

Diffusion Models

2022
Multicampus

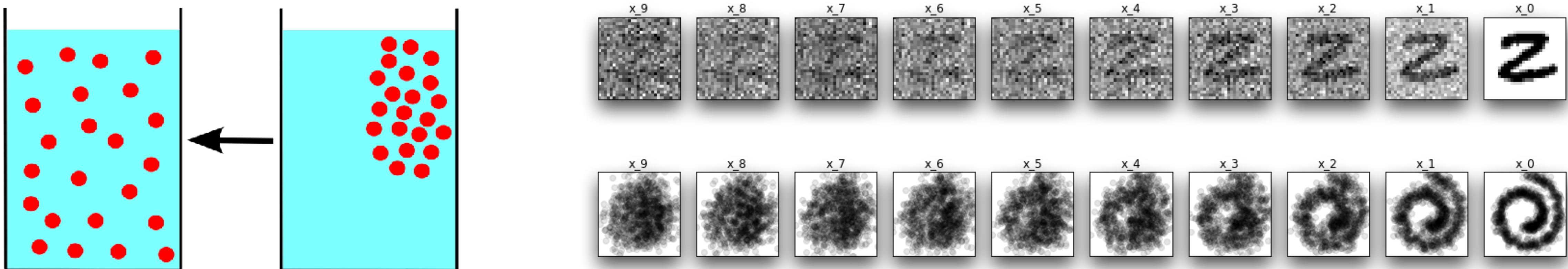
Soochul Park
Vivestudios, Research Scientist

Diffusion Model

DDPM

Diffusion Model

- Sohl-Dickstein, Jascha의 논문 Deep unsupervised learning using nonequilibrium thermodynamics에서 제안

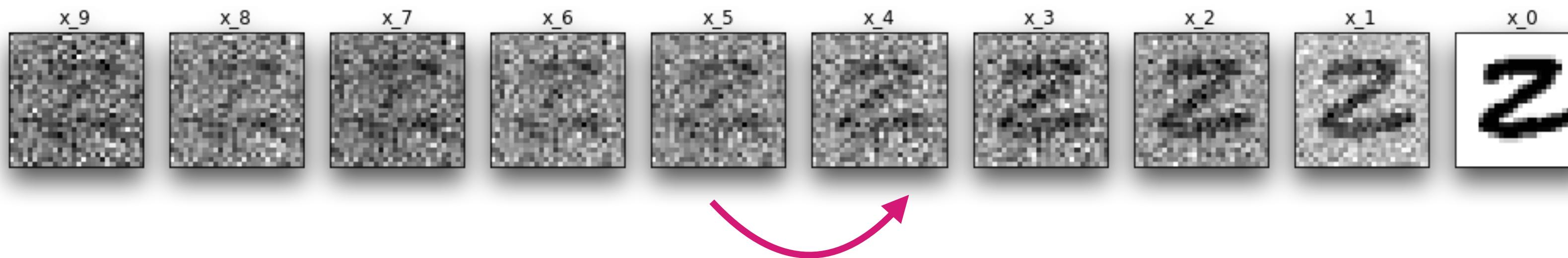


Diffusion

[https://en.wikipedia.org/wiki/Diffusion
\(Flipped\)](https://en.wikipedia.org/wiki/Diffusion_(Flipped))

Diffusion Model

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



Posterior $q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_{t-1}; \tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0), \tilde{\beta}_t \mathbf{I})$

where $\tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0) := \frac{\sqrt{\bar{\alpha}_{t-1}} \beta_t}{1 - \bar{\alpha}_t} \mathbf{x}_0 + \frac{\sqrt{\alpha_t} (1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \mathbf{x}_t$ and $\tilde{\beta}_t := \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t$

**Backward Process
(Neural Networks)**

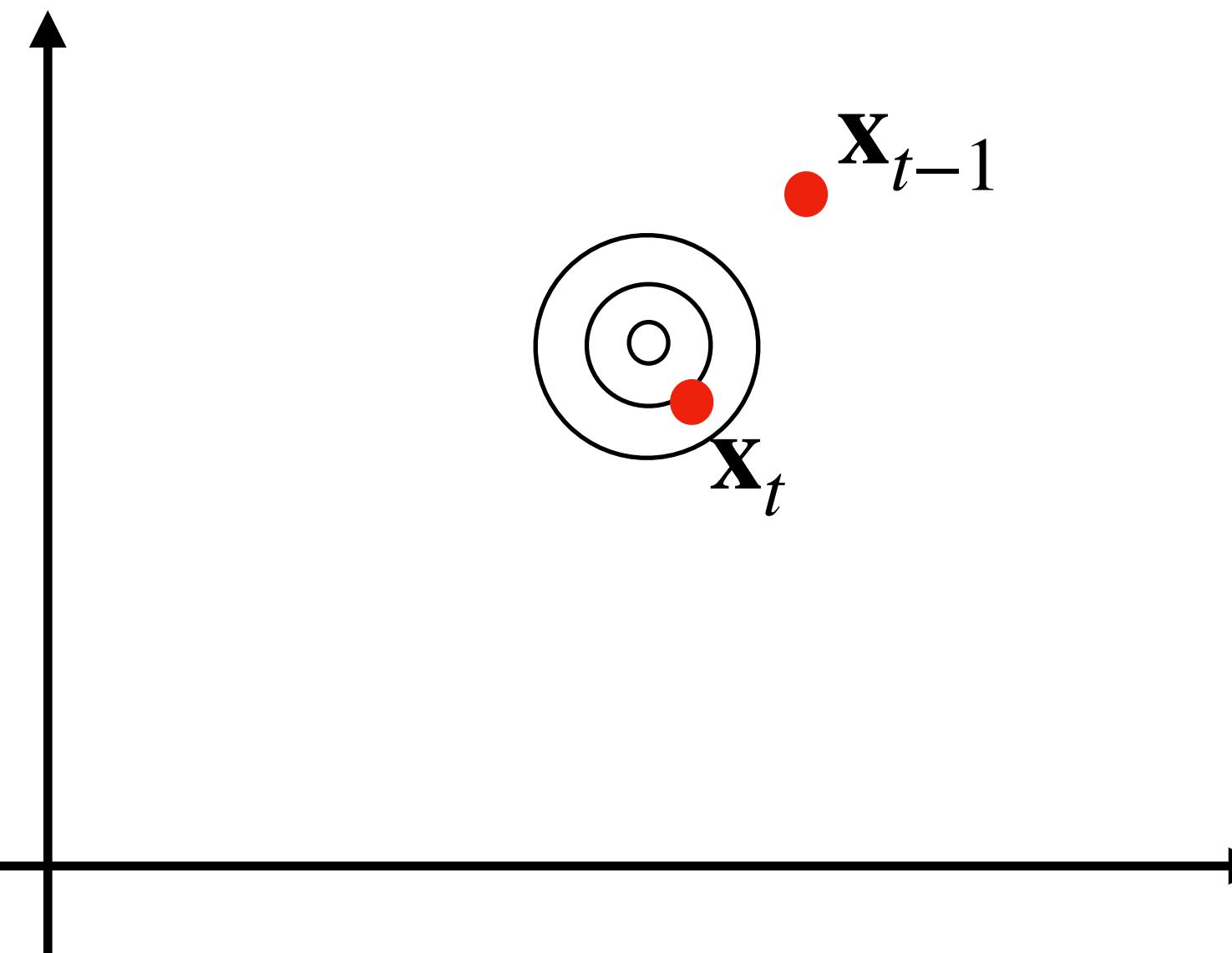
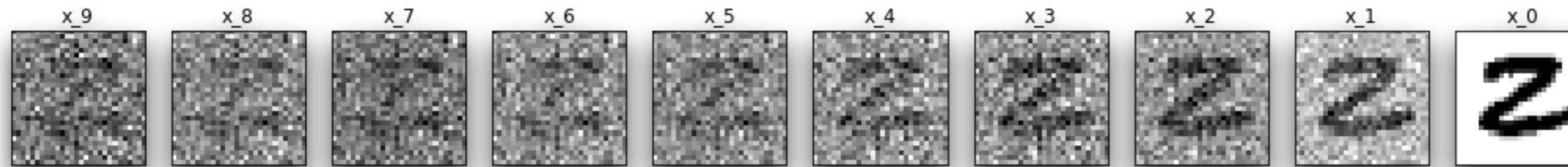
$p_\theta(\mathbf{x}_{t-1} \mid \mathbf{x}_t) := \mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta(\mathbf{x}_t, t), \Sigma_\theta(\mathbf{x}_t, t))$

Loss Function

$$D_{\text{KL}}(q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0) \parallel p_\theta(\mathbf{x}_{t-1} \mid \mathbf{x}_t))$$

Diffusion Model - Forward Process

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



Distribution of \mathbf{x}_t at an arbitrary timestep t in closed form

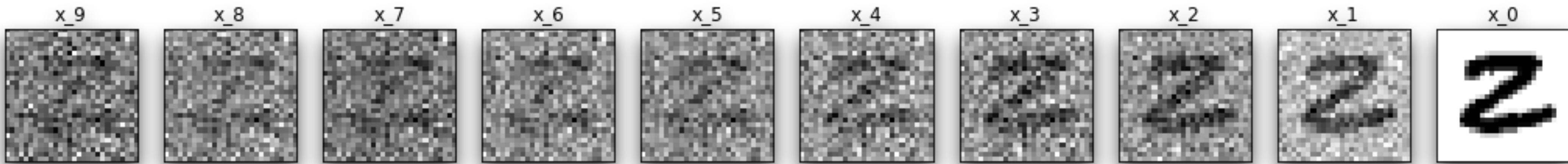
$$q(\mathbf{x}_t \mid \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I}), \text{ where } \alpha_t := 1 - \beta_t \text{ and } \bar{\alpha}_t := \prod_{s=1}^t \alpha_s$$

식 유도는 Lil'Log 참고

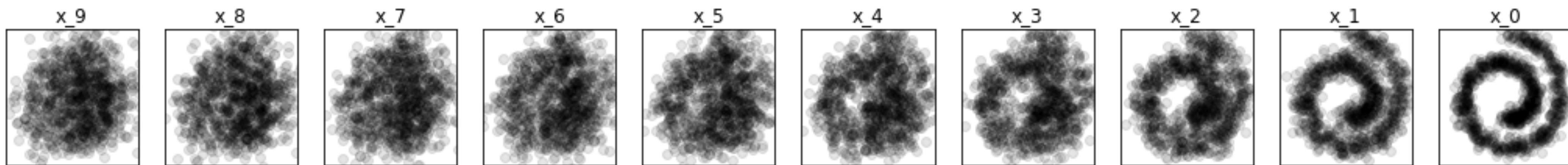
<https://lilianweng.github.io/posts/2021-07-11-diffusion-models/>

Diffusion Model - Forward Process

MNIST single data, $\beta = 0.2$, $T = 10$

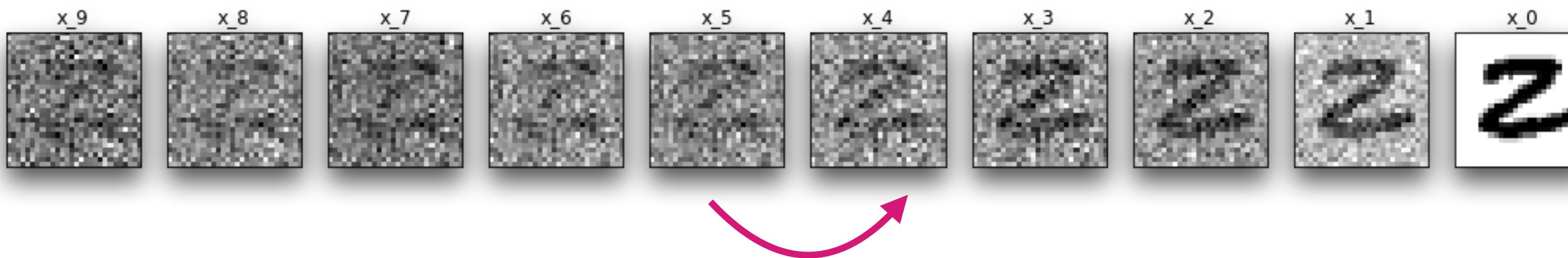


Swiss roll dataset, $\beta = 0.05$, $T = 10$



Diffusion Model - Posterior

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



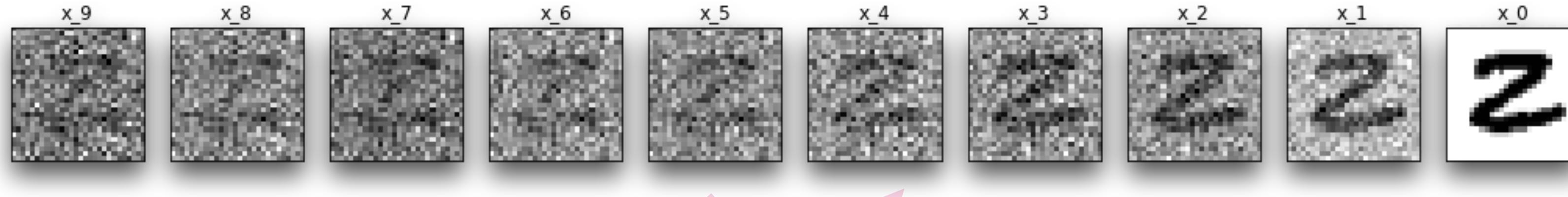
Posterior $q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_{t-1}; \tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0), \tilde{\beta}_t\mathbf{I})$

$$\text{where } \tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0) := \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1-\bar{\alpha}_t}\mathbf{x}_0 + \frac{\sqrt{\alpha_t}(1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_t}\mathbf{x}_t \quad \text{and} \quad \tilde{\beta}_t := \frac{1-\bar{\alpha}_{t-1}}{1-\bar{\alpha}_t}\beta_t$$

$$q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0) = \frac{q(\mathbf{x}_{t-1} \mid \mathbf{x}_0)q(\mathbf{x}_t \mid \mathbf{x}_{t-1}, \mathbf{x}_0)}{q(\mathbf{x}_t \mid \mathbf{x}_0)} \text{ by Bayes' Rule}$$

Diffusion Model - Backward Process

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



Posterior $q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_{t-1}; \tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0), \tilde{\beta}_t\mathbf{I})$

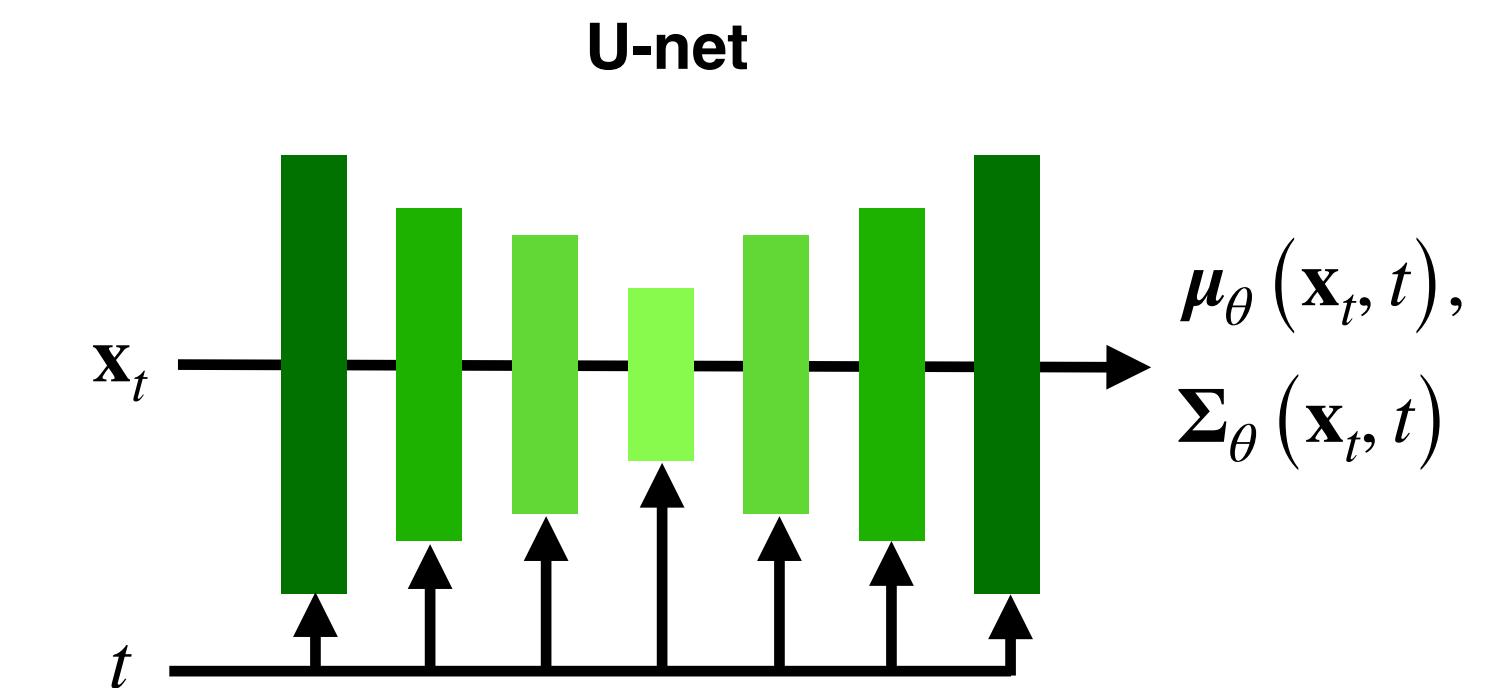
$$\text{where } \tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0) := \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1-\bar{\alpha}_t}\mathbf{x}_0 + \frac{\sqrt{\bar{\alpha}_t}(1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_t}\mathbf{x}_t \quad \text{and} \quad \tilde{\beta}_t := \frac{1-\bar{\alpha}_{t-1}}{1-\bar{\alpha}_t}\beta_t$$

**Backward Process
(Neural Networks)**

$p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t) := \mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta(\mathbf{x}_t, t), \Sigma_\theta(\mathbf{x}_t, t))$

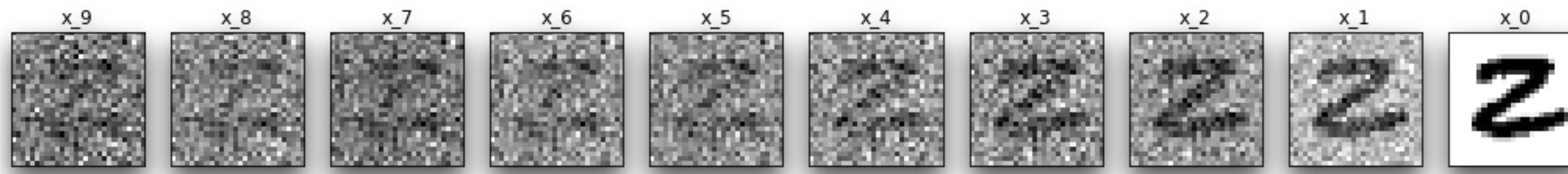
Loss Function

$$D_{\text{KL}}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \| p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t))$$



Diffusion Model - Loss Function

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



Posterior $q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_{t-1}; \tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0), \tilde{\beta}_t\mathbf{I})$

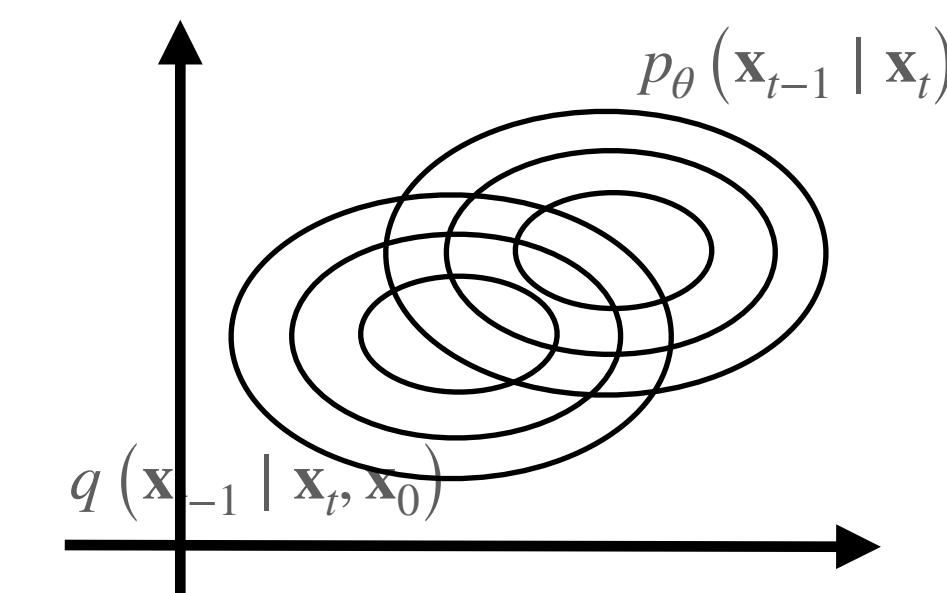
$$\text{where } \tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0) := \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1-\bar{\alpha}_t}\mathbf{x}_0 + \frac{\sqrt{\bar{\alpha}_t}(1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_t}\mathbf{x}_t \quad \text{and} \quad \tilde{\beta}_t := \frac{1-\bar{\alpha}_{t-1}}{1-\bar{\alpha}_t}\beta_t$$

**Backward Process
(Neural Networks)**

$p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t) := \mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta(\mathbf{x}_t, t), \Sigma_\theta(\mathbf{x}_t, t))$

Loss Function

$$D_{\text{KL}}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \| p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t))$$



Diffusion Model - Output Samples

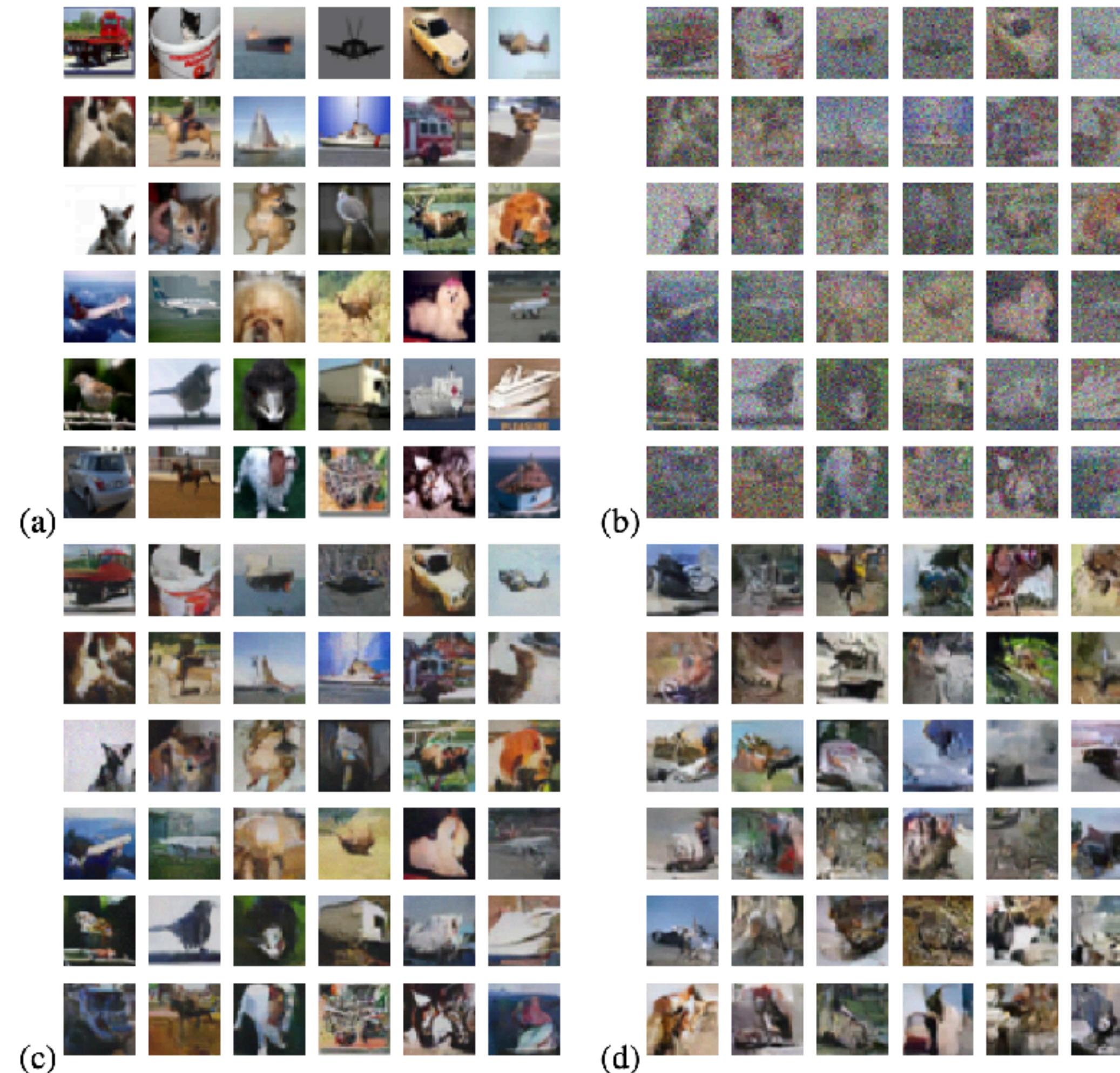


Figure 3. The proposed framework trained on the CIFAR-10 ([Krizhevsky & Hinton, 2009](#)) dataset. (a) Example holdout data (similar to training data). (b) Holdout data corrupted with Gaussian noise of variance 1 (SNR = 1). (c) Denoised images, generated by sampling from the posterior distribution over denoised images conditioned on the images in (b). (d) Samples generated by the diffusion model.

DDPM (Denoising Diffusion Probabilistic Models)

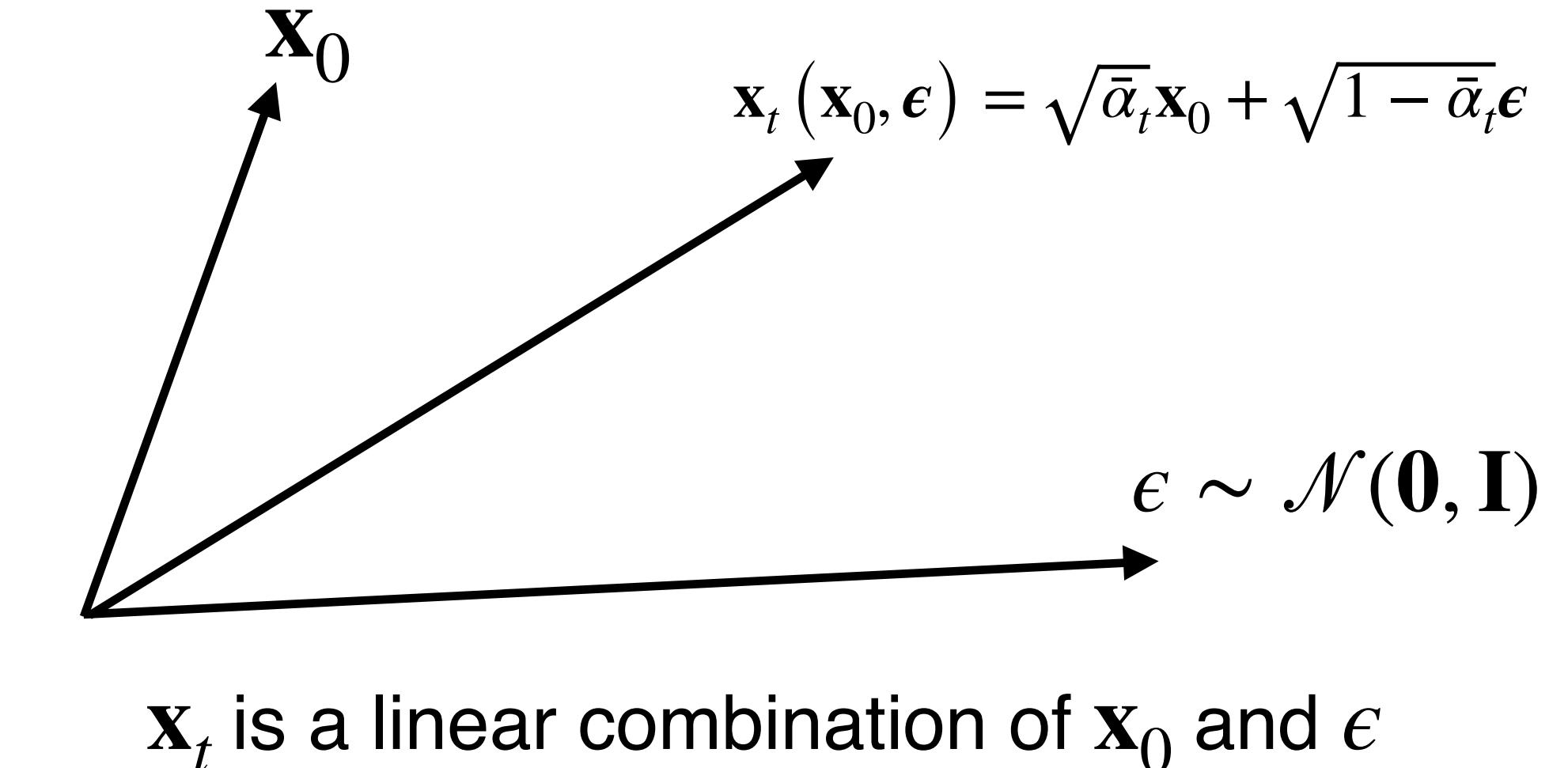
- Jonathan Ho의 논문 Denoising diffusion probabilistic models에서 제안

Posterior $q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}\left(\mathbf{x}_{t-1}; \tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0), \tilde{\beta}_t \mathbf{I}\right)$

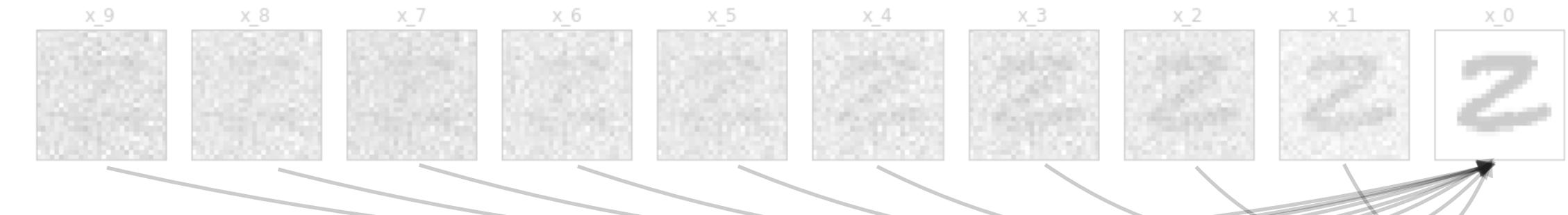
Distribution of \mathbf{x}_t at an arbitrary timestep t in closed form

$$q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}\left(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I}\right), \text{ where } \alpha_t := 1 - \beta_t \text{ and } \bar{\alpha}_t := \prod_{s=1}^t \alpha_s$$

Loss Function $L_{\text{simple}}(\theta) := \mathbb{E}_{t, \mathbf{x}_0, \epsilon} \left[\left\| \epsilon - \epsilon_\theta \left(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t \right) \right\|^2 \right]$



```
106     def p_loss(self, model, x_0, t, noise=None):
107         if noise is None:
108             noise = torch.randn_like(x_0)           generate epsilon
109
110         x_noise = self.q_sample(x_0, t, noise)    sample x_t
111         x_recon = model(x_noise, t)               predict epsilon
112
113         return F.mse_loss(x_recon, noise)
```



Predict ϵ (or \mathbf{x}_0) at each step

<https://github.com/robinly/denoising-diffusion-pytorch/blob/master/diffusion.py>

DDPM (Denoising Diffusion Probabilistic Models)

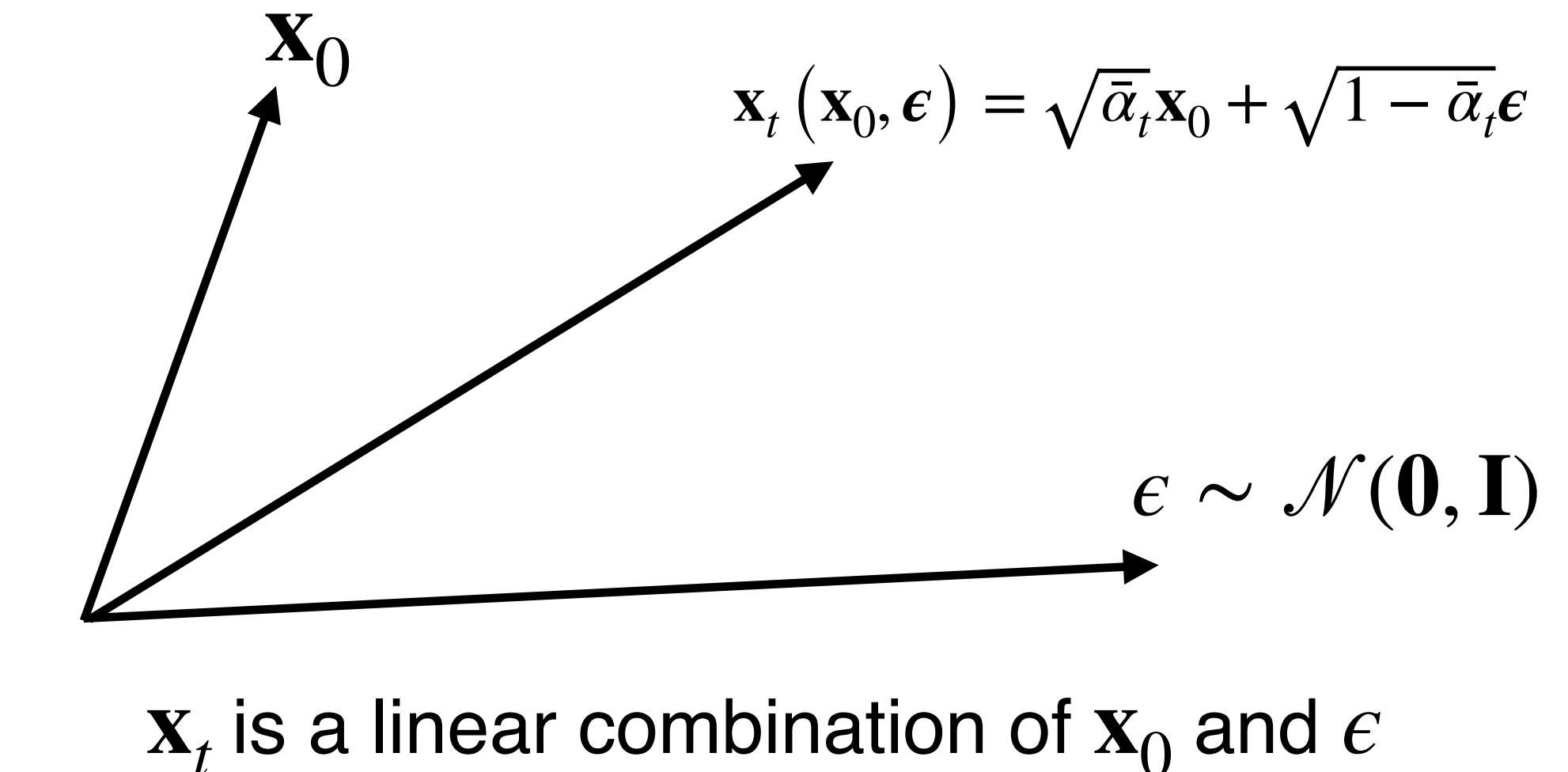
- Jonathan Ho의 논문 Denoising diffusion probabilistic models에서 제안

Posterior $q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}\left(\mathbf{x}_{t-1}; \tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0), \tilde{\beta}_t \mathbf{I}\right)$

Distribution of \mathbf{x}_t at an arbitrary timestep t in closed form

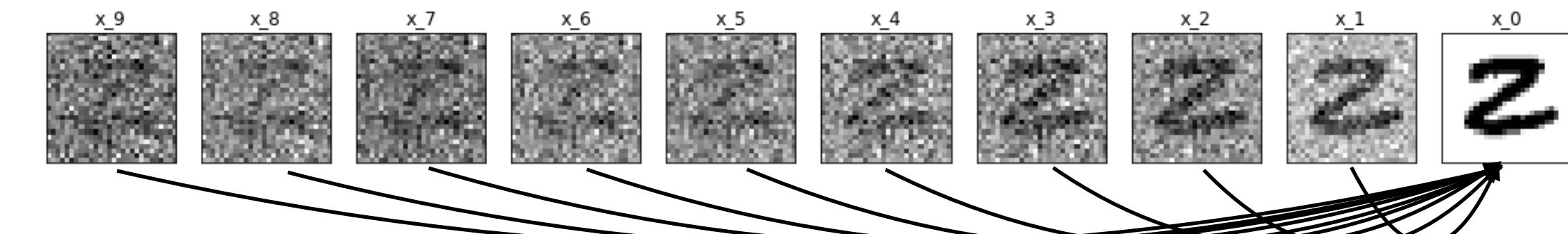
$$q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}\left(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I}\right), \text{ where } \alpha_t := 1 - \beta_t \text{ and } \bar{\alpha}_t := \prod_{s=1}^t \alpha_s$$

Loss Function $L_{\text{simple}}(\theta) := \mathbb{E}_{t, \mathbf{x}_0, \epsilon} \left[\left\| \epsilon - \epsilon_\theta \left(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t \right) \right\|^2 \right]$



```

106     def p_loss(self, model, x_0, t, noise=None):
107         if noise is None:
108             noise = torch.randn_like(x_0)           generate epsilon
109
110             x_noise = self.q_sample(x_0, t, noise) sample x_t
111             x_recon = model(x_noise, t)           predict epsilon
112
113             return F.mse_loss(x_recon, noise)
    
```



<https://github.com/robinly/denoising-diffusion-pytorch/blob/master/diffusion.py>

DDPM - Output Samples

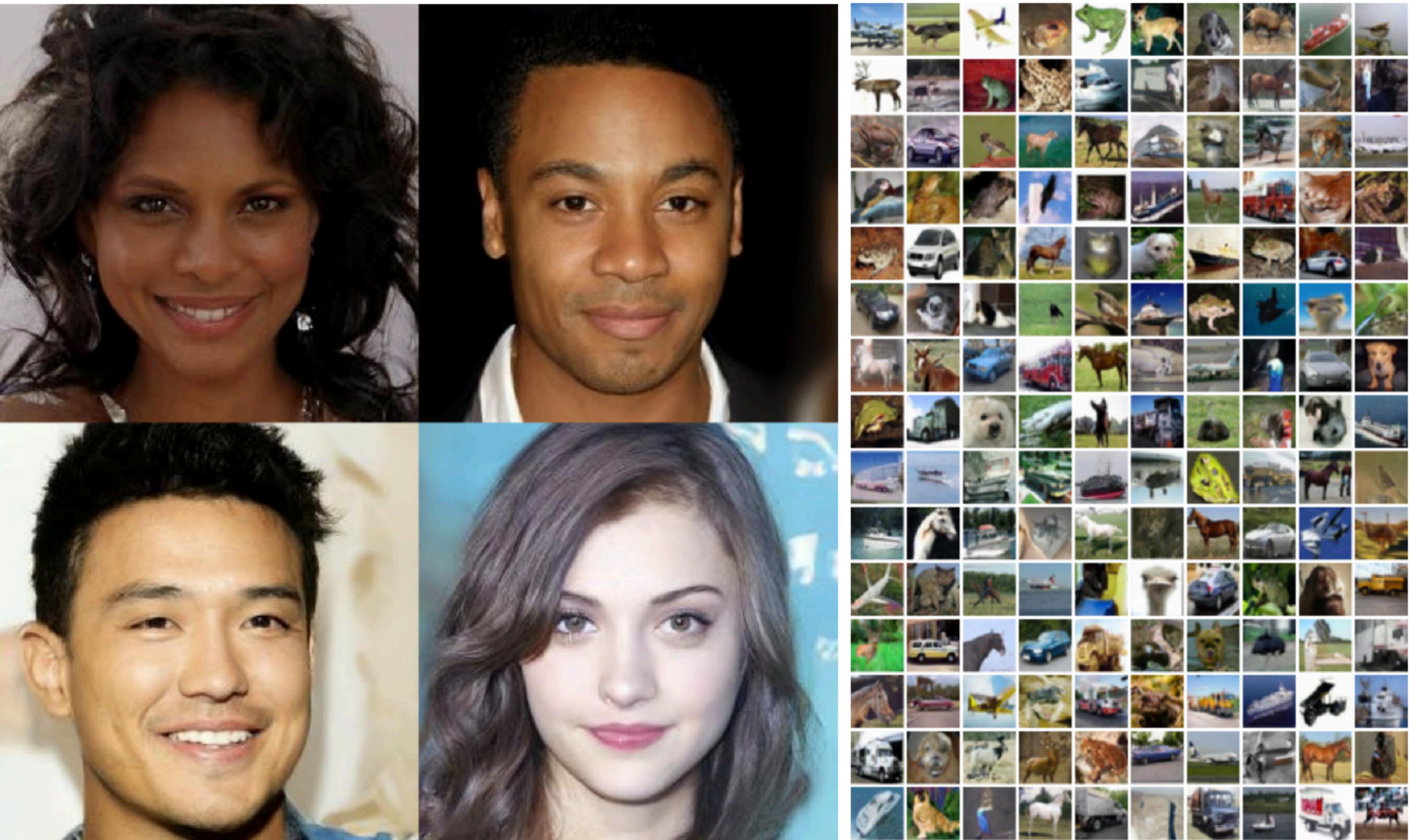
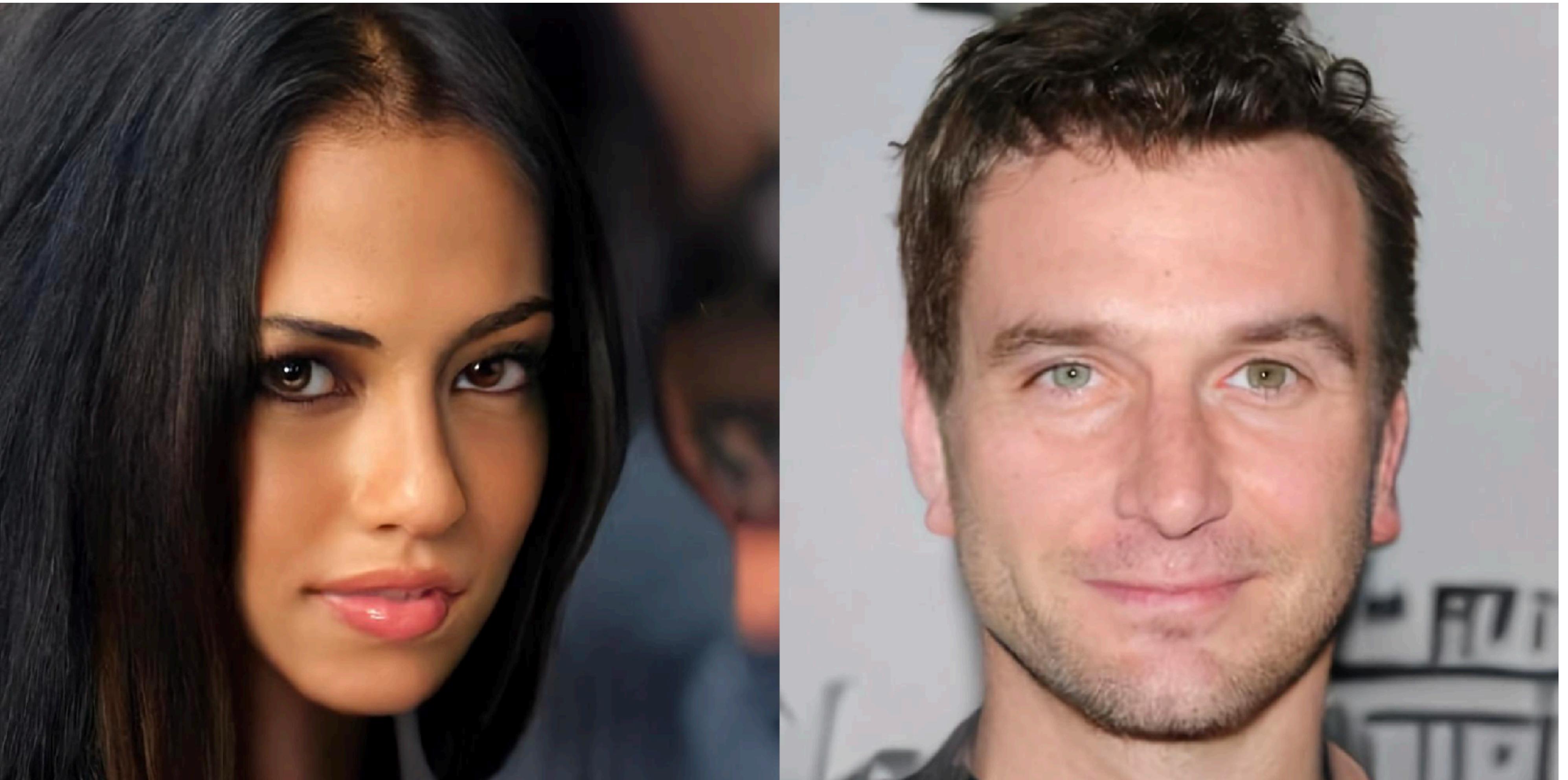


Figure 1: Generated samples on CelebA-HQ 256×256 (left) and unconditional CIFAR10 (right)

DDPM - Output Samples



Diffusion vs. ...

"We present a novel way to define probabilistic models that allows:

1. **extreme flexibility in model structure,**
2. exact sampling,
3. easy multiplication with other distributions, e.g. in order to compute a posterior, and
4. **the model log likelihood, and the probability of individual states, to be cheaply evaluated.**

Sohl-Dickstein, Jascha, et al. "Deep unsupervised learning using nonequilibrium thermodynamics." *International Conference on Machine Learning*. PMLR, 2015.

"We emphasize that our objective Eq. (6) requires **no adversarial training, no surrogate losses, and no sampling from the score network during training** (e.g., unlike contrastive divergence). Also, it **does not require $s\theta(x, \sigma)$ to have special architectures** in order to be tractable."

Song, Yang, and Stefano Ermon. "Generative modeling by estimating gradients of the data distribution." *Advances in Neural Information Processing Systems* 32 (2019).

Diffusion vs. ...

	Auto- Regressive	VAE	Flow	GAN	Diffusion
Tractability	Good	Good	Good	Not Good <small>Likelihood can't be evaluated</small>	Good
Flexibility	Not Good <small>Causal structure Fixed distribution</small>	Not Good <small>Dimension reduction Fixed distribution</small>	Not Good <small>Invertible structure Fixed distribution</small>	Good	Good

**GLIDE
CLIP
DALL·E 2**

GLIDE - Output Samples

- OpenAI에서 논문 Glide: Towards photorealistic image generation and editing with text-guided diffusion models을 통해 제안



“a hedgehog using a calculator”



“a corgi wearing a red bowtie and a purple party hat”



“robots meditating in a vipassana retreat”



“a fall landscape with a small cottage next to a lake”



“a surrealist dream-like oil painting by salvador dali of a cat playing checkers”



“a professional photo of a sunset behind the grand canyon”

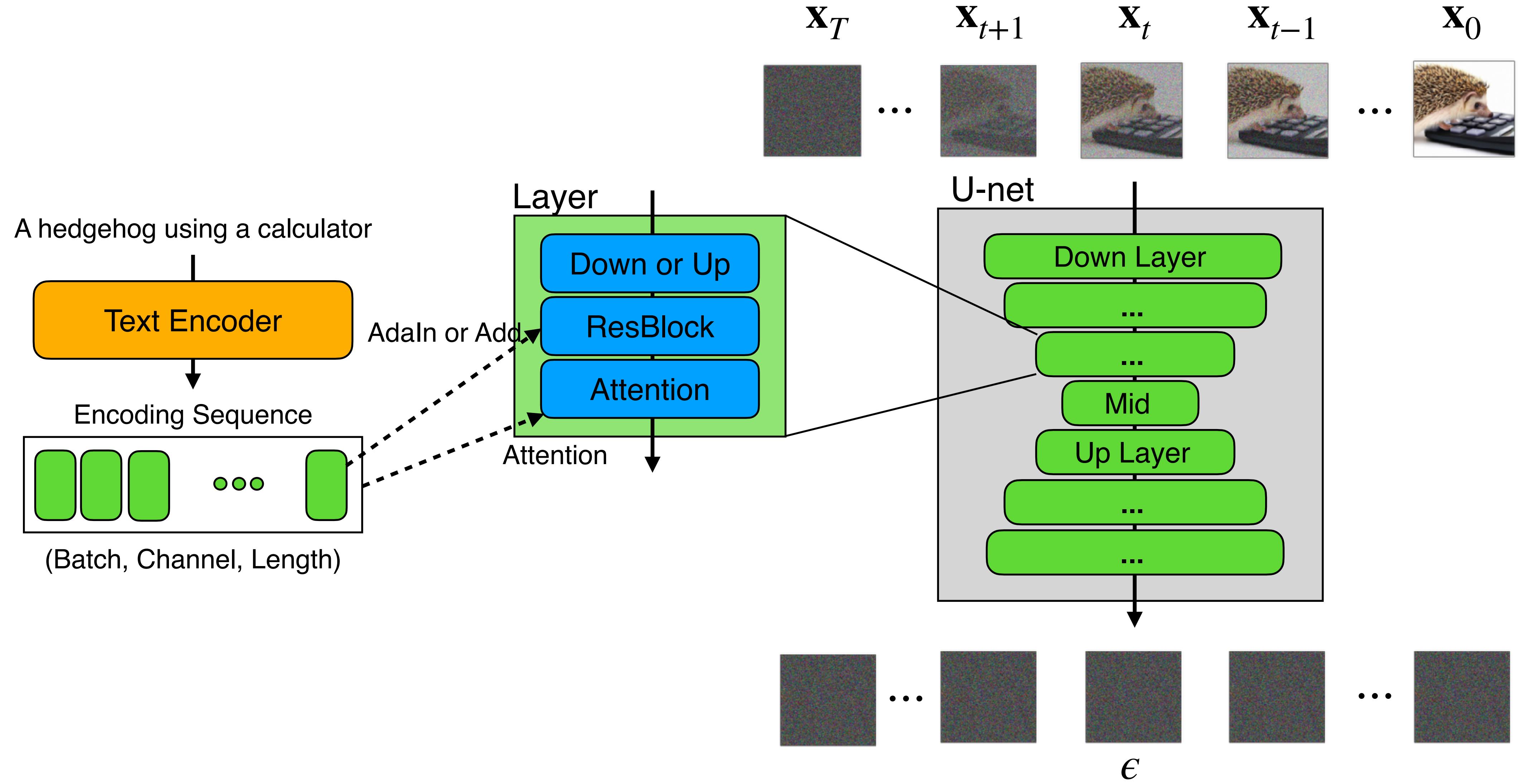


“a high-quality oil painting of a psychedelic hamster dragon”



“an illustration of albert einstein wearing a superhero costume”

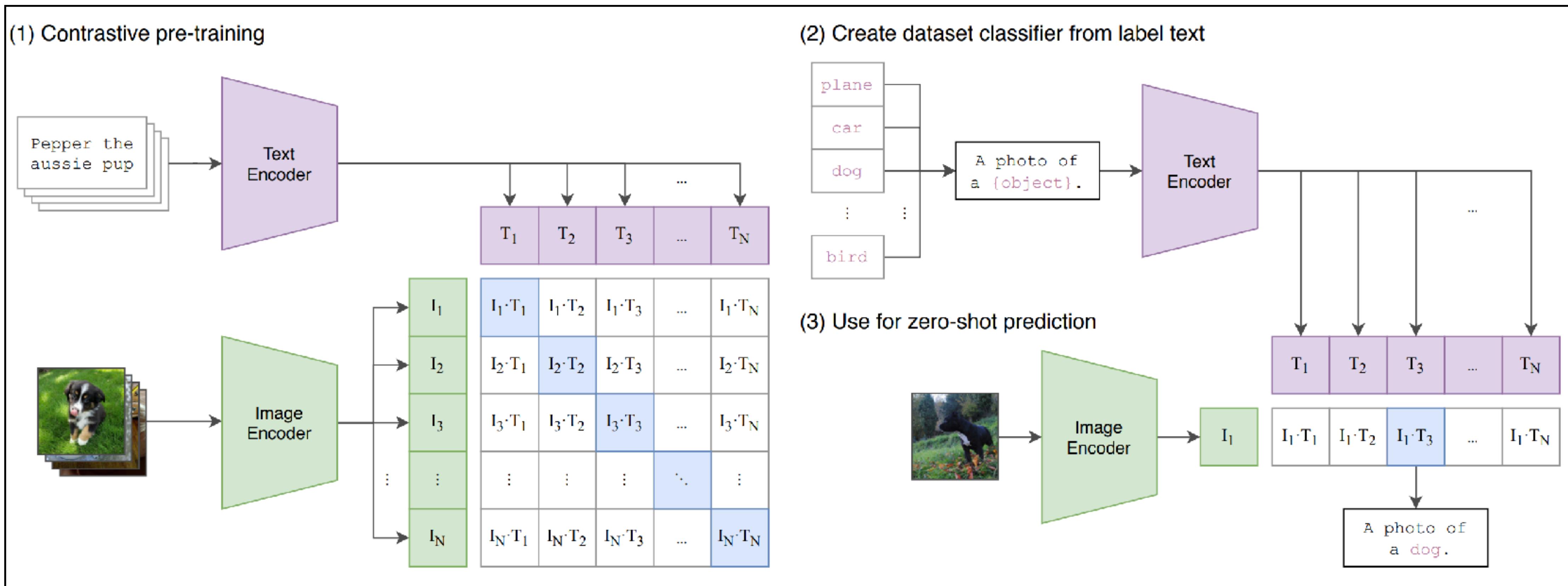
GLIDE - Overall Architecture



GLIDE source : <https://github.com/openai/glide-text2im>

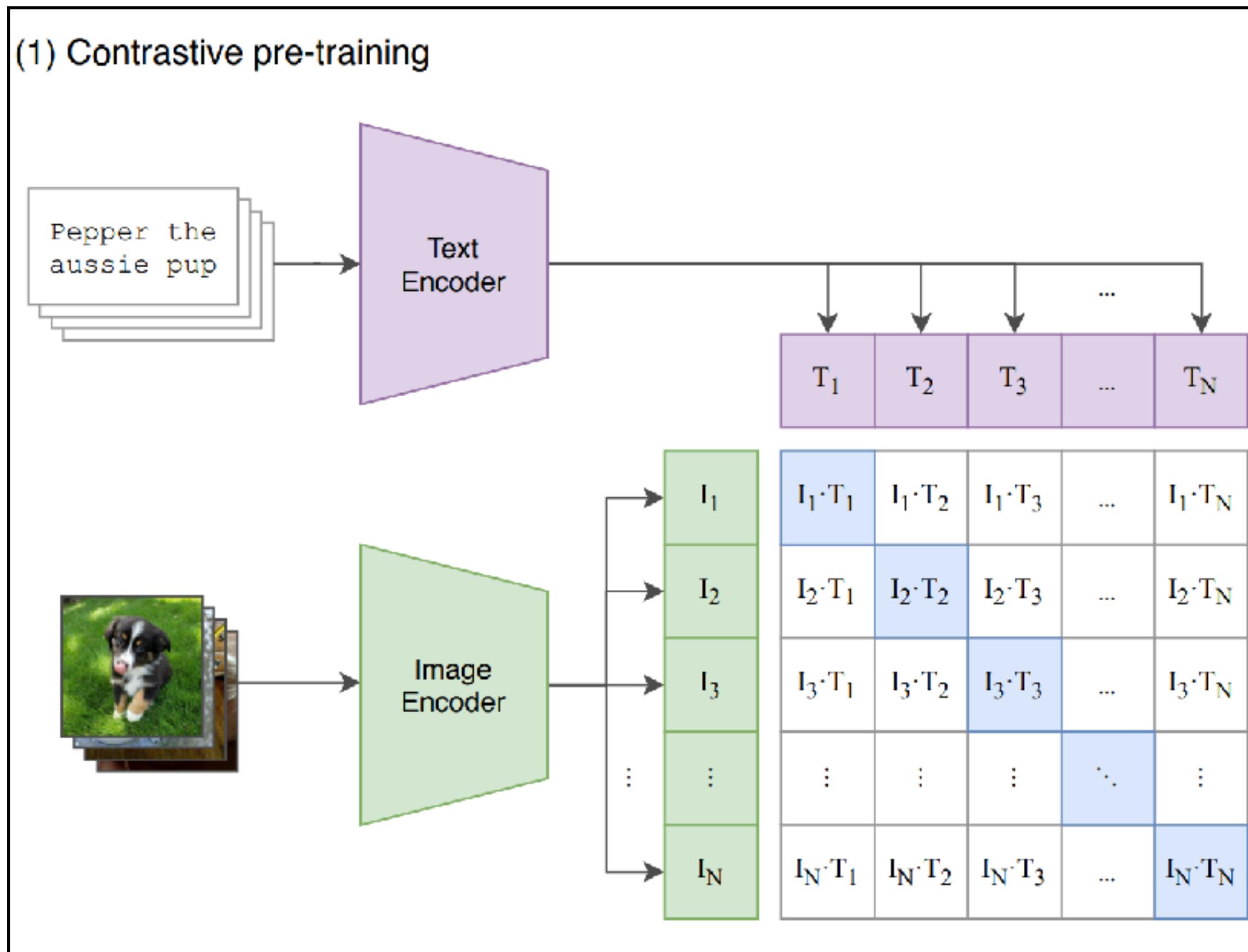
CLIP (Contrastive Language-Image Pre-training)

- (image, text) 쌍의 데이터로 self-supervised learning을 통해 text/image encoder를 학습



CLIP (Contrastive Language-Image Pre-training)

- (image, text) 쌍의 데이터로 self-supervised learning을 통해 text/image encoder를 학습



```
# image_encoder - ResNet or Vision Transformer
# text_encoder - CBOW or Text Transformer
# I[n, h, w, c] - minibatch of aligned images
# T[n, 1] - minibatch of aligned texts
# W_i[d_i, d_e] - learned proj of image to embed
# W_t[d_t, d_e] - learned proj of text to embed
# t - learned temperature parameter

# extract feature representations of each modality
I_f = image_encoder(I) #[n, d_i]
T_f = text_encoder(T) #[n, d_t]

# joint multimodal embedding [n, d_e]
I_e = l2_normalize(np.dot(I_f, W_i), axis=1)
T_e = l2_normalize(np.dot(T_f, W_t), axis=1)

# scaled pairwise cosine similarities [n, n]
logits = np.dot(I_e, T_e.T) * np.exp(t)

# symmetric loss function
labels = np.arange(n)
loss_i = cross_entropy_loss(logits, labels, axis=0)
loss_t = cross_entropy_loss(logits, labels, axis=1)
loss = (loss_i + loss_t)/2
```

CLIP (Contrastive Language-Image Pre-training)

FOOD101

guacamole (90.1%) Ranked 1 out of 101 labels



- ✓ a photo of **guacamole**, a type of food.
- ✗ a photo of **ceviche**, a type of food.
- ✗ a photo of **edamame**, a type of food.
- ✗ a photo of **tuna tartare**, a type of food.
- ✗ a photo of **hummus**, a type of food.

SUN397

television studio (90.2%) Ranked 1 out of 397



- ✓ a photo of a **television studio**.
- ✗ a photo of a **podium indoor**.
- ✗ a photo of a **conference room**.
- ✗ a photo of a **lecture room**.
- ✗ a photo of a **control room**.

DALL·E 2 - Output Samples

- OpenAI에서 Hierarchical text-conditional image generation with clip latents 논문을 통해 제안



vibrant portrait painting of Salvador Dalí with a robotic half face



a shiba inu wearing a beret and black turtleneck



a close up of a handpalm with leaves growing from it



an espresso machine that makes coffee from human souls, artstation

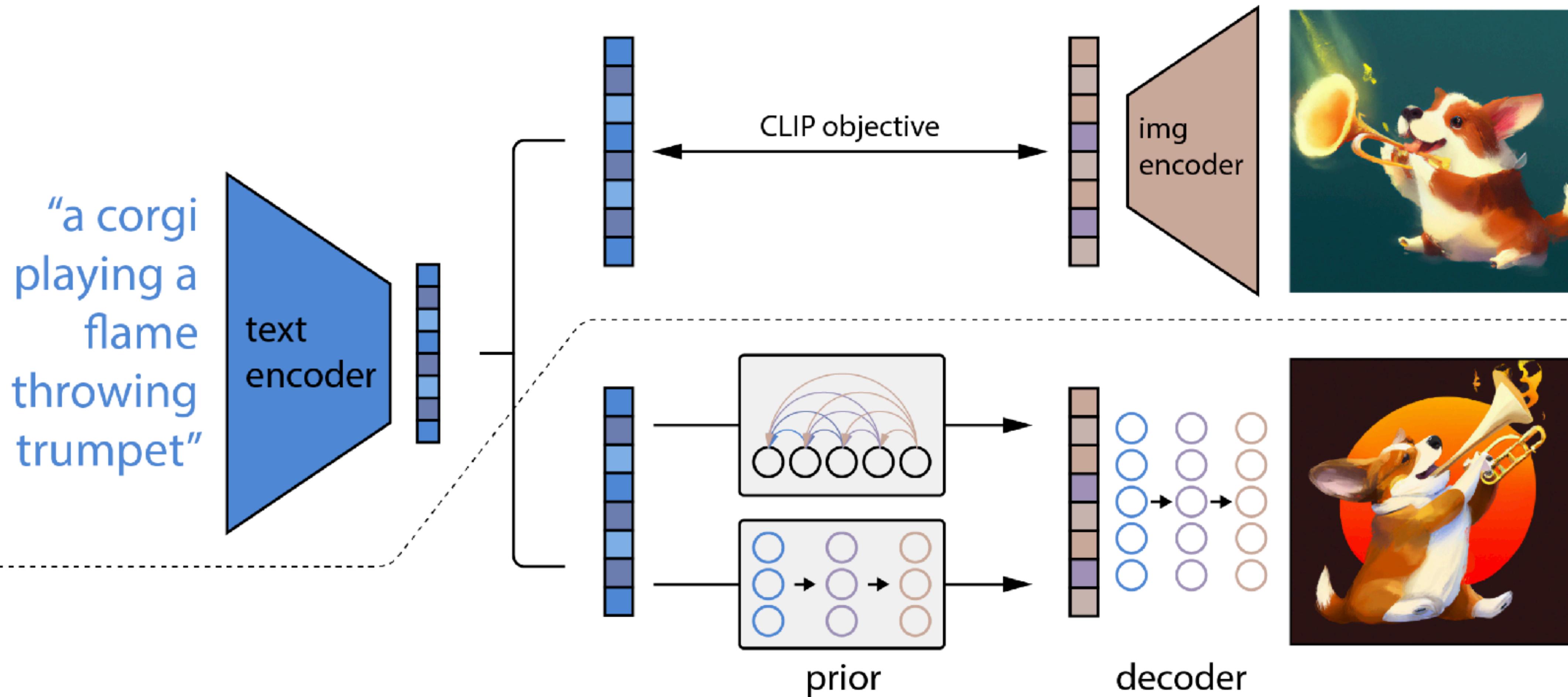


panda mad scientist mixing sparkling chemicals, artstation



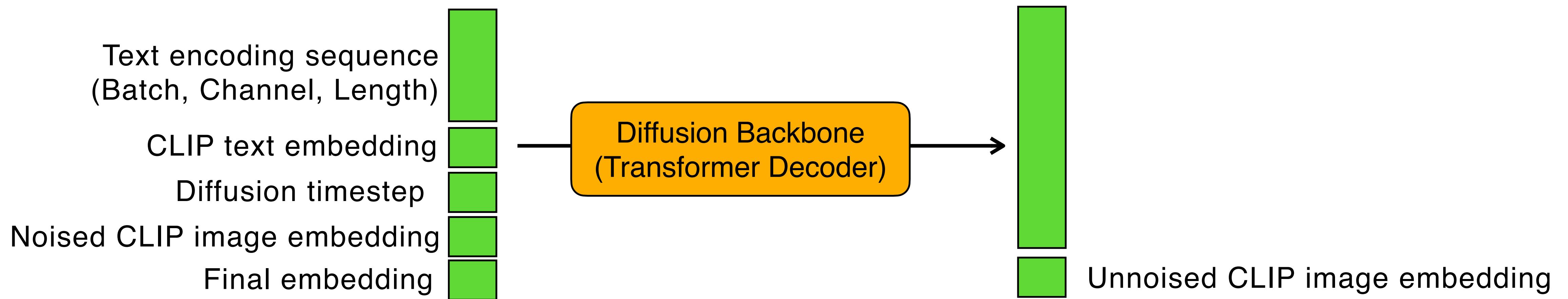
a corgi's head depicted as an explosion of a nebula

DALL·E 2 - Overall Architecture



DALL·E 2 - Prior

- GLIDE text embedding과 CLIP text embedding으로부터 CLIP image embedding을 생성



Guided Diffusion Sampling

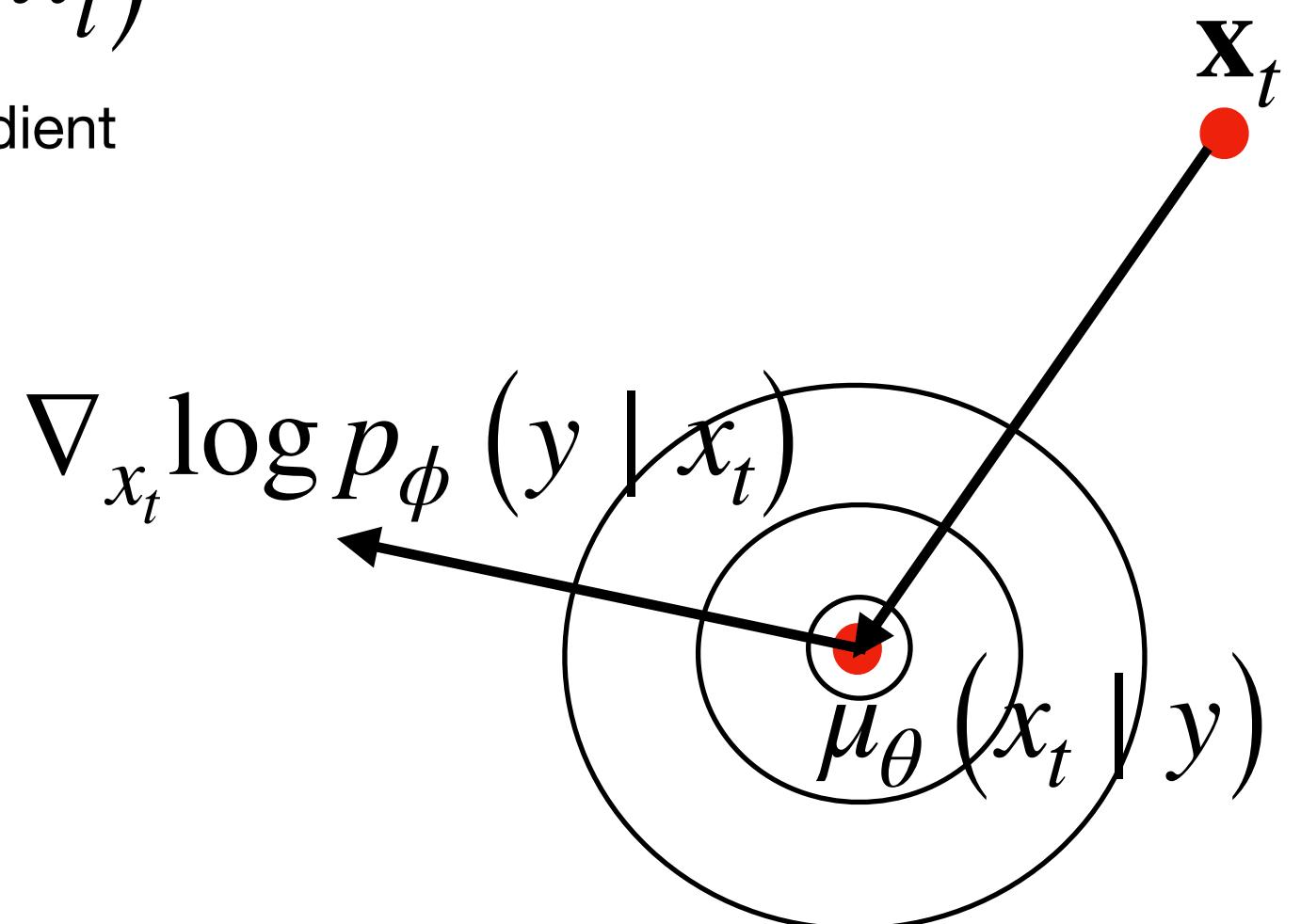
Classifier-free Diffusion Guidance

CLIP Guidance

Guided Diffusion Sampling

- Diffusion Models Beat GANs on Image Synthesis 논문에서 제안
 - Diffusion 모델 외 추가적인 image classifier를 학습시키고 sampling 과정에서 classifier로부터 gradient를 받아 sampling에 도움을 줌

$$\text{Posterior} \quad q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}\left(\mathbf{x}_{t-1}; \tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0), \tilde{\beta}_t \mathbf{I}\right)$$



Guided Diffusion Sampling



Figure 13: Samples from our best 512×512 model (FID: 3.85). Classes are 1: goldfish, 279: arctic fox, 323: monarch butterfly, 386: african elephant, 130: flamingo, 852: tennis ball.

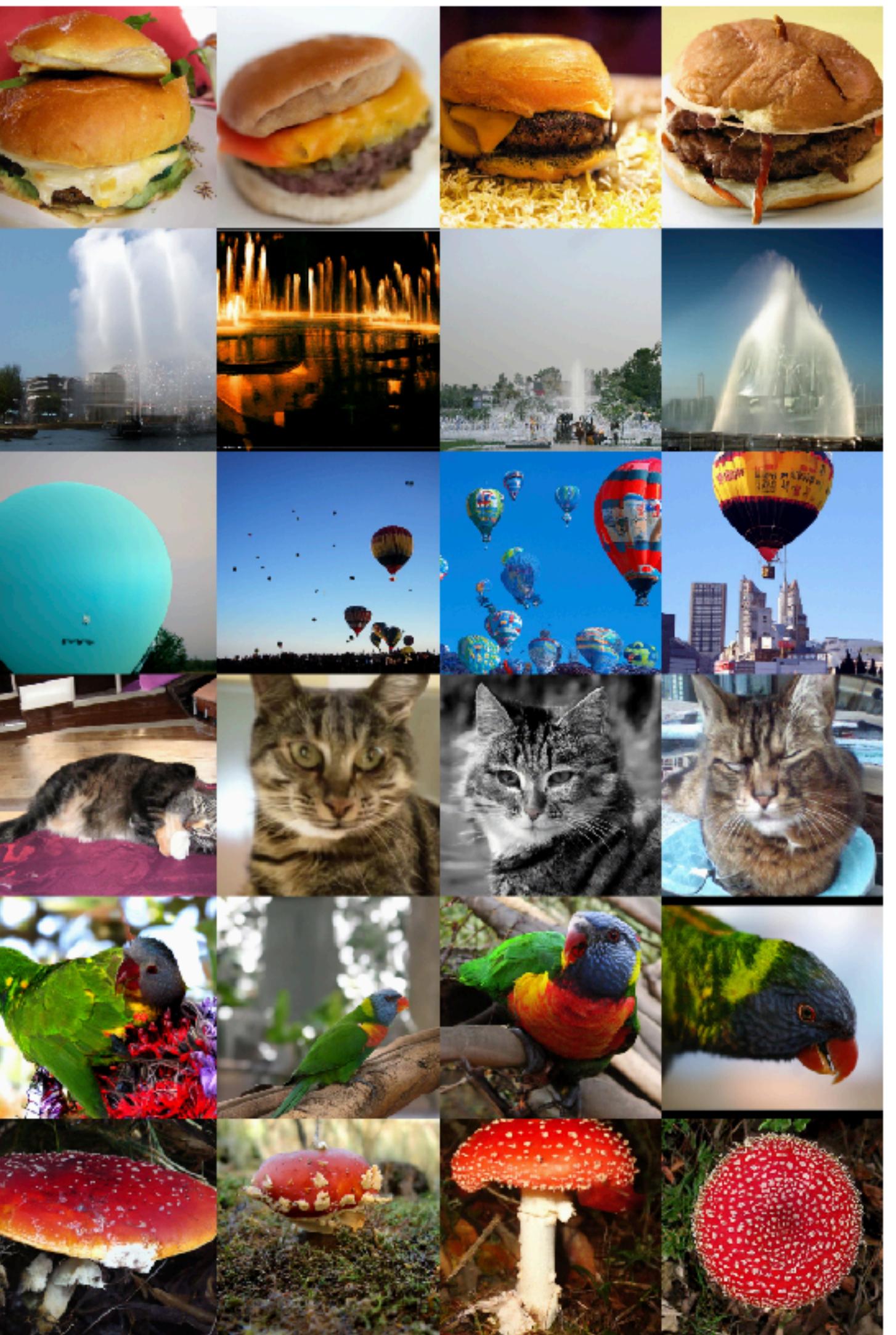


Figure 14: Samples from our best 512×512 model (FID: 3.85). Classes are 933: cheeseburger, 562: fountain, 417: balloon, 281: tahby cat, 90: lorikeet, 992: agaric.



Figure 15: Difficult class samples from our best 512×512 model (FID: 3.85). Classes are 432: bassoon, 468: cab, 424: barbershop, 444: bicycle-built-for-two, 981: ballplayer, 550: espresso maker.

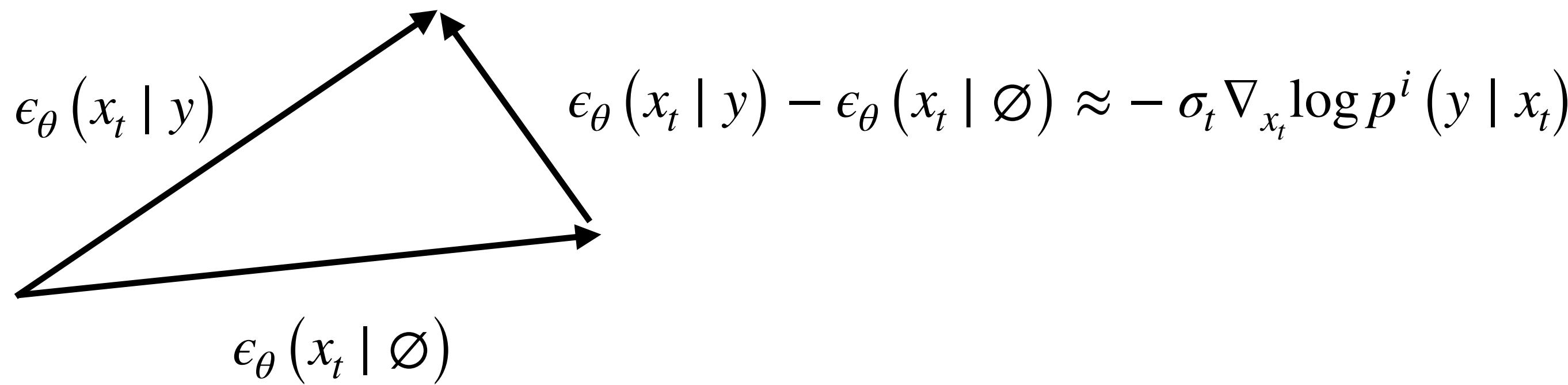
Classifier-free guidance

- Ho, Jonathan과 Tim Salimans의 논문 Classifier-free diffusion guidance에서 제안
- 추가적인 classifier를 트레이닝할 필요없이 diffusion 모델만 가지고 guided sampling을 가능하게 만듬

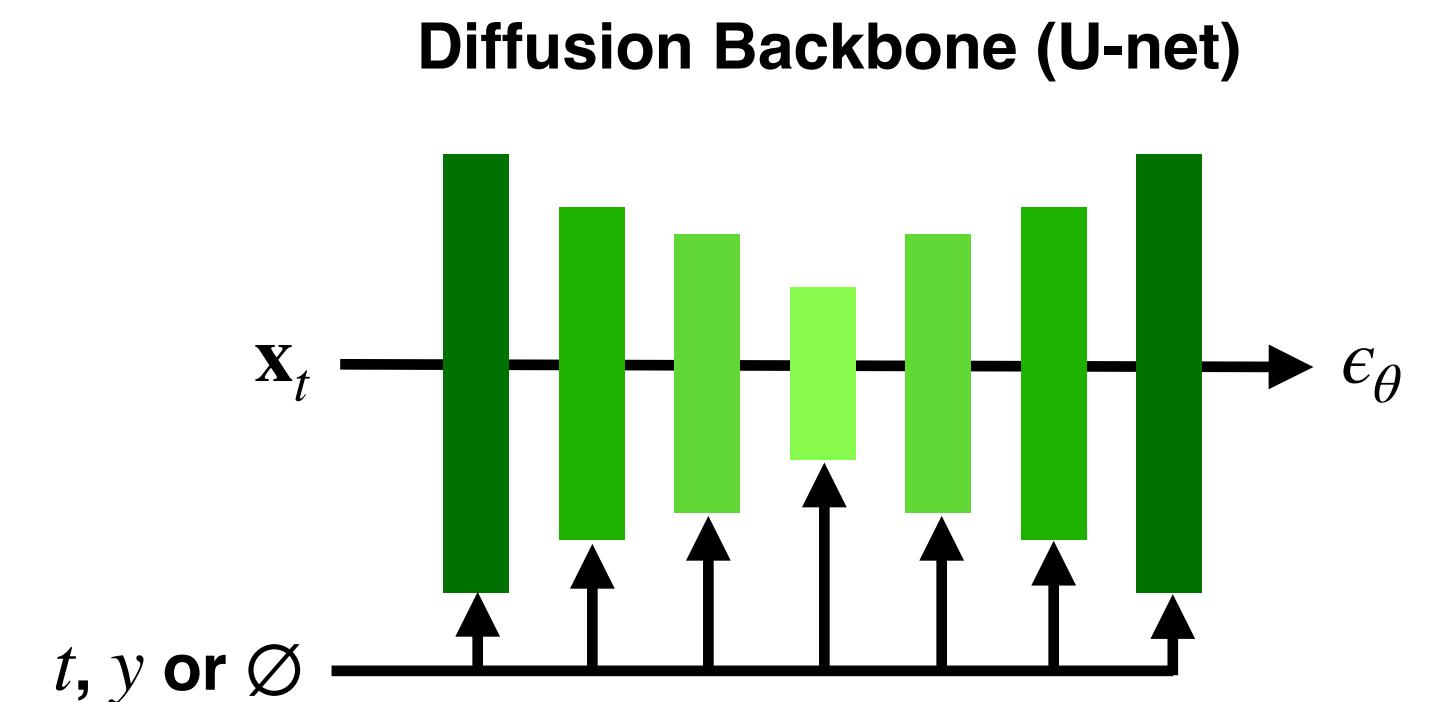
Nichol, Alex, et al.에서 재인용

$$\hat{\epsilon}_\theta(x_t | y) = \epsilon_\theta(x_t | y) + s \cdot (\epsilon_\theta(x_t | y) - \epsilon_\theta(x_t | \emptyset))$$

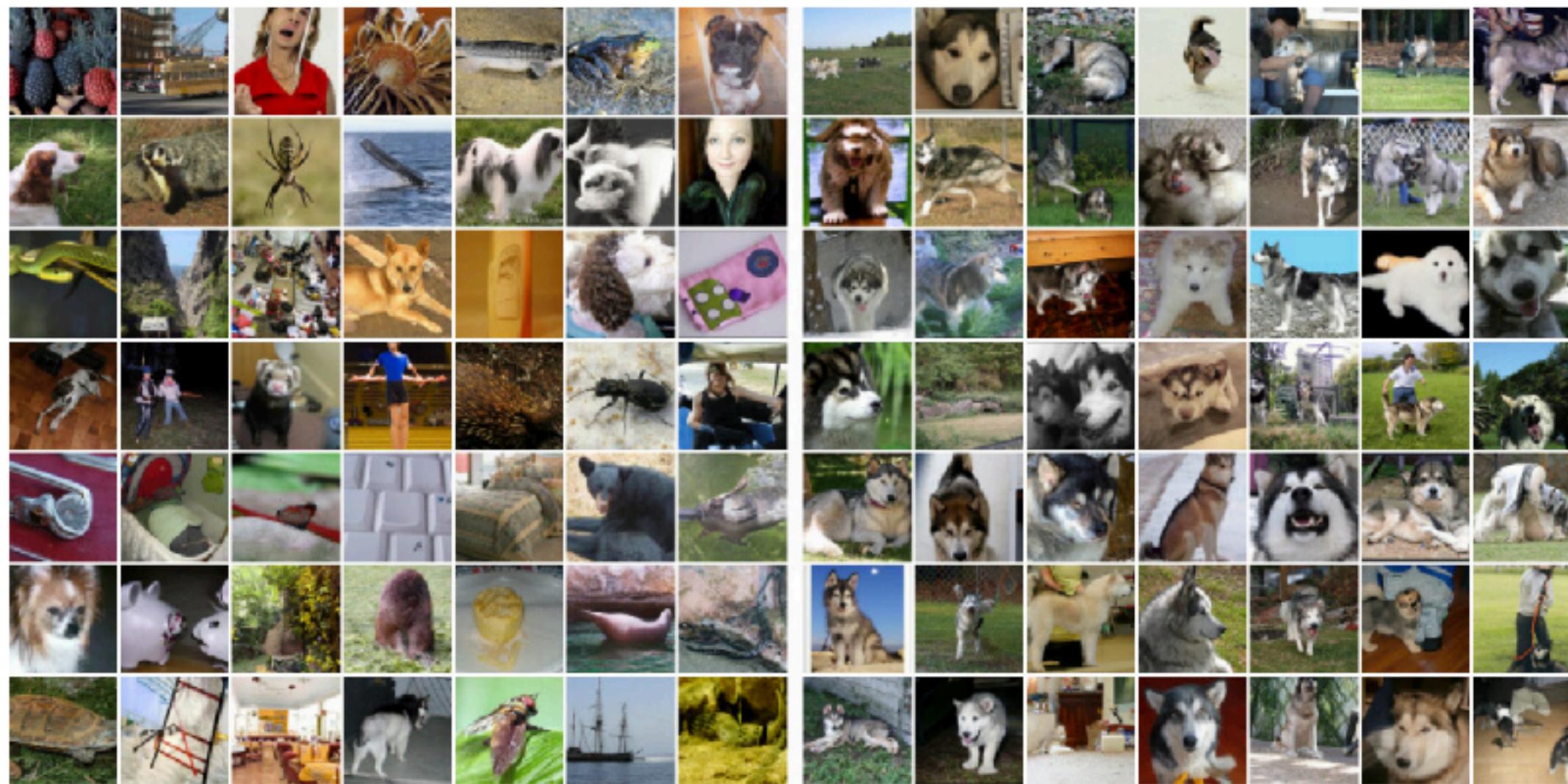
수정된 score conditional predicted score guidance scale unconditional predicted score



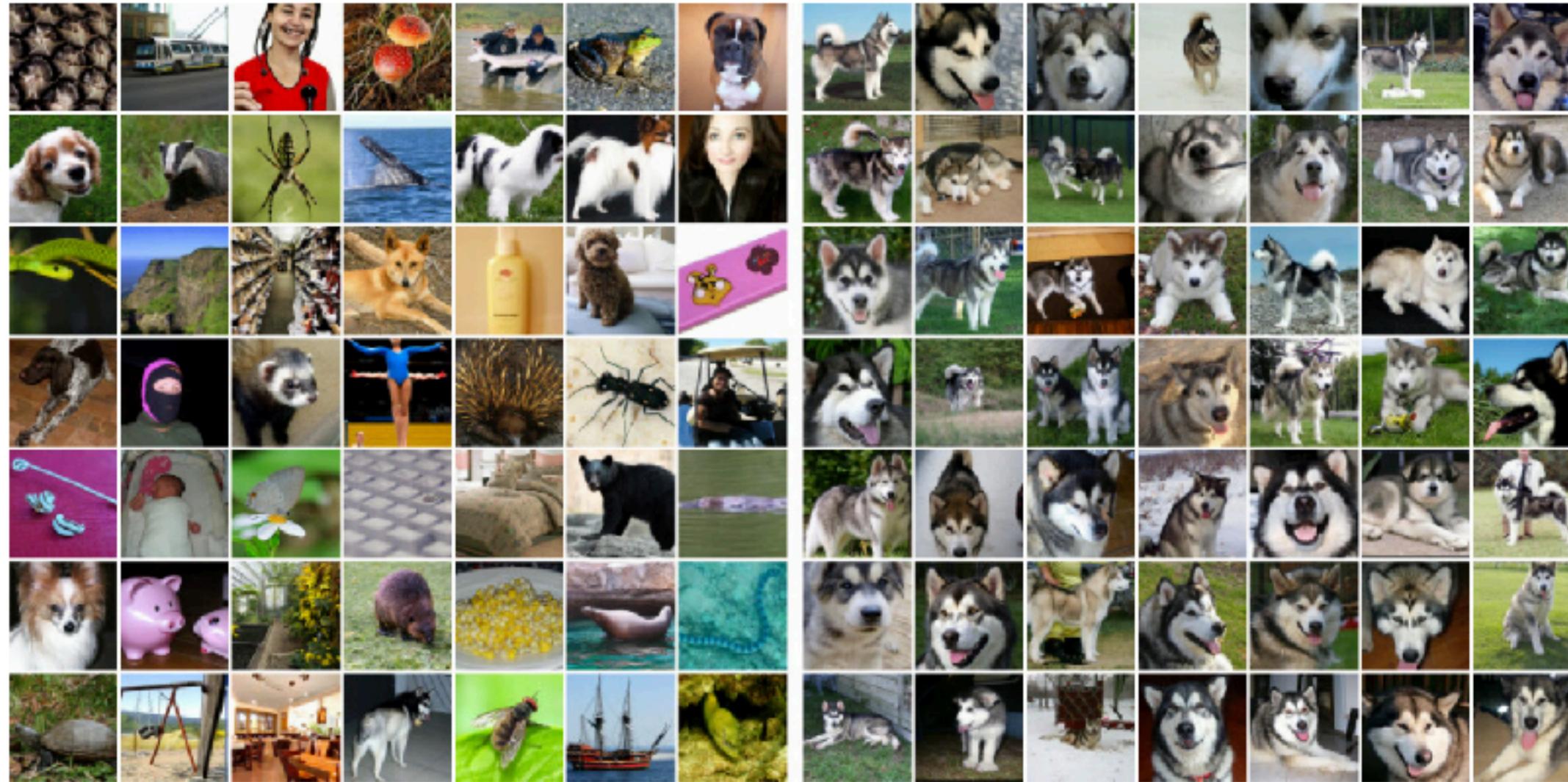
$$\epsilon_\theta(x_t | y) - \epsilon_\theta(x_t | \emptyset) \approx -\sigma_t \nabla_{x_t} \log p^i(y | x_t)$$



Classifier-free guidance



(a) Non-guided conditional sampling: FID=1.80, IS=53.71



(b) Classifier-free guidance with $w = 1.0$: FID=12.6, IS=170.1

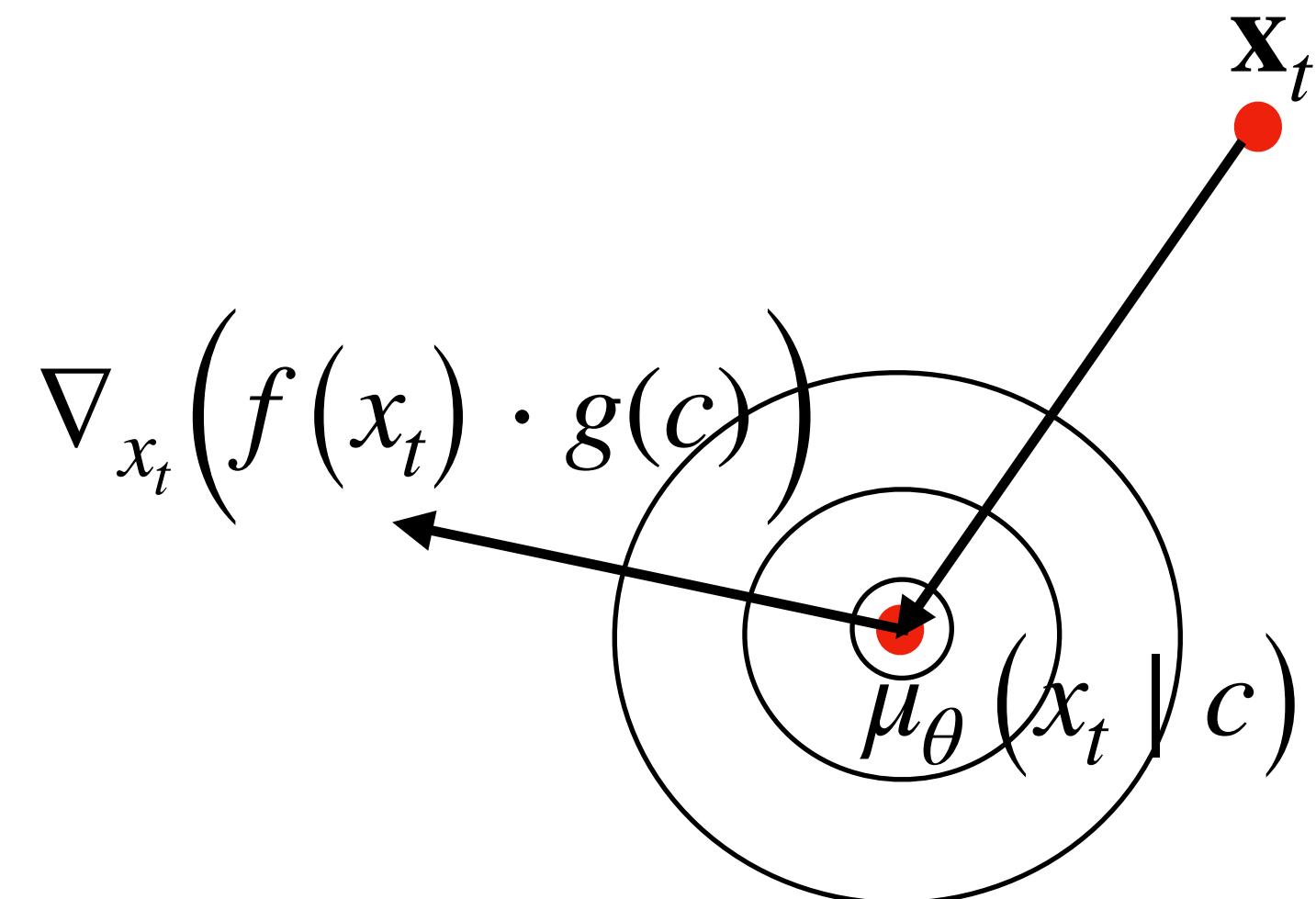
CLIP Guidance

- Classifier를 이용한 guidance 방법과 비슷하게 CLIP 모델을 이용하여 sampling 단계에서 도움을 줌

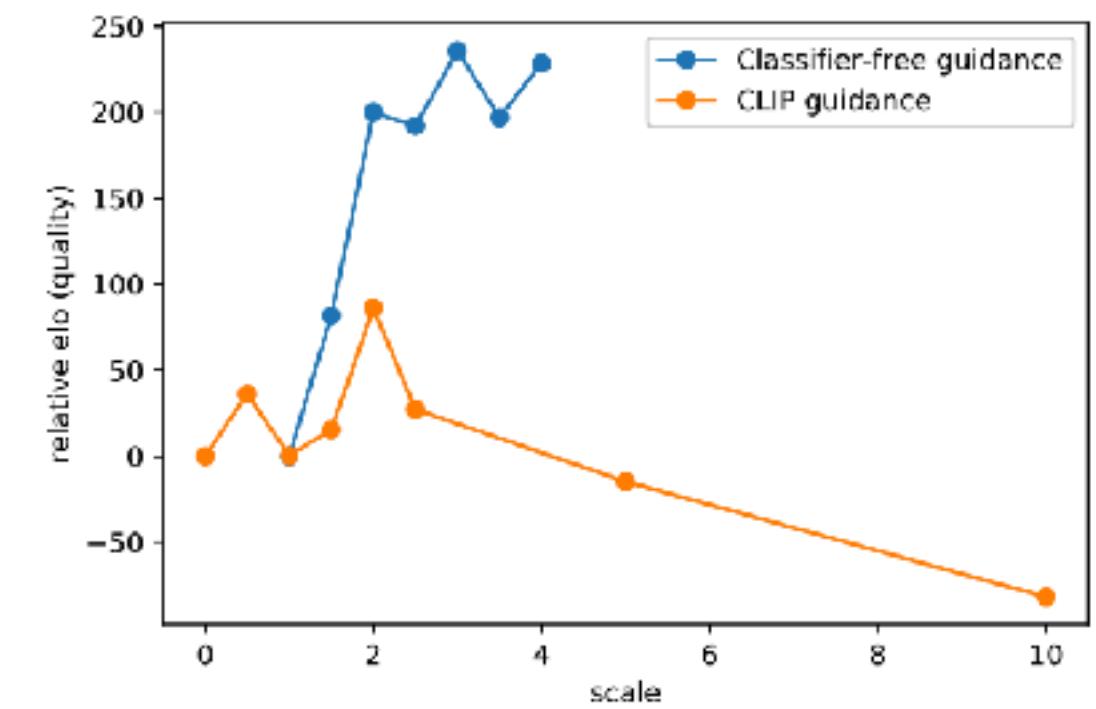
$$\hat{\mu}_\theta(x_t | c) = \mu_\theta(x_t | c) + s \cdot \Sigma_\theta(x_t | c) \nabla_{x_t} (f(x_t) \cdot g(c))$$

CLIP
image encoding

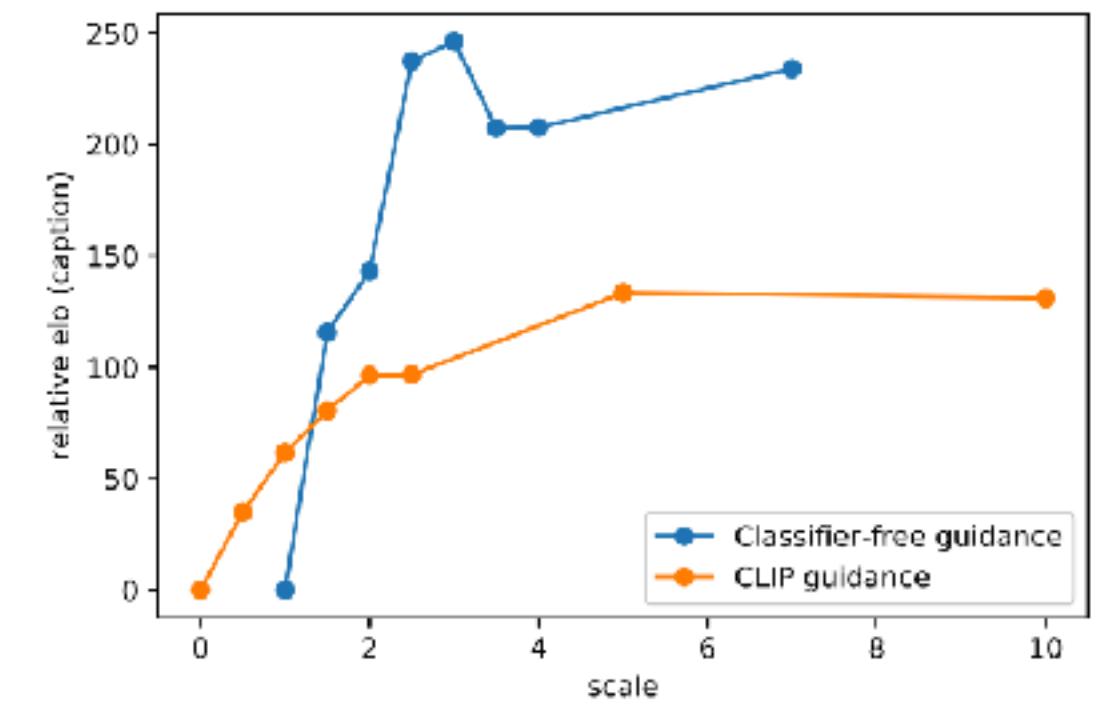
CLIP
text encoding



Classifier-free Guidance vs. CLIP Guidance in GLIDE



(a) Photorealism



(b) Caption Similarity