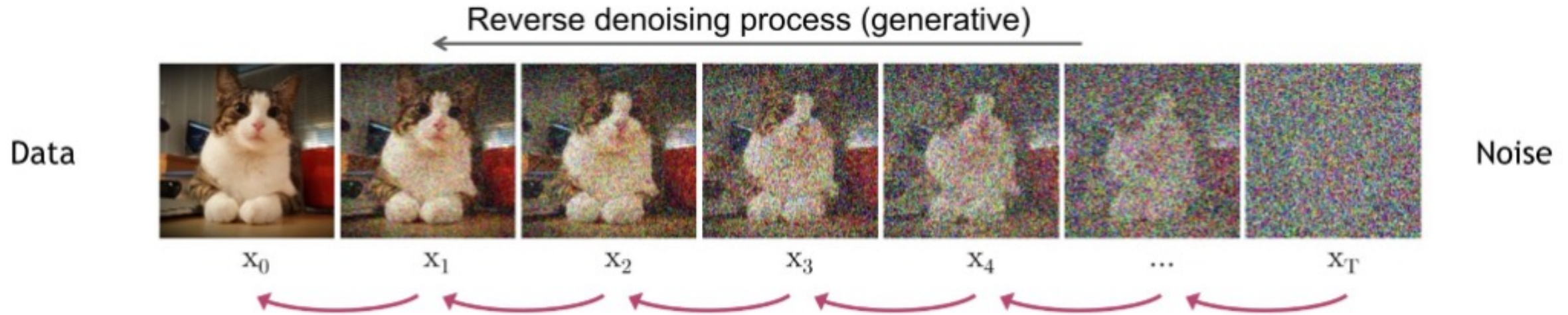
# DDPM | Reverse Diffusion

Reverse denoising process (generative)

Data

Noise

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

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:

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

Slides compiled by my student Aman Shrivastava

# How do we train?

Reverse denoising process (generative)

Data $\quad \mathbf{x}_0 \quad\quad \mathbf{x}_1 \quad\quad \mathbf{x}_2 \quad\quad \mathbf{x}_3 \quad\quad \mathbf{x}_4 \quad\quad \ldots \quad\quad \mathbf{x}_T \quad$ Noise

**Algorithm 1** Training

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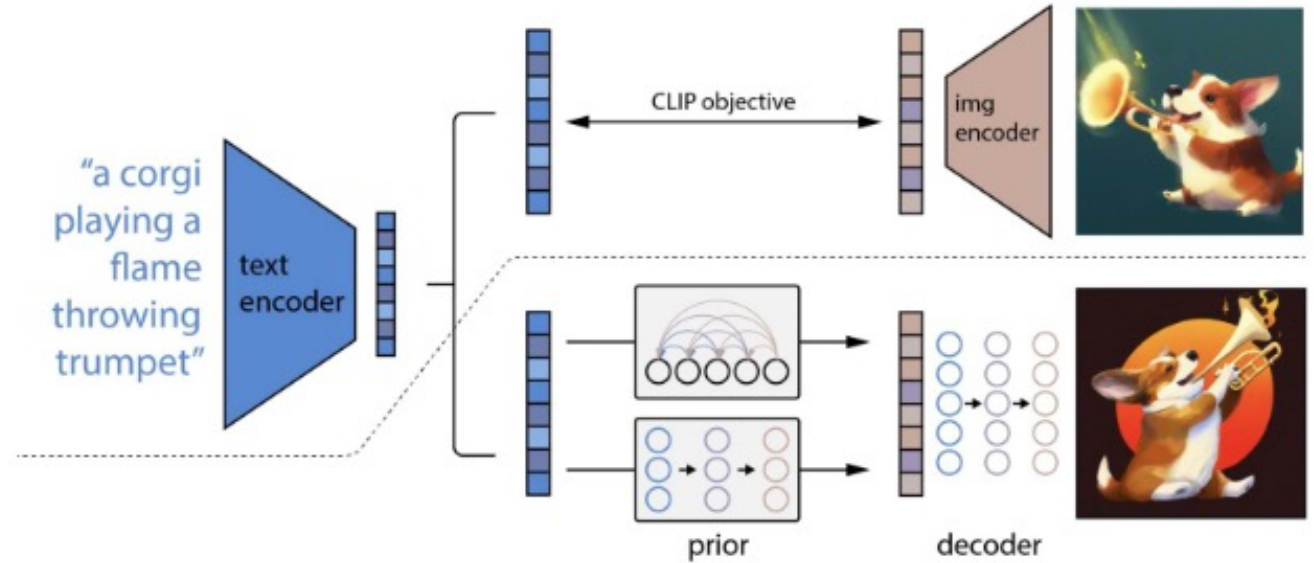6: **until** converged

**Algorithm 2** Sampling

1: $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
2: **for** $t = T, \ldots, 1$ **do**
3: $\quad \mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 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** $\mathbf{x}_0$

Slides compiled by my student Aman Shrivastava

# DALL.E 2 | Open AI

Conditioning on CLIP-embeddings

- Helps capture multimodal representations

- The bi-partite latent enables several text-controlled image manipulation tasks



Ramesh, Aditya, et al. "Hierarchical text-conditional image generation with clip latents." *arXiv preprint arXiv:2204.06125* (2022).

# DALL.E 2 | OpenAI

- 1kx1k text-conditioned image generation
- Uses a **prior** to produce CLIP embeddings conditioned on the text-caption
- Uses a **decoder** to produce images conditioned on the CLIP embeddings



a shiba inu wearing a beret and black turtleneck

a close up of a handpalm with leaves growing from it

panda mad scientist mixing sparkling chemicals, artstation

a corgi's head depicted as an explosion of a nebula

# Questions