

Diffusion

Philipp Krähenbühl, UT Austin

Generative models

- Two tasks of a generative model $P(X)$
 - Sampling: $x \sim P(X)$
 - Density estimation: $P(X = x)$



Deep Network

$P(X)$



Deep Network



Generative models

Two kinds of models

Sampling based $x \sim P(X)$

- Sample $z \sim P(Z)$
- Learn transformation
- $P(x|z)$ or $f: z \rightarrow x$

z

Deep
Network



Density estimation based $P(X)$

- Learn special form of $P(X)$
- Model specific sampling / generation

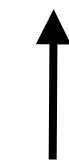


Deep
Network

$P(X)$

Generative modeling is hard

- Density estimation $P(X = x)$
 - How to ensure $\sum_x P(x) = 1$ for all x
- Impossible to compute (in general)



Deep Network

$P(X)$



Deep Network



Auto-regressive models

$$P(x) = P(x_1)P(x_2 | x_1)P(x_3 | x_1, x_2)P(x_4 | x_1 \dots x_3) \dots$$

- $P(x_i | x_1 \dots x_{i-1}) = \text{softmax}(f(x_1 \dots x_{i-1}))$

- Basis of most LLM models

- Easy estimation of $P(x)$

- Easy sampling

$$x_1 \sim P(X_1); x_2 \sim P(X_2 | x_1)$$

- Slow sampling



[1] WaveNet: A Generative Model for Raw Audio. Aaron van den Oord, et al. 2016

[2] Long Video Generation with Time-Agnostic VQGAN and Time-Sensitive Transformer. Songwei Ge, et al. 2022

Auto-regressive models

Issues

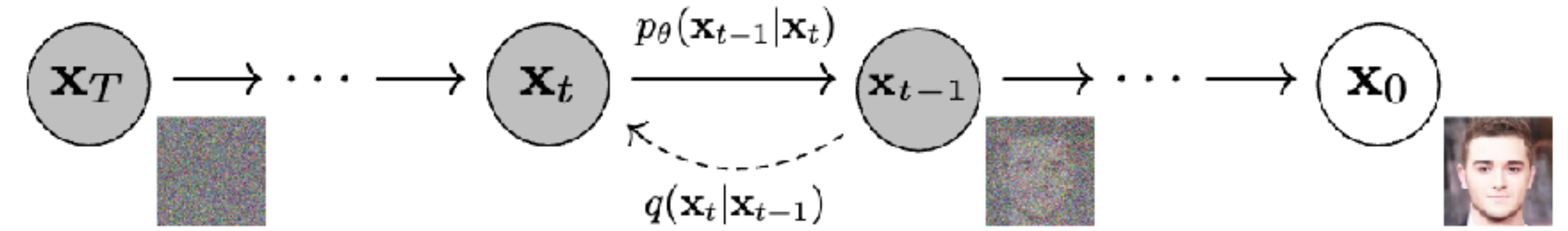
$$P(x) = P(x_1)P(x_2 | x_1)P(x_3 | x_1, x_2)P(x_4 | x_1 \dots x_3) \dots$$

- Difficult learning problem for long sequences (requires good model)
- Solution: Tokenization/Vector-Quantization (next class)
 - More complex x_i
 - Shorter sequence



Diffusion Process

Forward process



- Make an image noisy
 - Start with an image x_0
 - Add noise $q(x_t | x_{t-1}) = \mathcal{N}(\sqrt{1 - \beta_t}x_{t-1}, \beta_t I)$
 - β_t increases linearly with t

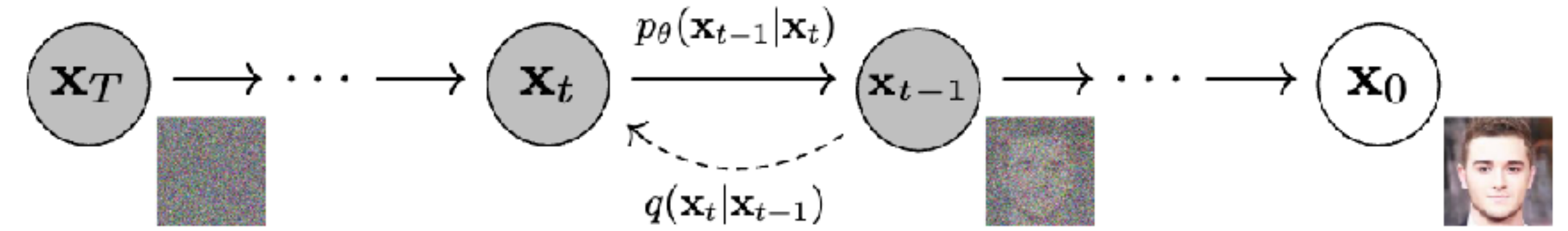
$$q(x_t | x_0) = \int \prod_{i=1}^t q(x_i | x_{i-1}) dx_{1\dots t-1}$$

$$\bullet \quad = \mathcal{N}(\sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)I)$$

$$\text{where } \bar{\alpha}_t = \prod_{i=1}^t (1 - \beta_i)$$

Diffusion Process

Reverse process



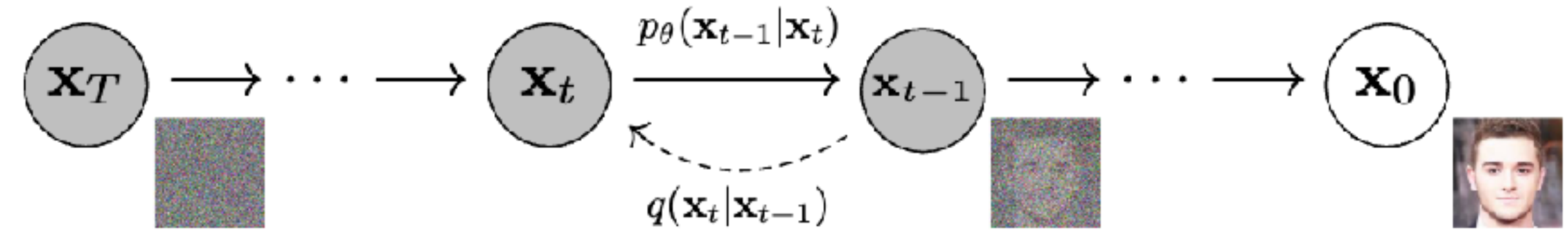
- Learn to predict image progressively
 - Start $P(x_T) = \mathcal{N}(\mathbf{0}, I)$
 - Denoise $P(x_{t-1} | x_t) = \mathcal{N}(\mu_\theta(x_t), \Sigma_\theta(x_t))$

- Reverse process

$$P(x_{0\dots T}) = P(x_T) \prod_{t=1}^T P(x_{t-1} | x_t)$$

- $P(x_0) = \int P(x_{0\dots T}) dx_{1\dots T}$

Diffusion Process



- Maximize Evidence lower bound (ELBO)

$$\log P(x_0) \geq E_q \left[\log \frac{P(x_{0\dots T})}{q(x_{1\dots T}|x_0)} \right]$$

- (Lot's of math later)
- Relatively simply training and sampling algorithms
 - $\epsilon(\mathbf{x}_t, t)$ is a noise-prediction network

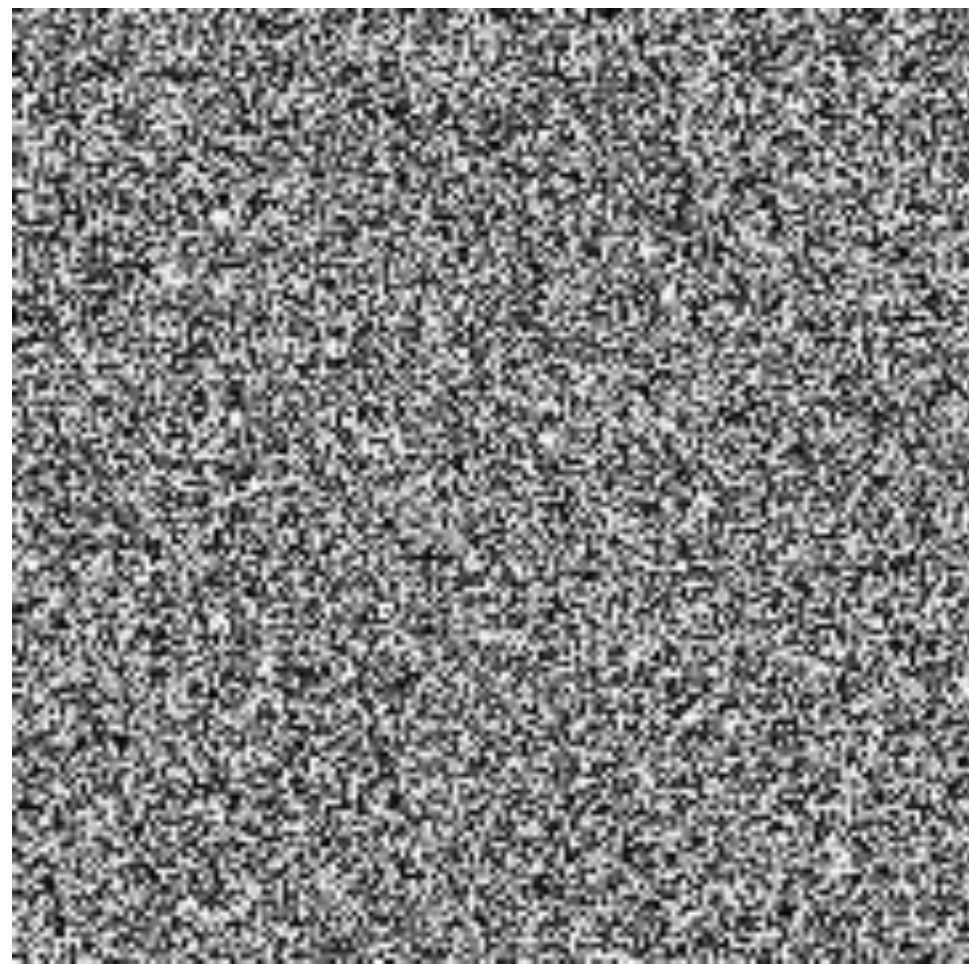
Algorithm 1 Training

- 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
$$\nabla_{\theta} \|\epsilon - \epsilon_{\theta}(\sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon, t)\|^2$$
 - 6: **until** converged
-

Algorithm 2 Sampling

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

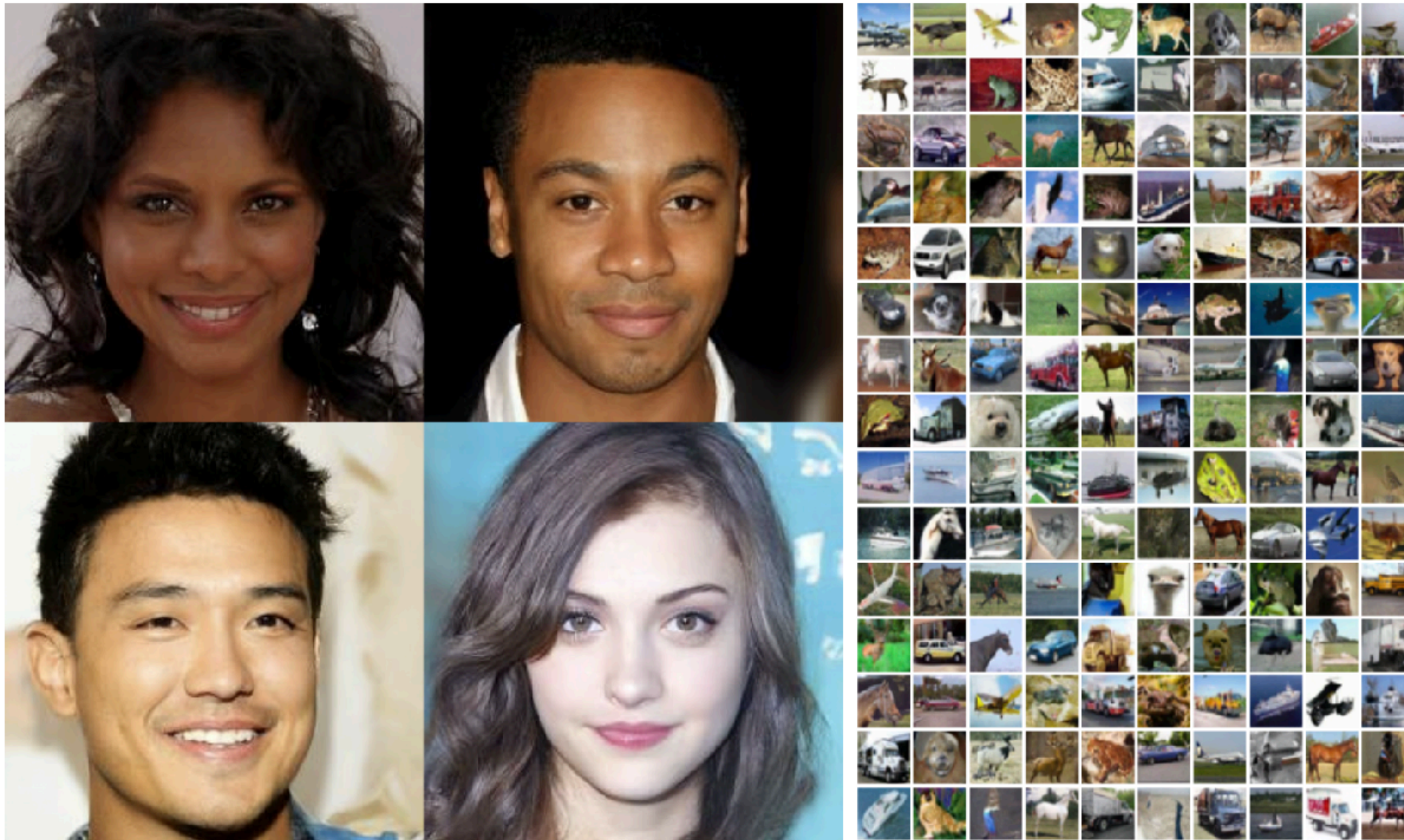
Diffusion Model



U-Net



Diffusion - Results

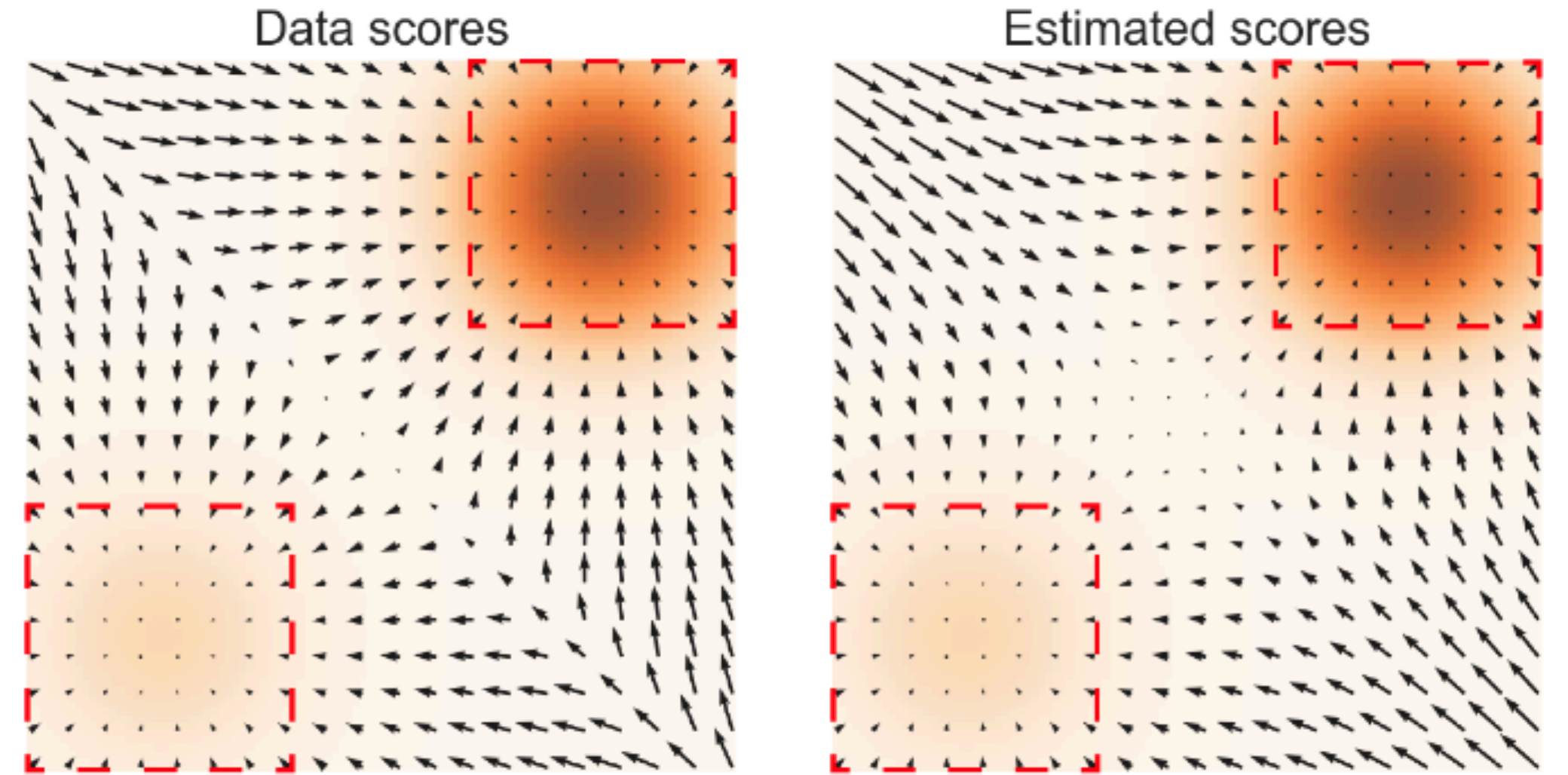


Score-based models

- $P(x) / \log P(x)$
 - is hard to learn
- $\nabla \log P(x)$ (score function)
 - is easier to learn/estimate

$$E_{x \sim P} \left[\left| s(x) - \nabla P(x) \right|^2 \right] =$$

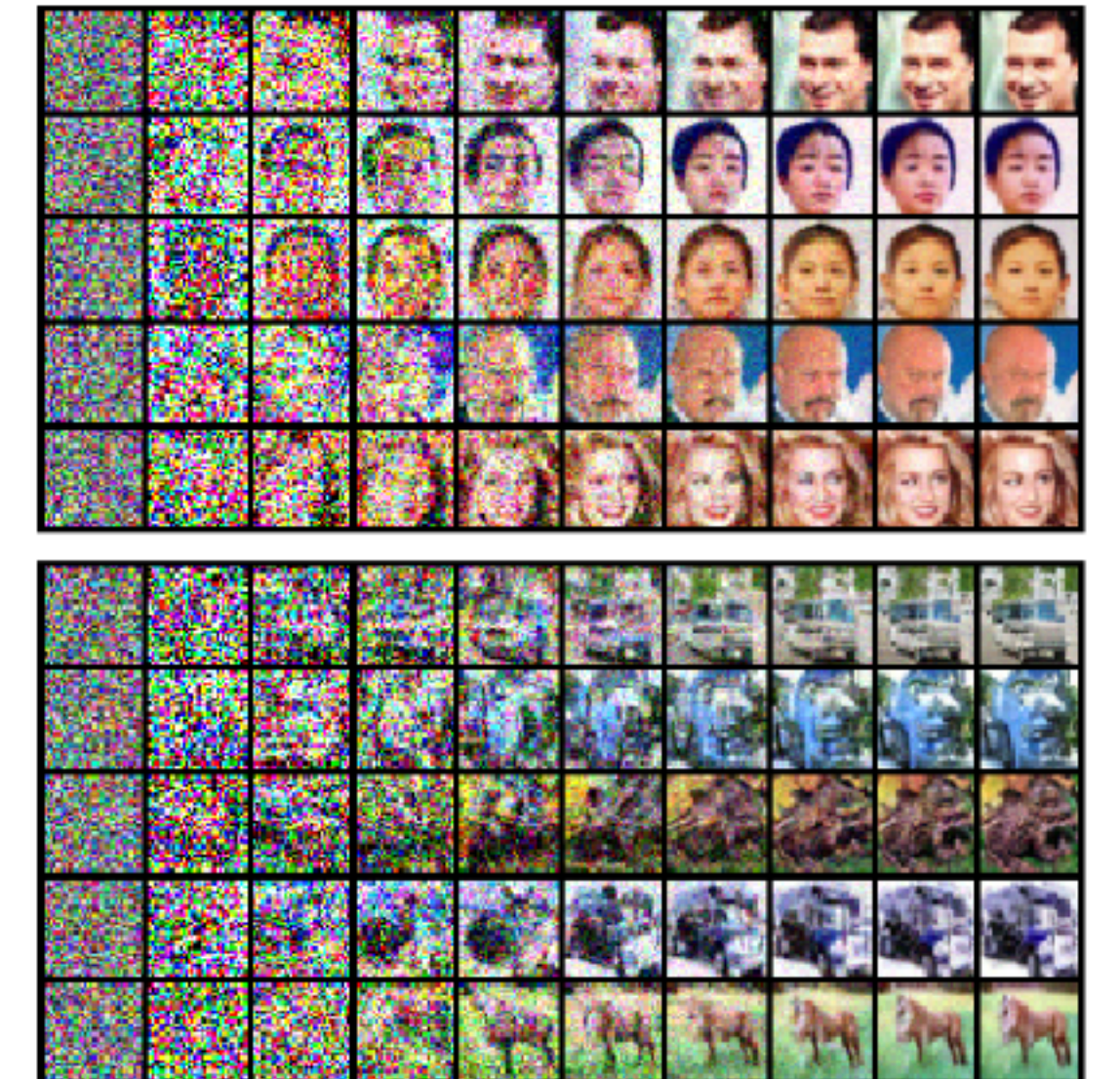
$$\bullet E_{x \sim P} \left[\text{tr}(\nabla_x s(x)) + \frac{1}{2} \left| s(x) \right|^2 \right] + \text{const}$$



Algorithm 1 Annealed Langevin dynamics.

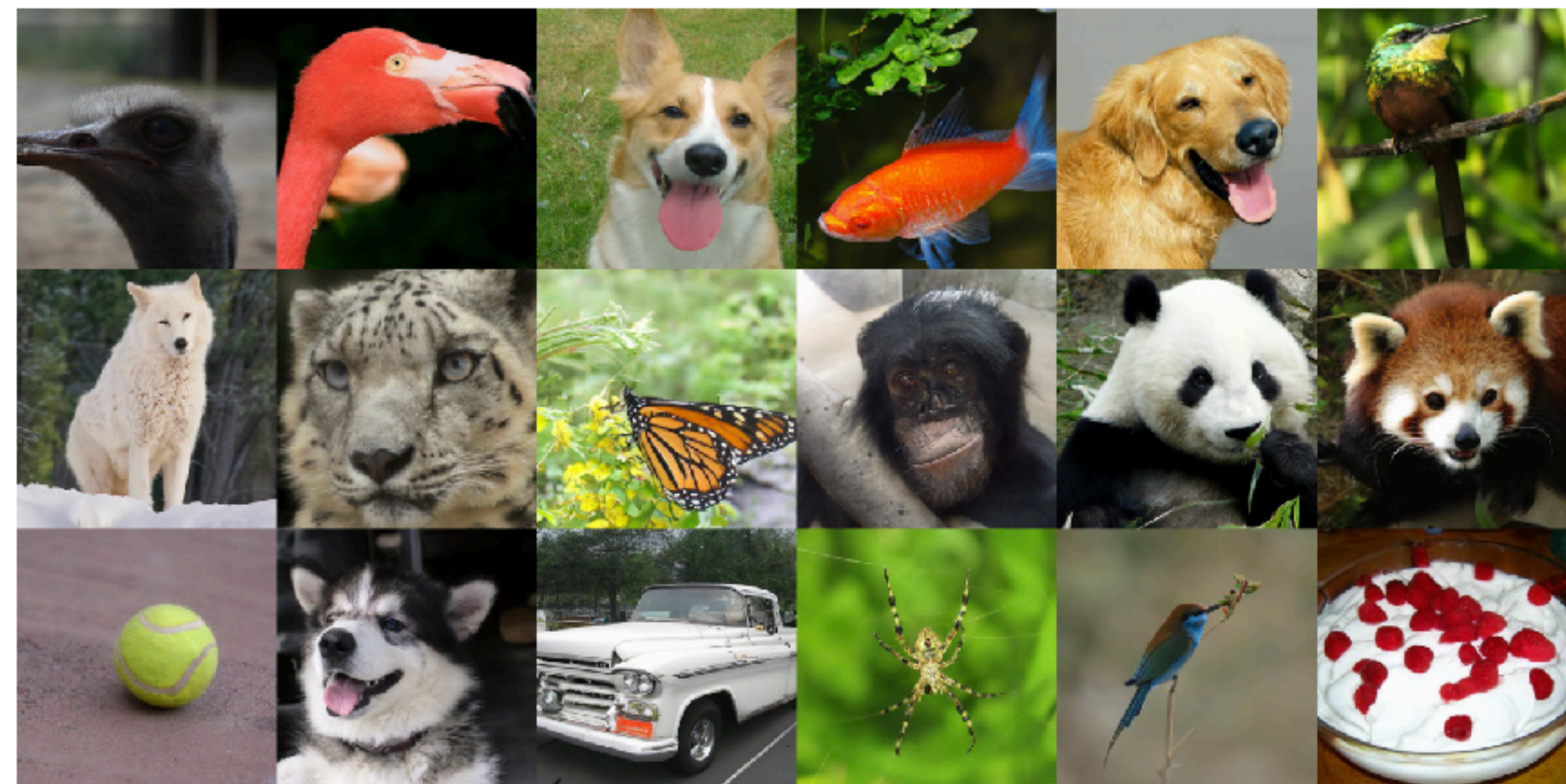
Require: $\{\sigma_i\}_{i=1}^L, \epsilon, T$.

- 1: Initialize $\tilde{\mathbf{x}}_0$
 - 2: **for** $i \leftarrow 1$ to L **do**
 - 3: $\alpha_i \leftarrow \epsilon \cdot \sigma_i^2 / \sigma_L^2$ $\triangleright \alpha_i$ is the step size.
 - 4: **for** $t \leftarrow 1$ to T **do**
 - 5: Draw $\mathbf{z}_t \sim \mathcal{N}(0, I)$
 - 6: $\tilde{\mathbf{x}}_t \leftarrow \tilde{\mathbf{x}}_{t-1} + \frac{\alpha_i}{2} \mathbf{s}_\theta(\tilde{\mathbf{x}}_{t-1}, \sigma_i) + \sqrt{\alpha_i} \mathbf{z}_t$
 - 7: **end for**
 - 8: $\tilde{\mathbf{x}}_0 \leftarrow \tilde{\mathbf{x}}_T$
 - 9: **end for**
- return** $\tilde{\mathbf{x}}_T$
-



Guided Diffusion

- Learn variance $\Sigma(x_t)$
- Better architecture
 - Deeper, more attention heads, attention on multiple blocks, ...
- Classifier guidance (conditioning)



Diffusion

- Very good image quality
- Not easily controllable
- Computationally quite expensive
 - Multiple sampling steps
 - Fairly high resolution inputs and outputs required (original image size)



References

- [1] WaveNet: A Generative Model for Raw Audio. Aaron van den Oord, et al. 2016
- [2] Long Video Generation with Time-Agnostic VQGAN and Time-Sensitive Transformer. Songwei Ge, et al. 2022
- [3] Denoising Diffusion Probabilistic Models. Jonathan Ho, et al. 2020.
- [4] Generative Modeling by Estimating Gradients of the Data Distribution. Yang Song, et al. 2019.
- [5] Diffusion Models Beat GANs on Image Synthesis. Prafulla Dhariwal, et al. 2021.