

Recap.

... in a nutshell

noise

data



x_T

...

x_t

x_{t-1}

...

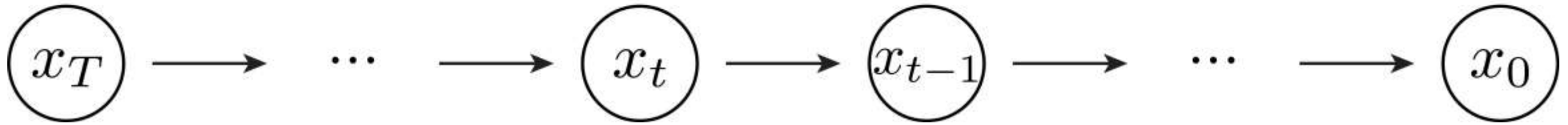
x_0



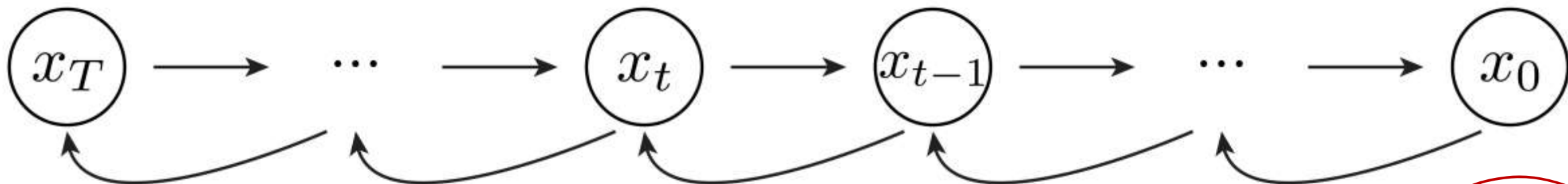
... in a nutshell

noise

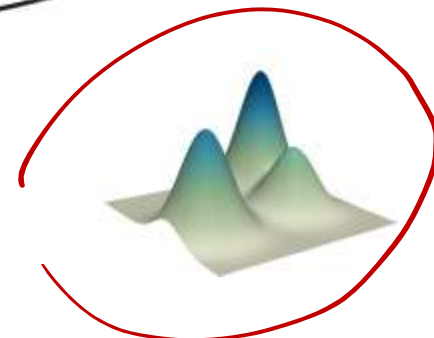
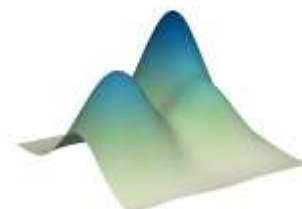
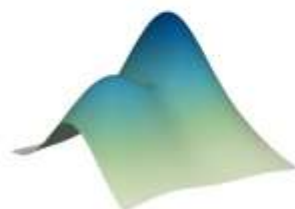
data



What is noise?

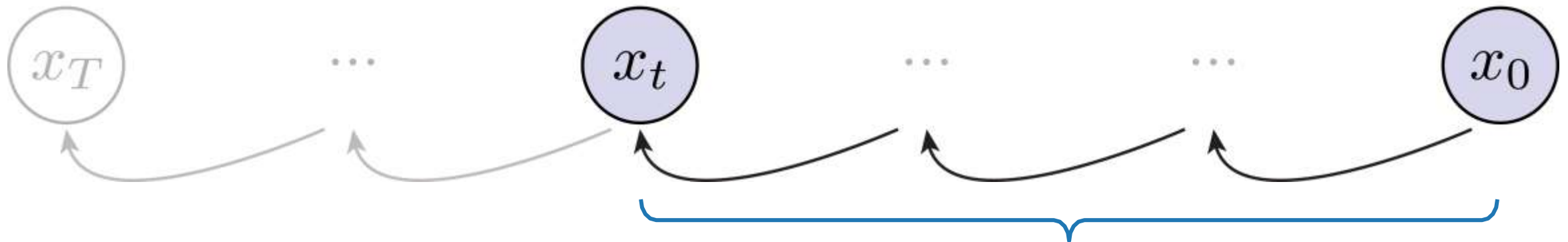


latent
distribution



data
distribution

Forward Process



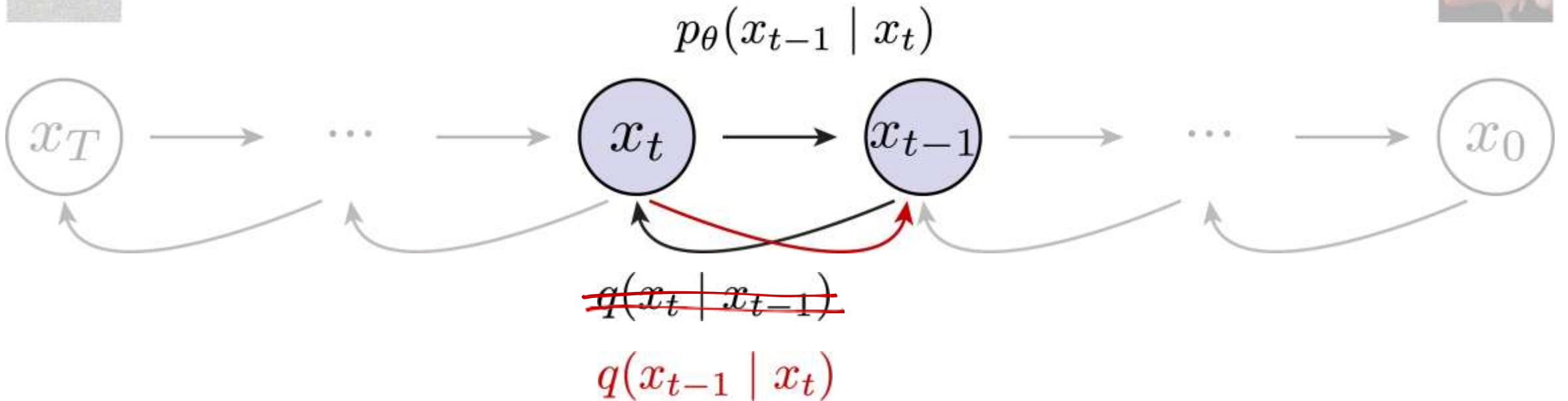
- sampling without simulation
- x_t from x_0 in closed form

$$q(x_t | x_0) = \mathcal{N}(x_t | \sqrt{\bar{\alpha}_t} x_0, (1 - \bar{\alpha}_t) \mathbf{I})$$

coefficients
given by β

$$\alpha_t := 1 - \beta_t$$
$$\bar{\alpha}_t := \prod_{s=1}^t \alpha_s$$

Reverse Process



- our target
- but unknown

Reverse Process

$$q(x_t|x_0) \quad \text{[noisy cat]} = \sqrt{\bar{\alpha}_t} \quad \text{[clear cat]} + \sqrt{1 - \bar{\alpha}_t} \quad \text{[noise]}$$

$$q(x_{t-1}|x_0) \quad \text{[less noisy cat]} = \sqrt{\bar{\alpha}_{t-1}} \quad \text{[clear cat]} + \sqrt{1 - \bar{\alpha}_{t-1}} \quad \text{[noise]}$$

$$q(x_t|x_{t-1}) \quad \text{[noisy cat]} = \sqrt{1 - \beta_t} \quad \text{[less noisy cat]} + \sqrt{\beta_t} \quad \text{[noise]}$$

$$q(x_{t-1}|x_t, x_0)$$

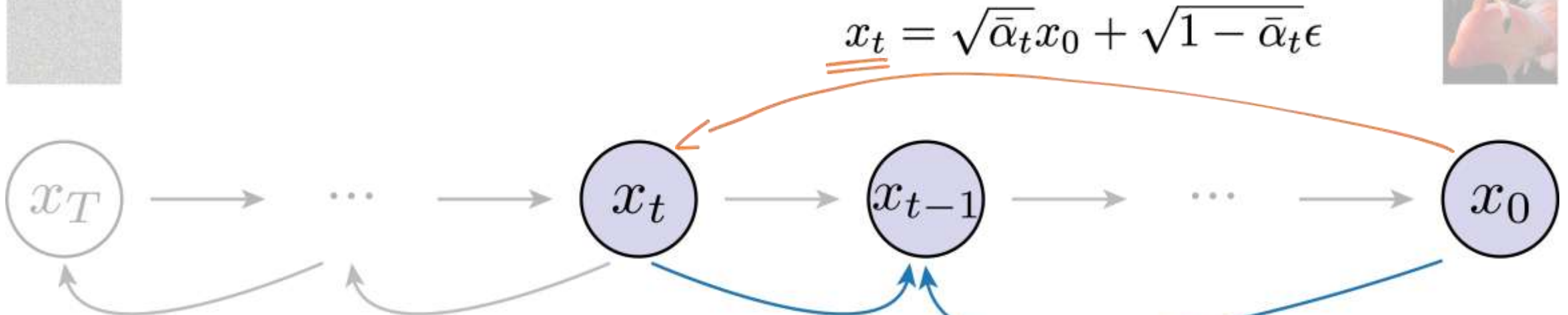
$$= \frac{q(x_{t-1}, x_t, x_0)}{q(x_t, x_0)} = \frac{q(x_t|x_{t-1})q(x_{t-1}|x_0)q(x_0)}{q(x_t|x_0)q(x_0)} = \frac{q(x_t|x_{t-1})q(x_{t-1}|x_0)}{q(x_t|x_0)}$$

- known
- Gaussian
- known
- Gaussian
- known
- Gaussian

Reverse Process

tl; dr:

- outcome of the dependency graph
- some linear combinations



$$q(x_{t-1} | x_t, x_0)$$

$$= \mathcal{N}(x_{t-1} | \tilde{\mu}(x_t, x_0), \tilde{\beta}_t \mathbf{I})$$

$$\tilde{\mu}_t(x_t, x_0) := \frac{\sqrt{\alpha_t} \beta_t}{1 - \alpha_t} x_0 + \frac{\sqrt{\alpha_t} (1 - \alpha_{t-1})}{1 - \bar{\alpha}_t} x_t$$

mean

$$= \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon \right) \text{ "noise"}$$

$$\tilde{\beta}_t := \frac{1 - \alpha_{t-1}}{1 - \bar{\alpha}_t} \beta_t$$

var



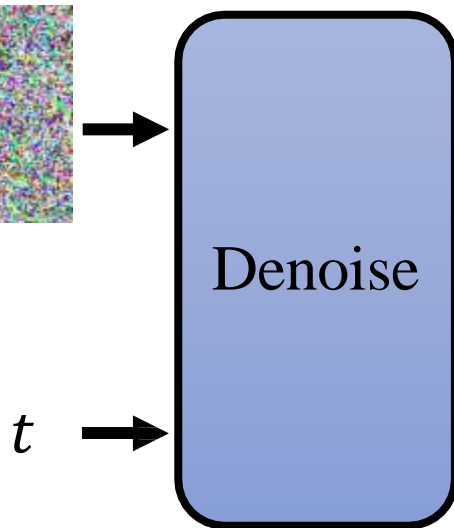
x_0



x_t



Sample x_t



$$\frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \begin{matrix} \boxed{\varepsilon} \\ \vdots \end{matrix} \right)$$

To predict

Algorithm 2 Sampling

$$x_t = \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \varepsilon$$

$$x_t - \sqrt{1 - \bar{\alpha}_t} \varepsilon = \sqrt{\bar{\alpha}_t} x_0$$

$$\frac{x_t - \sqrt{1 - \bar{\alpha}_t} \varepsilon}{\sqrt{\bar{\alpha}_t}} = x_0$$

- 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: $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\varepsilon}_{\theta}(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$
 - 5: **end for**
 - 6: **return** \mathbf{x}_0
-

Diffusion Models

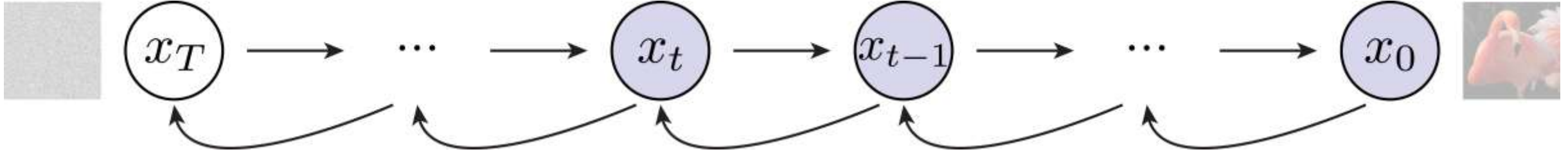
- Forward process
 - add noise to data

- Reverse process
 - learn to denoise

- Training objective
 - from Hierarchical VAE to L2 loss

- Noise Conditional Network
 - represent a distribution

Training Objective



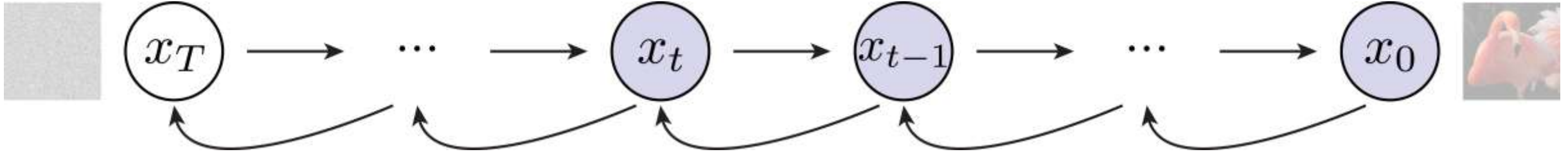
$$\mathcal{L}_{\text{VLB}} := \mathcal{L}_T + \mathcal{L}_{T-1} + \dots + \mathcal{L}_0$$

$$\mathcal{L}_T := D_{\text{KL}}\left(q(x_T | x_0) \parallel p_{\theta}(x_T)\right)$$

$$\mathcal{L}_{t-1} := D_{\text{KL}}\left(q(x_{t-1} | x_t, x_0) \parallel p_{\theta}(x_{t-1} | x_t)\right)$$

$$\mathcal{L}_0 := -\log p_{\theta}(x_0 | x_1)$$

Training Objective



- variational lower bound
- like ELBO

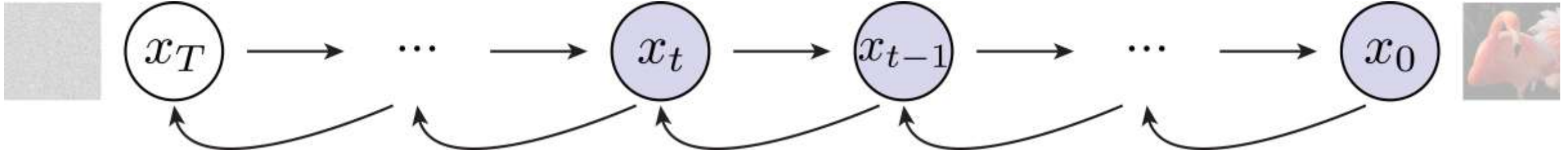
$$\mathcal{L}_{\text{VLB}} := \mathcal{L}_T + \mathcal{L}_{T-1} + \dots + \mathcal{L}_0$$

$$\mathcal{L}_T := D_{\text{KL}}\left(q(x_T | x_0) \parallel p_{\theta}(x_T)\right)$$

$$\mathcal{L}_{t-1} := D_{\text{KL}}\left(q(x_{t-1} | x_t, x_0) \parallel p_{\theta}(x_{t-1} | x_t)\right)$$

$$\mathcal{L}_0 := -\log p_{\theta}(x_0 | x_1)$$

Training Objective – to VAE



- variational lower bound
- like ELBO

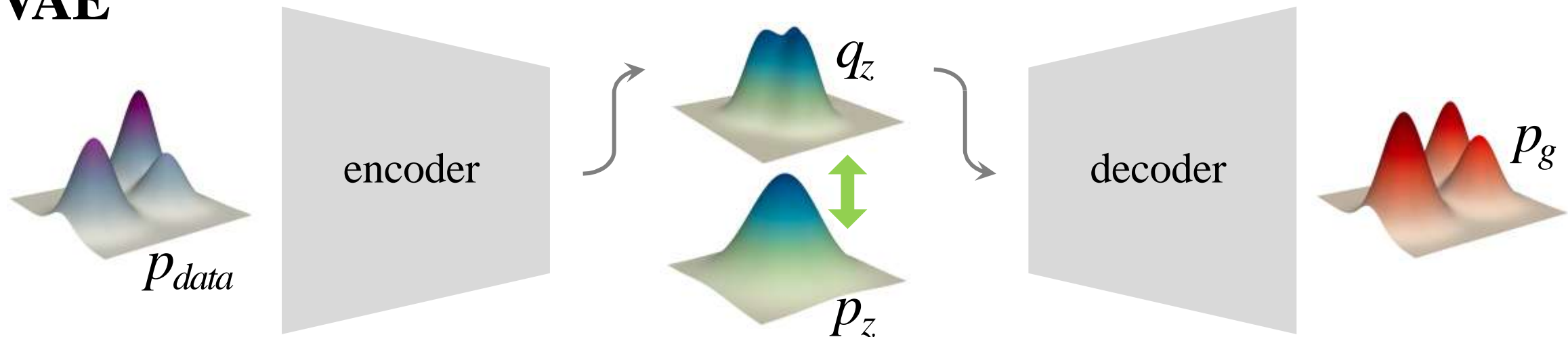
$$\mathcal{L}_{\text{VLB}} := \mathcal{L}_T + \mathcal{L}_{T-1} + \dots + \mathcal{L}_0$$

$$\mathcal{L}_T := D_{\text{KL}}\left(q(\overset{z}{\cancel{x_T}} \mid x_0) \parallel p_{\theta}(\overset{z}{\cancel{x_T}})\right) \quad \text{it's ELBO if one step}$$

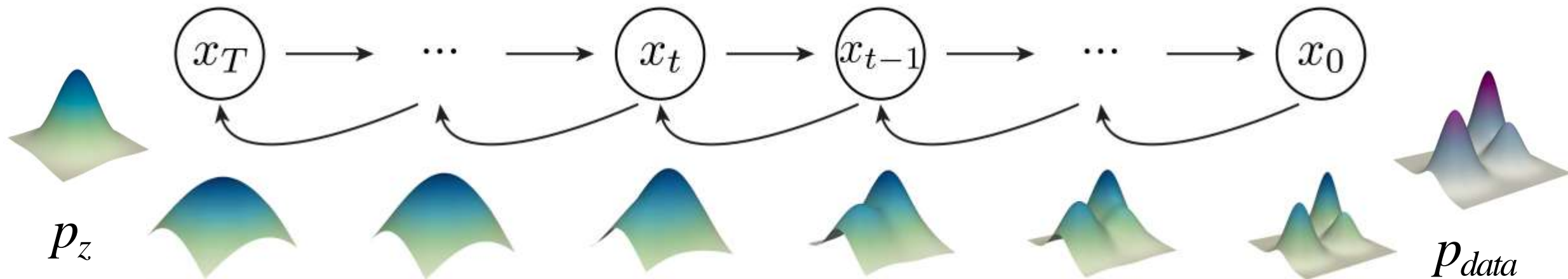
~~$$\mathcal{L}_{t-1} := D_{\text{KL}}\left(q(x_{t-1} \mid x_t, x_0) \parallel p_{\theta}(x_{t-1} \mid x_t)\right)$$~~

$$\mathcal{L}_0 := -\log p_{\theta}(x_0 \mid \overset{z}{\cancel{x_1}})$$

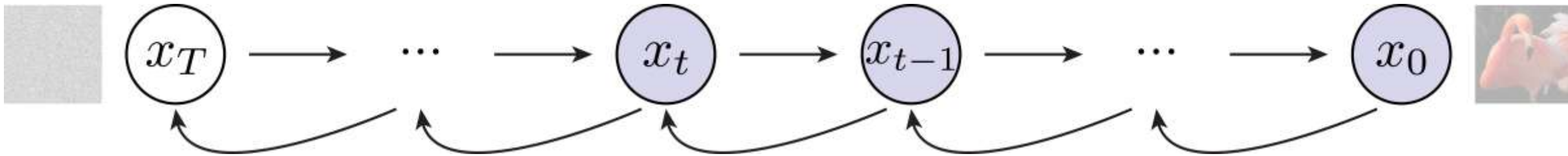
VAE



DM/Flow



Training Objective



$$\mathcal{L}_{\text{VLB}} := \mathcal{L}_T + \mathcal{L}_{T-1} + \dots + \mathcal{L}_0$$

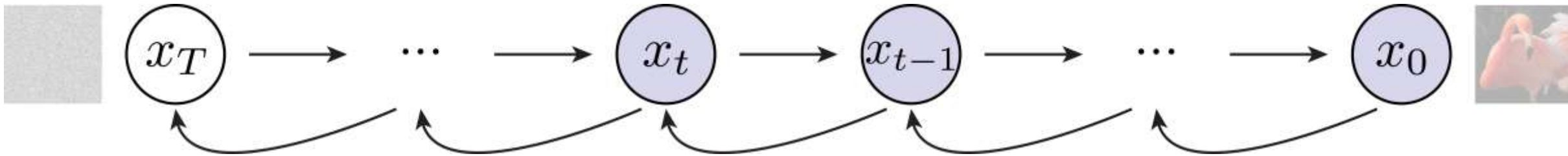
no parameter, unlike VAE's q_ϕ

$$\mathcal{L}_T := D_{\text{KL}} \left(\underbrace{q(x_T | x_0)}_{\text{constant}} \parallel \underbrace{p_\theta(x_T)}_{\text{Gaussian}} \right)$$

$$\mathcal{L}_{t-1} := D_{\text{KL}} \left(q(x_{t-1} | x_t, x_0) \parallel p_\theta(x_{t-1} | x_t) \right)$$

$$\mathcal{L}_0 := -\log p_\theta(x_0 | x_1)$$

Training Objective



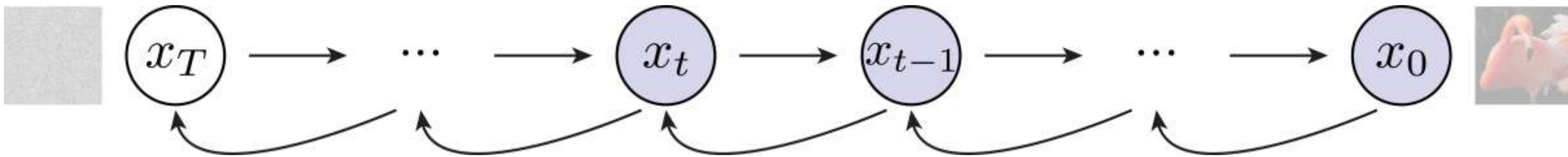
$$\mathcal{L}_{\text{VLB}} := \mathcal{L}_T + \mathcal{L}_{T-1} + \dots + \mathcal{L}_0$$

$$\mathcal{L}_T := D_{\text{KL}}\left(q(x_T | x_0) \parallel p_{\theta}(x_T)\right)$$

$$\mathcal{L}_{t-1} := D_{\text{KL}}\left(q(x_{t-1} | x_t, x_0) \parallel p_{\theta}(x_{t-1} | x_t)\right)$$

reconstruction loss
like VAE $\mathcal{L}_0 := -\log p_{\theta}(x_0 | x_1)$

Training Objective



$$\mathcal{L}_{\text{VLB}} := \mathcal{L}_T + \mathcal{L}_{T-1} + \dots + \mathcal{L}_0$$

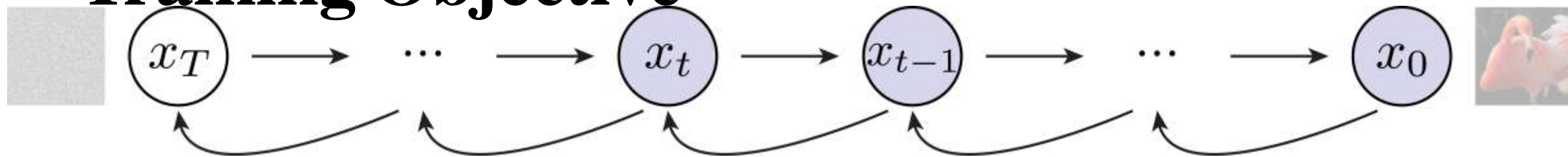
$$\mathcal{L}_T := D_{\text{KL}}\left(q(x_T | x_0) \parallel p_{\theta}(x_T)\right)$$

L2 loss on
noise

$$\mathcal{L}_{t-1} := D_{\text{KL}}\left(q(x_{t-1} | x_t, x_0) \parallel p_{\theta}(x_{t-1} | x_t)\right)$$

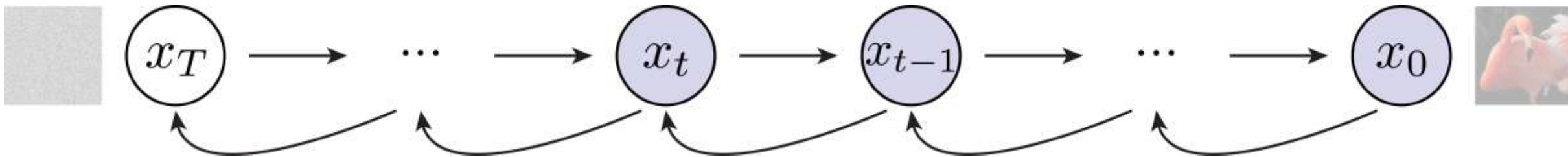
$$\mathcal{L}_0 := -\log p_{\theta}(x_0 | x_1)$$

Training Objective



$$\mathcal{L} = \mathbb{E}_{x_0, t, \epsilon} \left[w_t \|\epsilon - \epsilon_\theta(x_t, t)\|^2 \right]$$

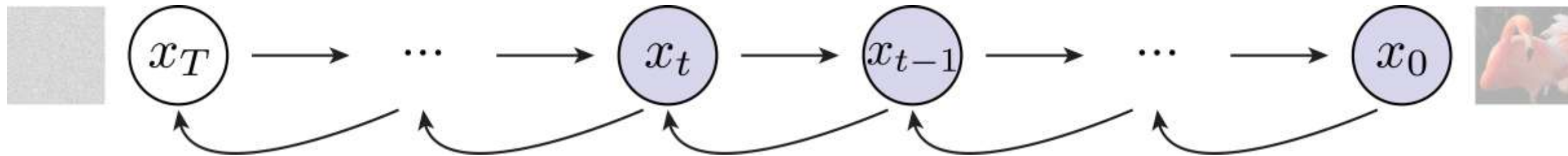
Training Objective



$$\mathcal{L} = \mathbb{E}_{x_0, t, \epsilon} \left[w_t \|\epsilon - \epsilon_\theta(x_t, t)\|^2 \right]$$

over p_{data} over $[1, T]$ over $\mathcal{N}(0, \mathbf{I})$

Training Objective



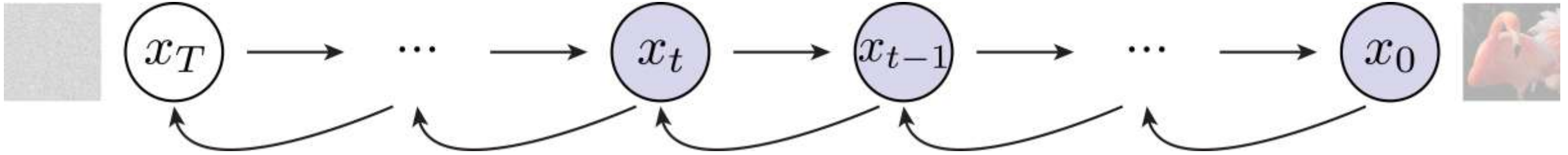
$$\mathcal{L} = \mathbb{E}_{x_0, t, \epsilon} \left[w_t \|\epsilon - \epsilon_\theta(x_t, t)\|^2 \right]$$

set as 1 (**critical**)

Objective	IS	FID
L , learned diagonal Σ	–	–
L , fixed isotropic Σ	7.67 ± 0.13	13.51
$\ \tilde{\epsilon} - \epsilon_\theta\ ^2$ (L_{simple})	9.46 ± 0.11	3.17

[Ho et al. 2020]; see more in [Salimans & Ho, 2022]

Training Objective



$$\mathcal{L} = \mathbb{E}_{x_0, t, \epsilon} \left[\left\| w_t \left\| \epsilon - \underbrace{\epsilon_\theta}_{\text{network to predict noise}}(x_t, \underbrace{t}_{\text{conditioned on noise level (critical)}}) \right\|^2 \right]$$

Diffusion Models

- Forward process
 - add noise to data

- Reverse process
 - learn to denoise

- Training objective
 - from Hierarchical VAE to L2 loss

- Noise Conditional Network
 - represent a distribution

Noise Conditional Network

- Diffusion models decompose a distribution into **many** simpler ones.
- We need the same # networks to fit **all** of them.
- We can **combine** all into one “powerful” network.
- This network is conditioned on noise level t .

- **Noise Conditional Network** [Song & Ermon 2019]: things made work

Noise Conditional Network

How to represent $p_{\theta}(x_{t-1} | x_t)$

- network input: x_t
- network output: μ and σ of a distribution

- parametrize μ by: $\epsilon_{\theta}(x_t, t)$

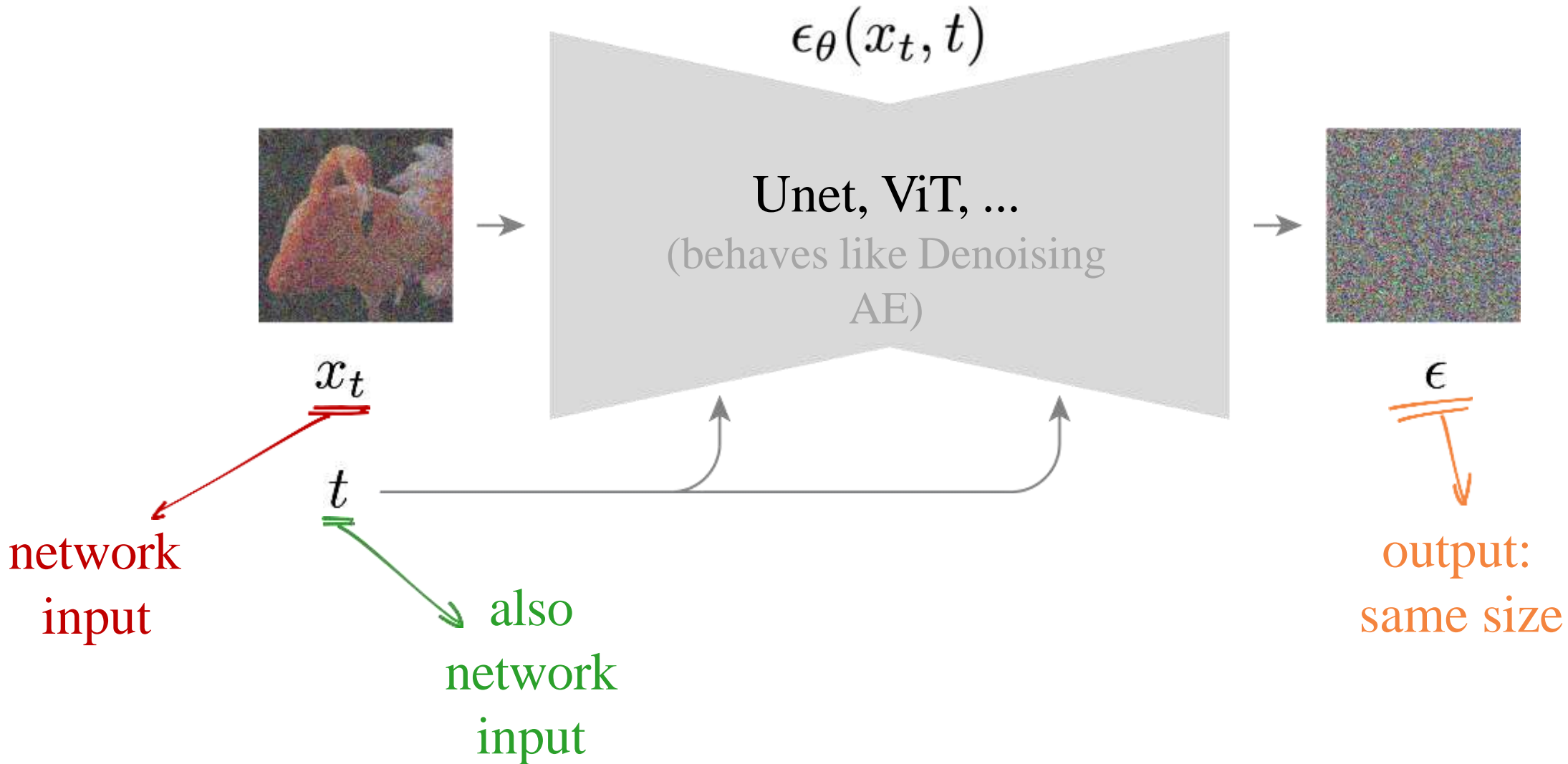
noisy image:

- condition
- network input

noise level:

- condition
- network input

Noise Conditional Network



Diffusion algorithm annotated:

Algorithm 1 Training

- 1: **repeat**
 - 2: $\mathbf{x}_0 \sim q(\mathbf{x}_0)$
 - 3: $t \sim \text{Uniform}(\{1, \dots, T\})$
 - 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
-

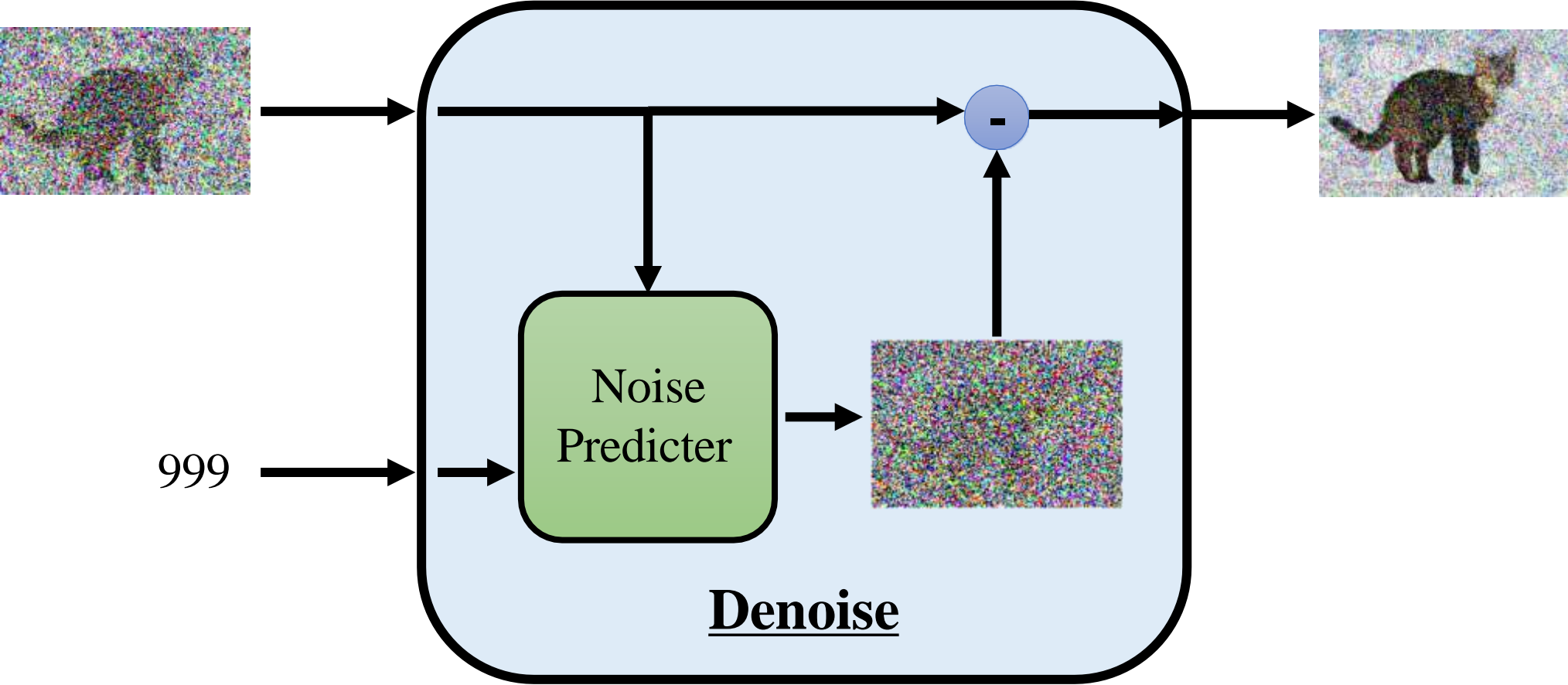
Algorithm 2 Sampling

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

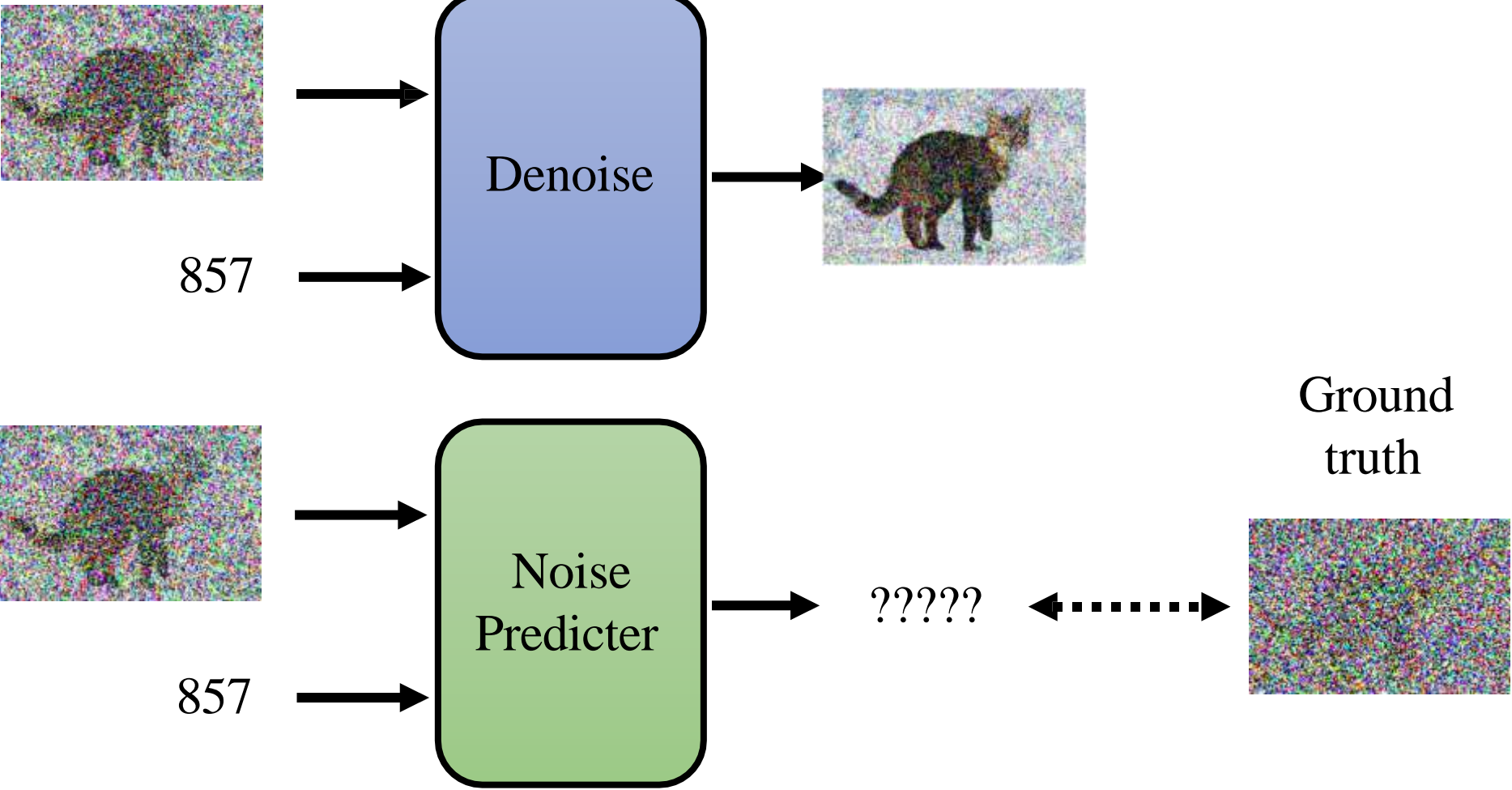
estimated μ

sampling from
estimated distribution

Noise Conditional Network



Noise Conditional Network



Diffusion algorithm annotated:

Algorithm 1 Training

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

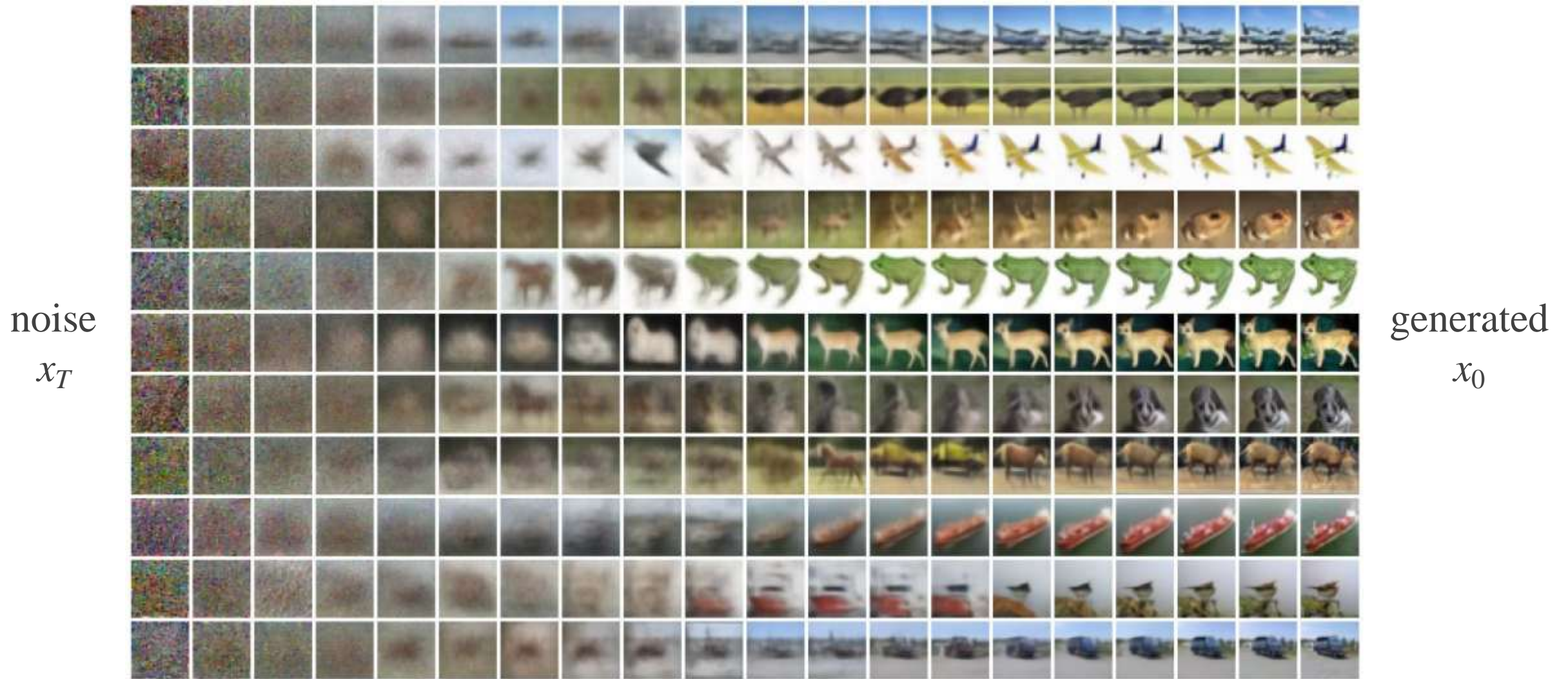
Algorithm 2 Sampling

```
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:  $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$   
5: end for  
6: return  $\mathbf{x}_0$ 
```

tl; dr: noising and denoising

- Turns out to be extremely simple
- Being “simple and effective” moves the needle

Example: Unconditional Generation on CIFAR-10

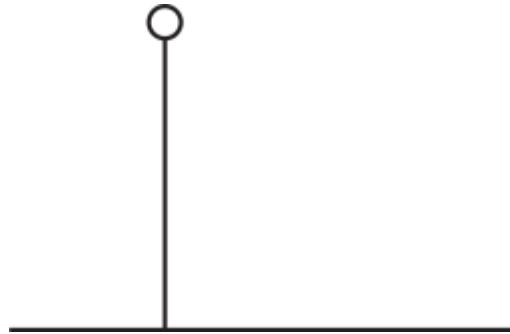


Example: shared intermediate latents

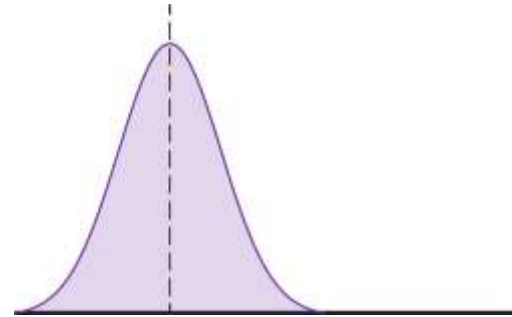


Recap.

- Adding Gaussian noise \Leftrightarrow sampling $x \sim \mathcal{N}(x | x_0, \sigma)$



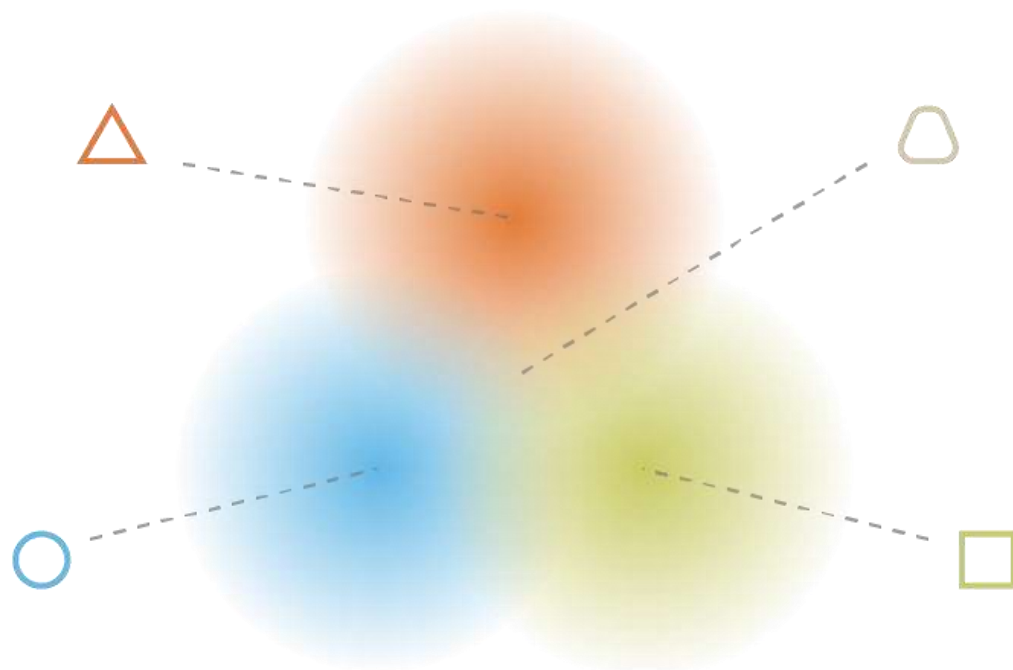
$$p(x) = \delta(x - x_0)$$



$$p(x) = \mathcal{N}(x | x_0, \sigma)$$



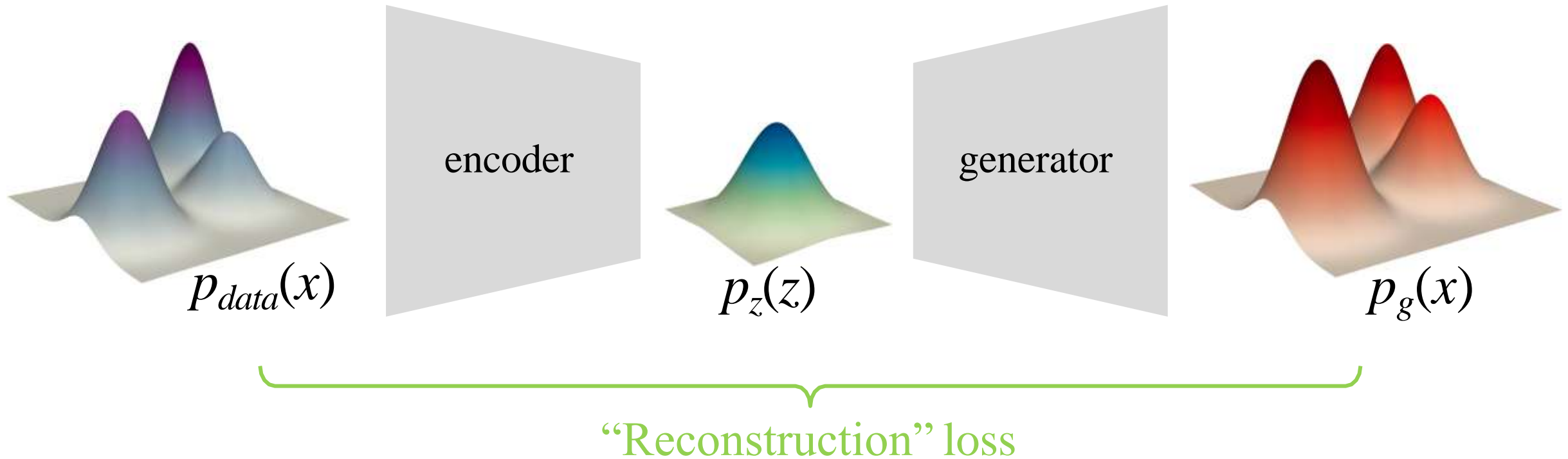
Problems



Recap: Variational Autoencoder (VAE)

Autoencoding distributions:

“Encoding” data distribution p_{data} into latent distribution p_z



Summary

Forward process

- add noise to data

Reverse process

- learn to denoise

Training objective

- from Hierarchical VAE to L2 loss

Noise Conditional Network

- represent distributions

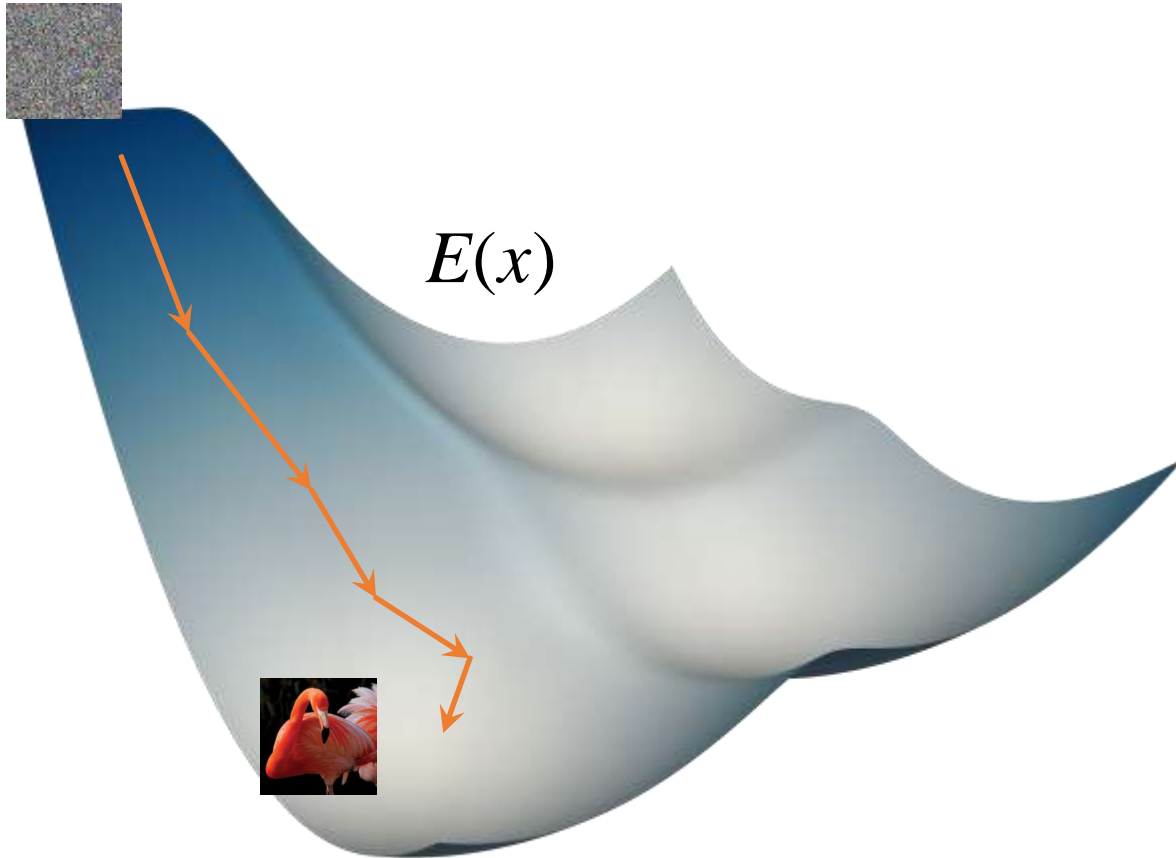
Energy-based Models and Score Matching

Diffusion and Score Matching

- Diffusion Models are closely related to **Score Matching**.
- Score Matching is one solution to **Energy-based Models**.
- Energy-based Models:
 - can be probabilistic or non-probabilistic
 - can be generative or discriminative
- Many useful concepts in diffusion co-evolved w/ score matching
 - Annealed importance sampling [Neal 1998]
 - Denoising score matching [Vincent 2011]
 - Noise Conditional Score Network [Song & Ermon 2019]

Energy-based Models

- Define a scalar function, called “energy”.
- At inference time, find x that minimizes energy

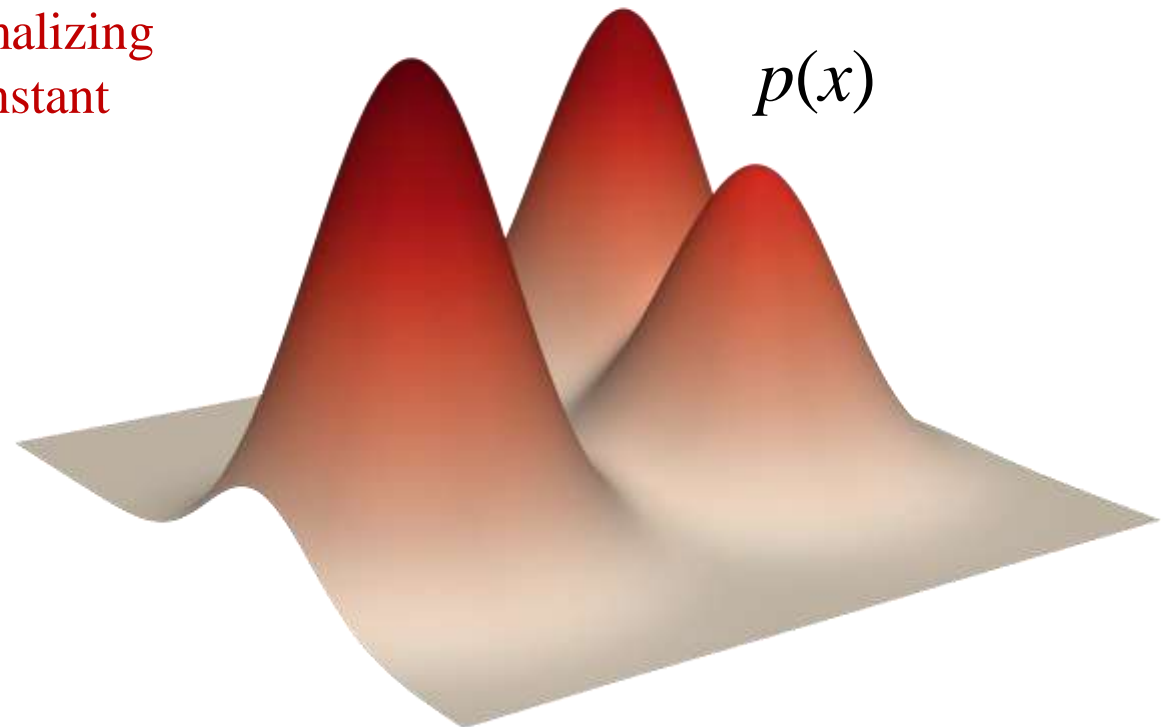
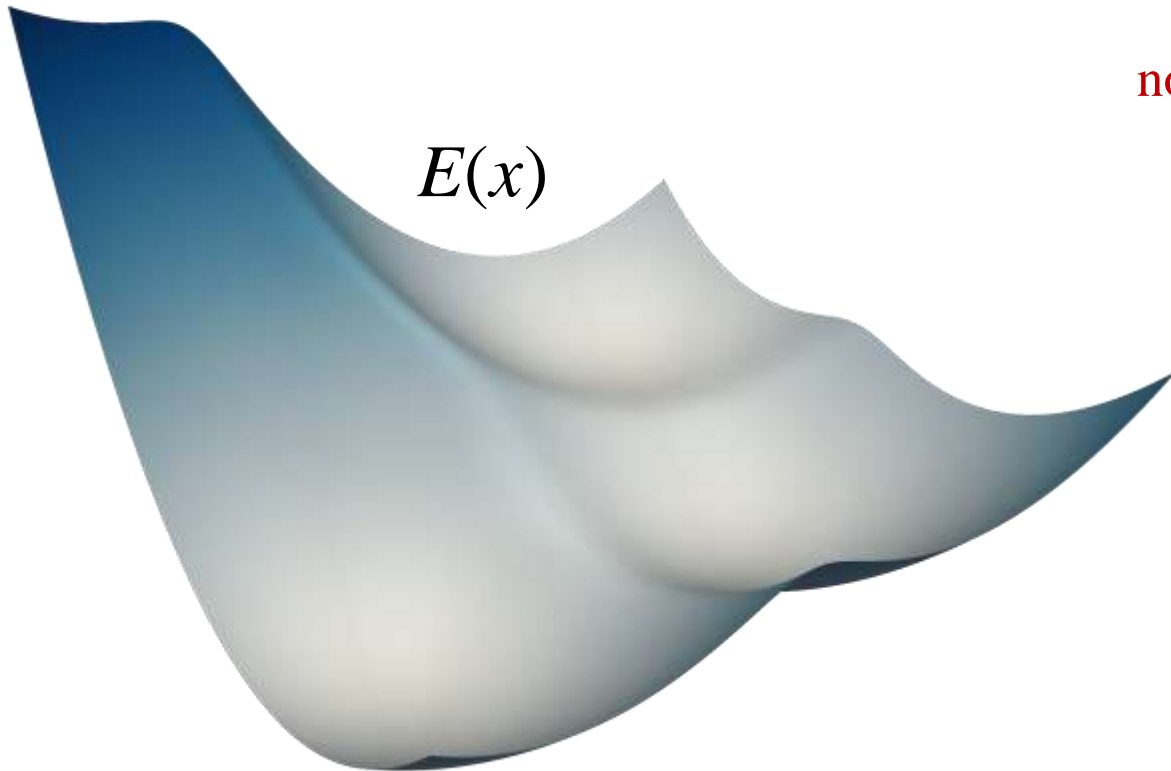


Energy-based Models

- We can use an energy to model a probability distribution

$$p(x) = \frac{\exp(-E(x))}{Z}$$

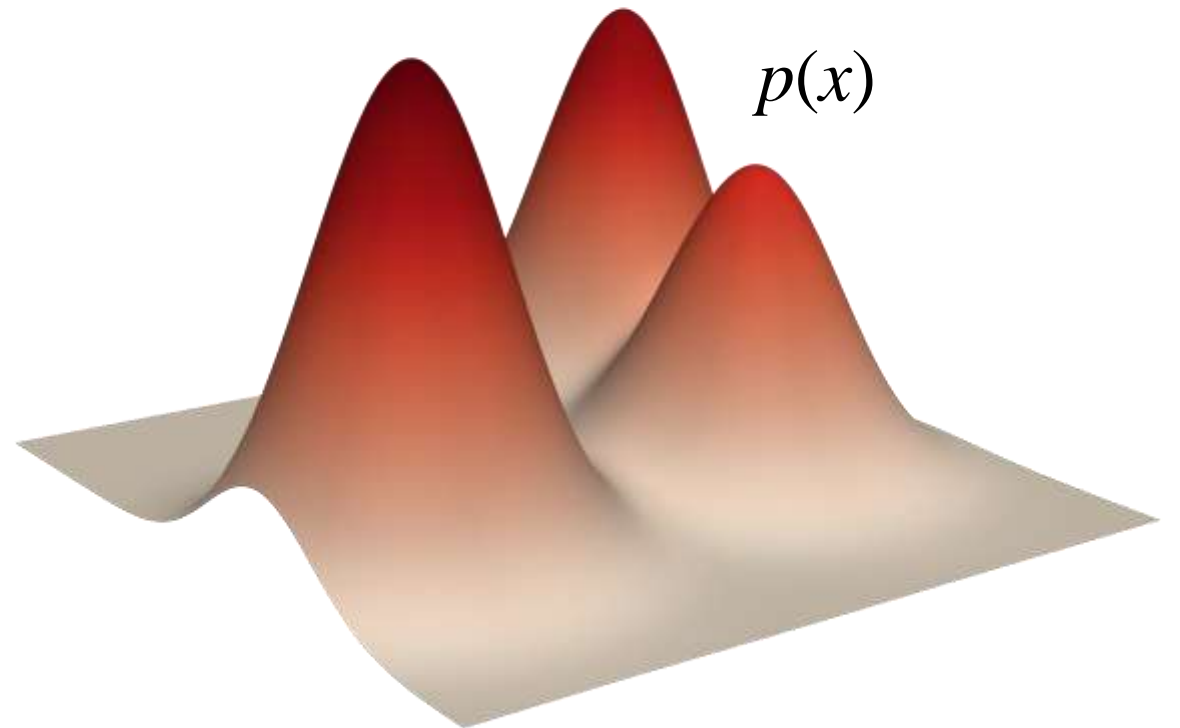
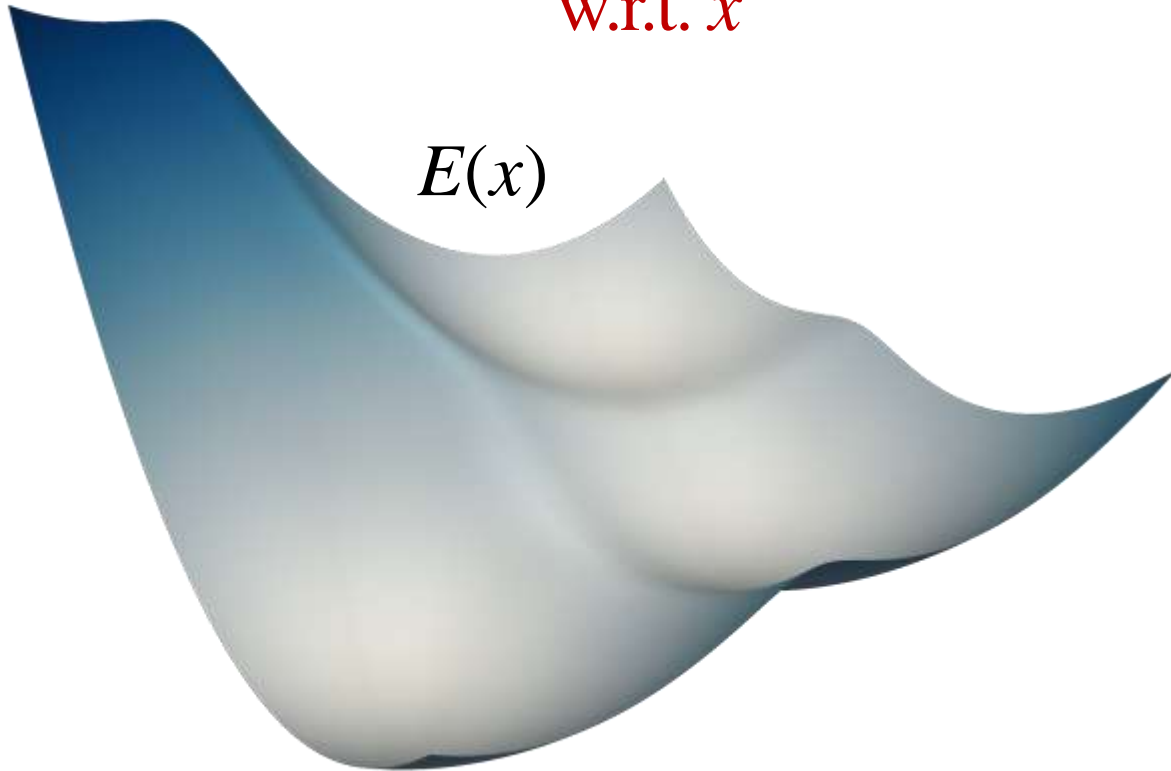
normalizing
constant



Energy-based Models

- “Score function”: gradient of log-probability.

$$\underbrace{\nabla_x}_{\text{w.r.t. } x} \log p(x) = -\nabla_x E(x) - \underbrace{\nabla_x \log Z}_{=0}$$

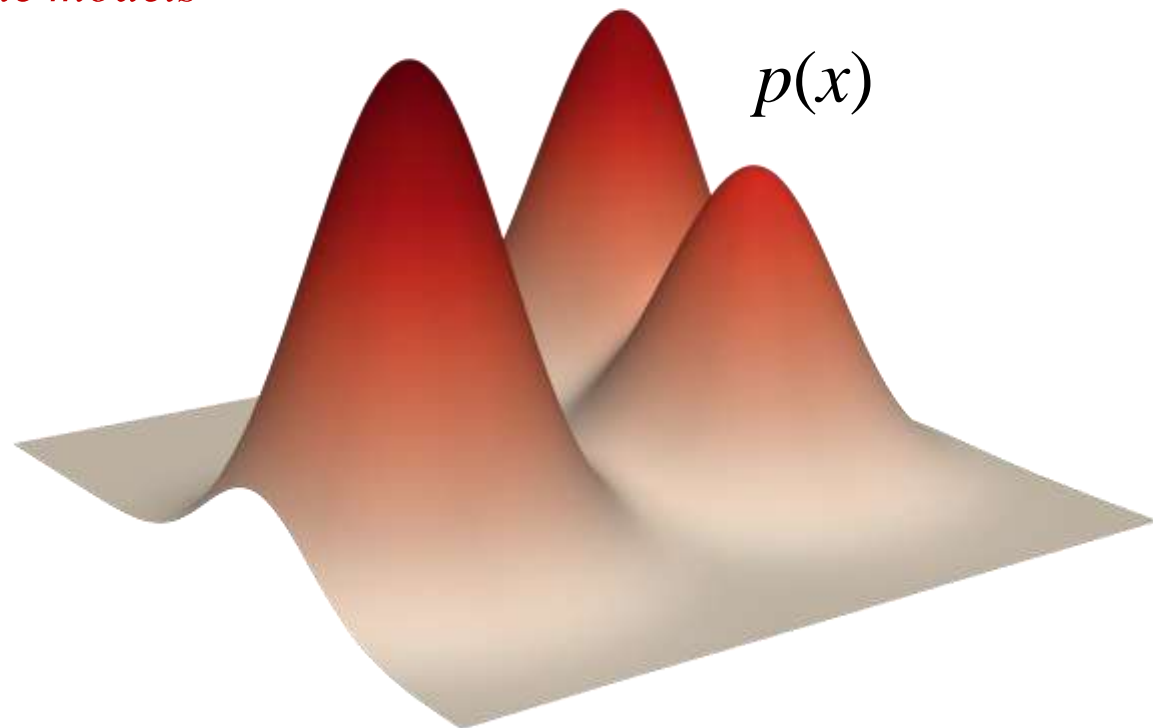
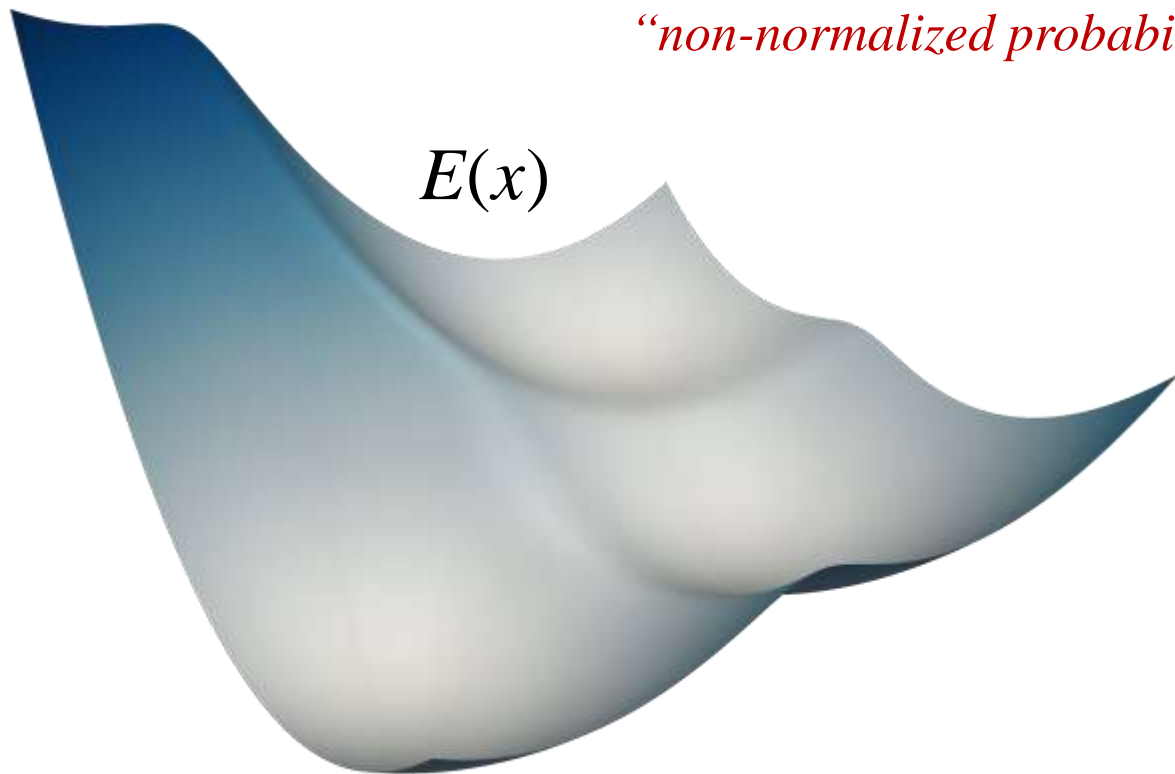


Energy-based Models

- “**Score function**”: gradient of log-probability

$$\nabla_x \log p(x) = -\nabla_x E(x)$$

“non-normalized probabilistic models”

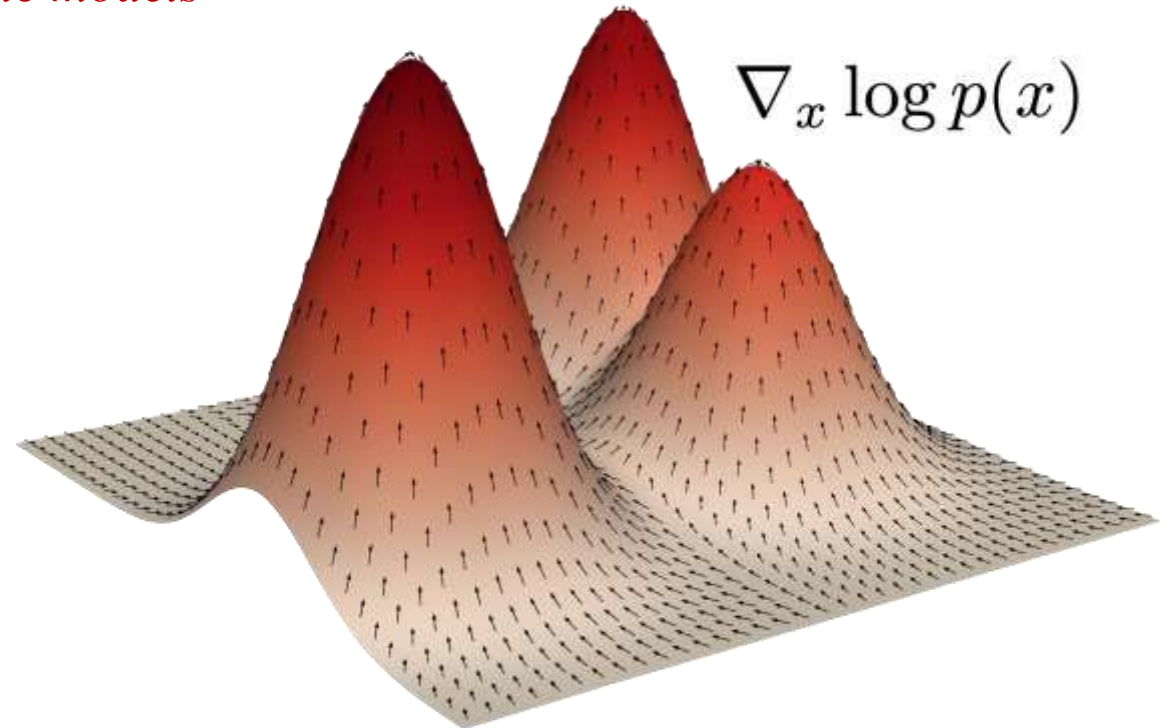
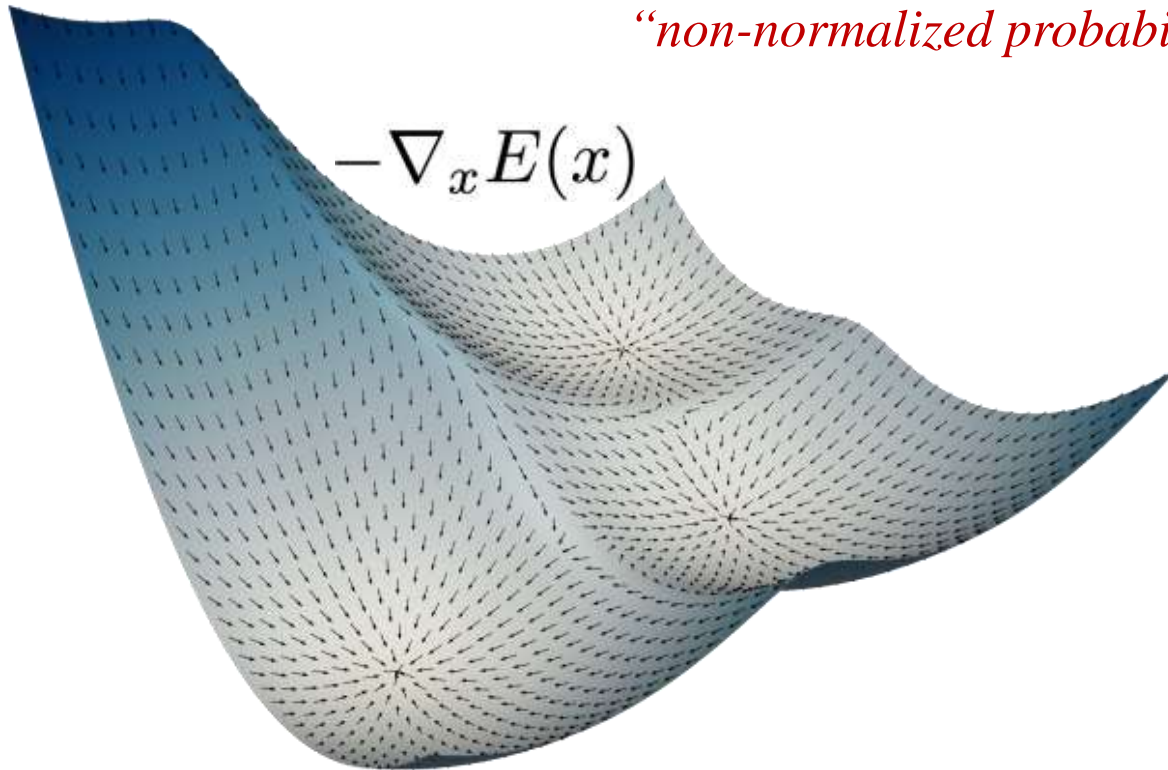


Energy-based Models

- “**Score function**”: gradient of log-probability

$$\nabla_x \log p(x) = -\nabla_x E(x)$$

“non-normalized probabilistic models”

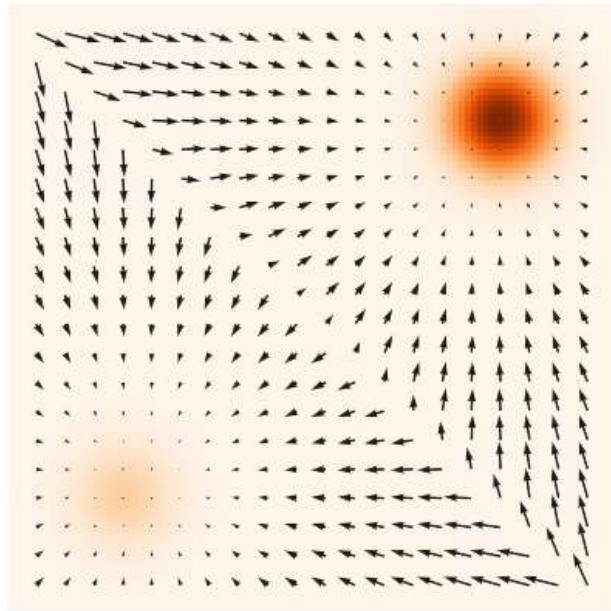


*only visualize

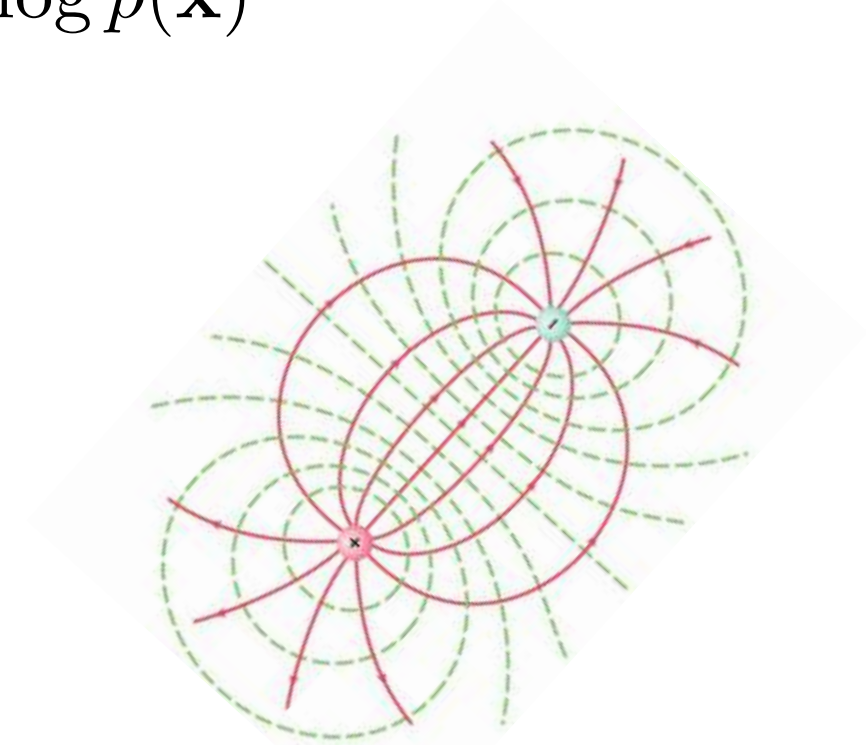
How to represent probability distributions?

- When the pdf is differentiable, we can compute the gradient of a probability density.

Score function $\nabla_{\mathbf{x}} \log p(\mathbf{x})$



(pdf and score)

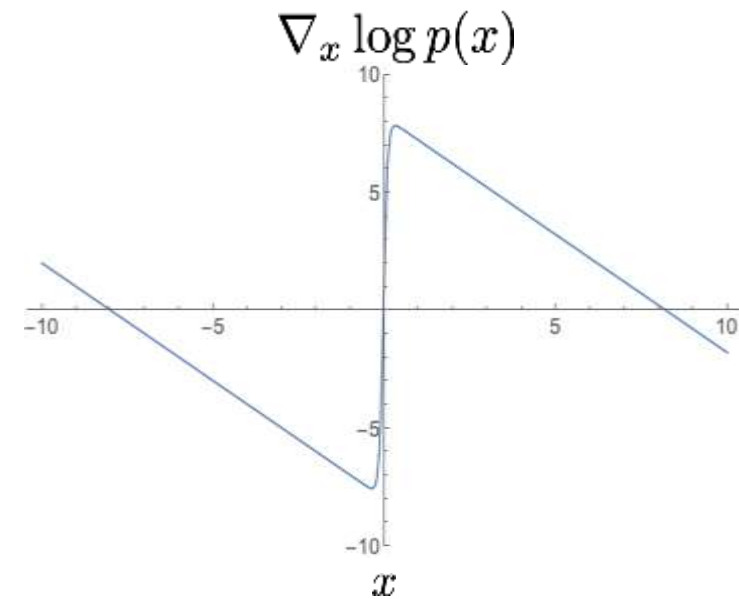
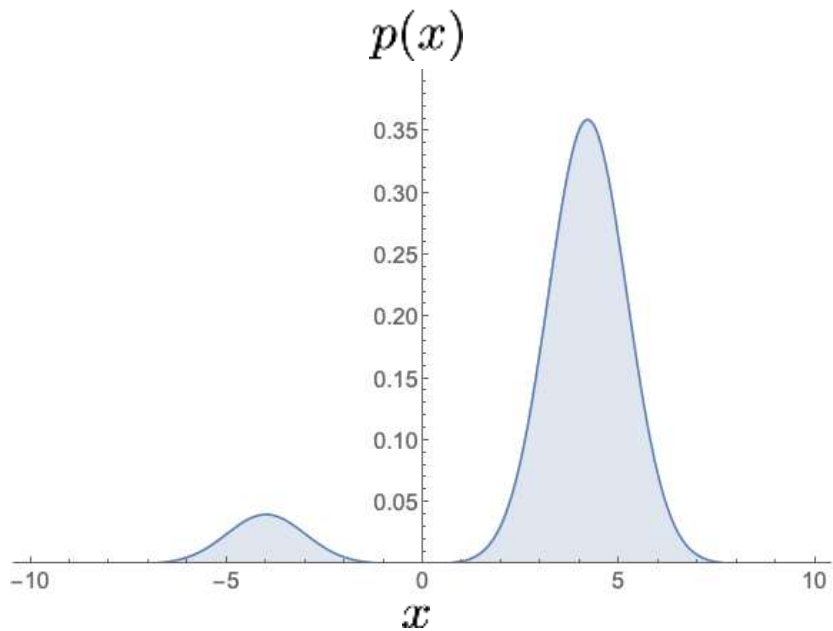


(Electrical potentials and fields)

How to represent probability distributions?

- When the pdf is differentiable, we can compute the gradient of a probability density.

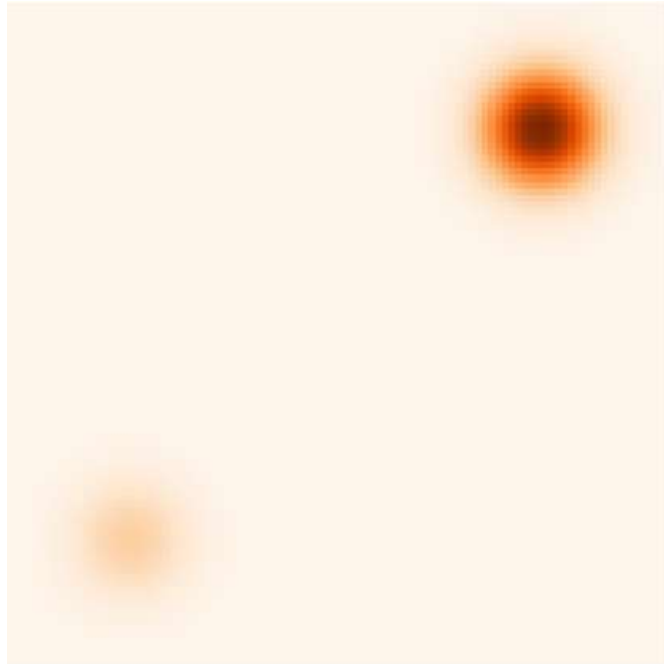
Score function $\nabla_{\mathbf{x}} \log p(\mathbf{x})$



Score estimation by training score-based models

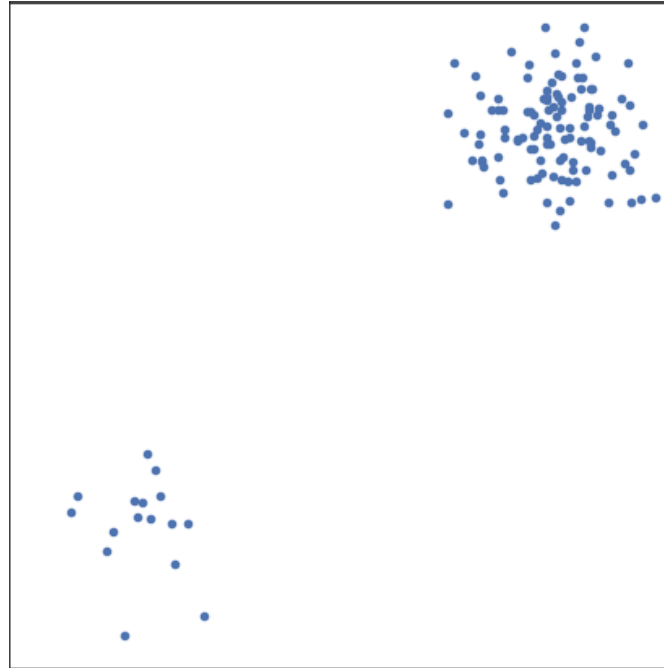
Probability density

$$p_{\text{data}}(\mathbf{x})$$



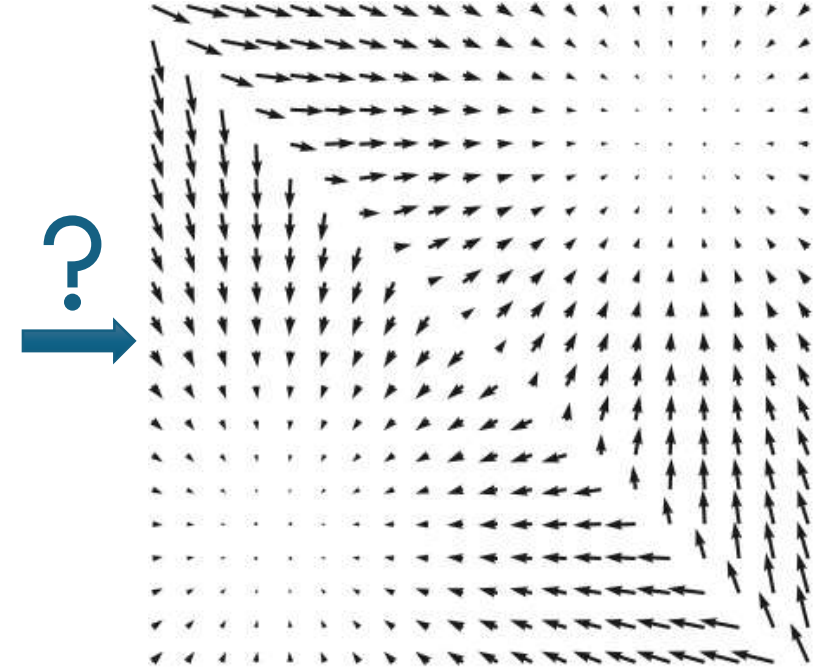
i.i.d. samples

$$\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$$



Score function

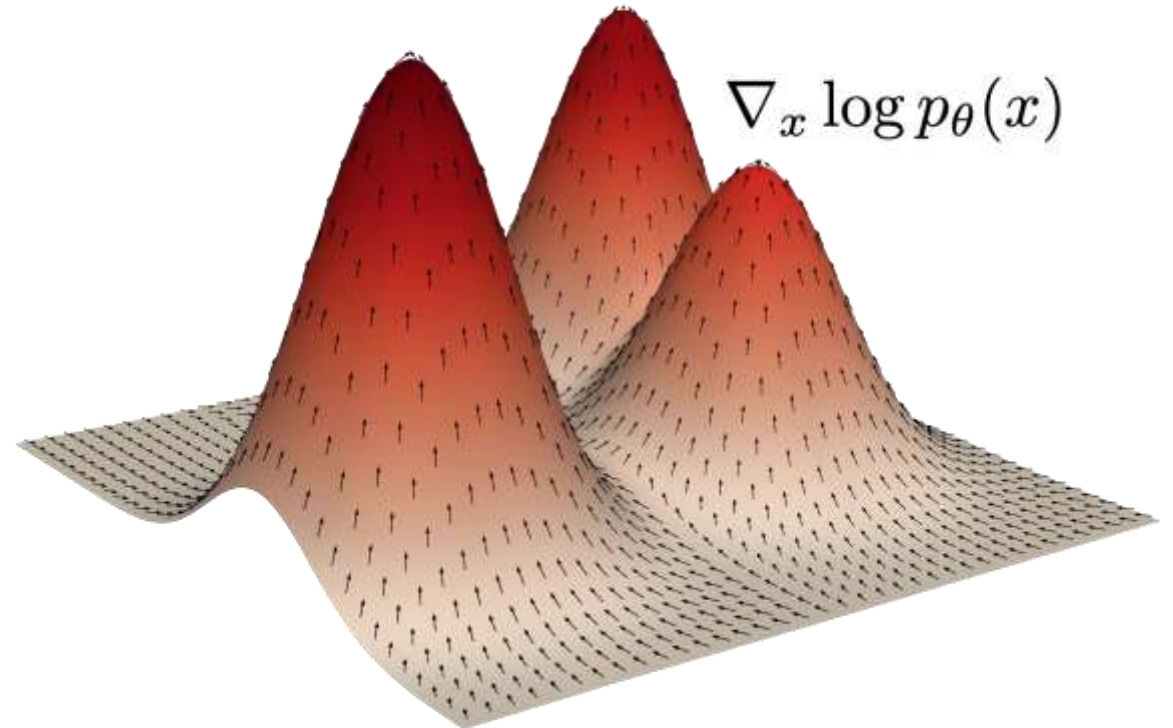
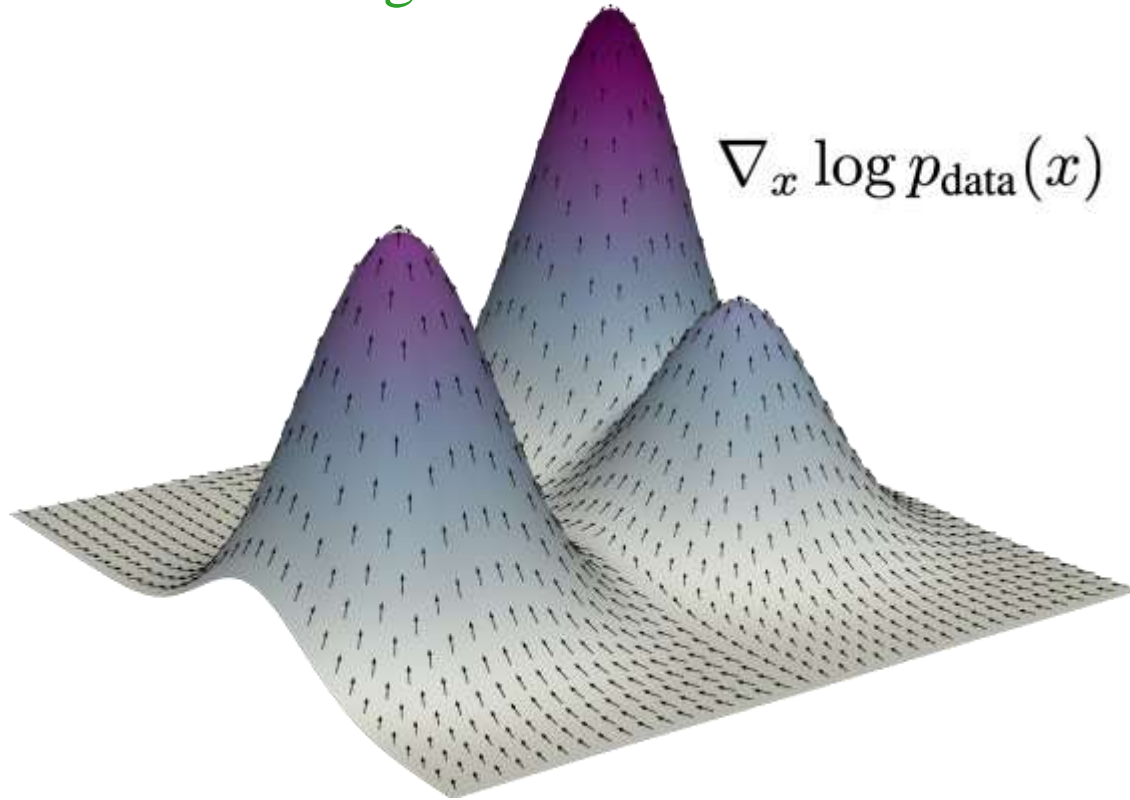
$$\mathbf{s}_\theta(\mathbf{x}) \approx \nabla_{\mathbf{x}} \log p_{\text{data}}(\mathbf{x})$$



Score Matching

- Instead of parametrizing p , we can parametrize the score

$$\underbrace{D_F(p_{\text{data}}(x) \parallel p_{\theta}(x))}_{\text{Fisher divergence}} = \mathbb{E}_{p_{\text{data}}(x)} \left[\frac{1}{2} \left\| \underbrace{\nabla_x \log p_{\text{data}}(x)}_{\text{score of data}} - \underbrace{\nabla_x \log p_{\theta}(x)}_{\text{parameterized score}} \right\|^2 \right]$$

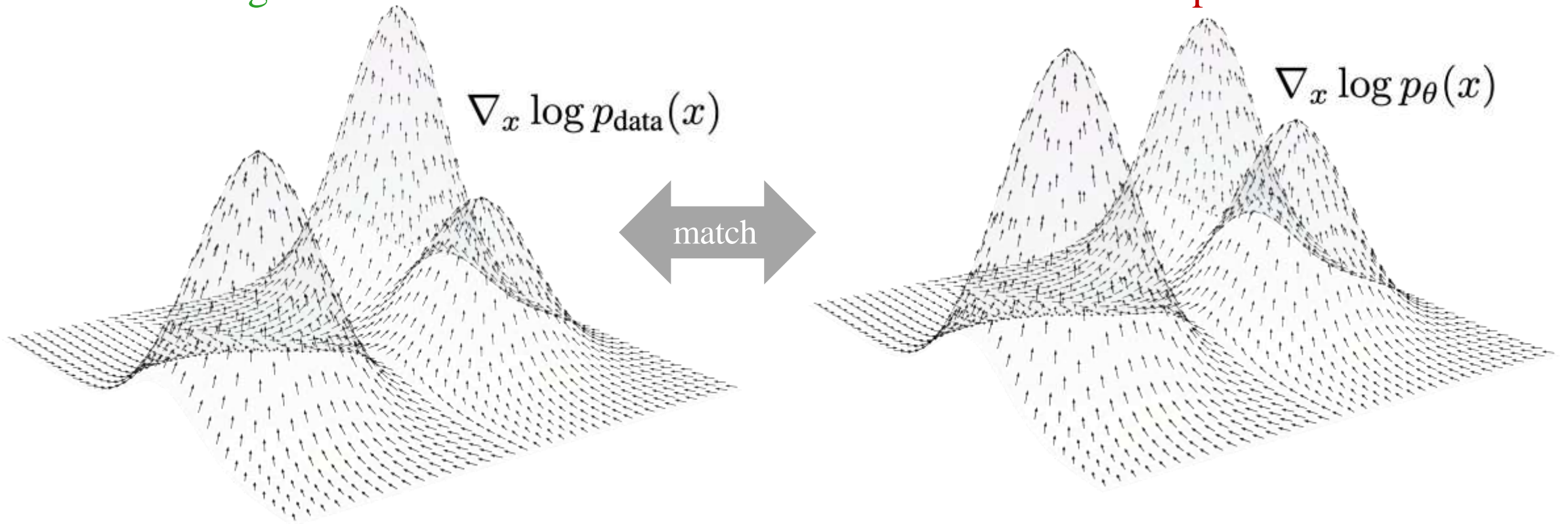


*only visualize

Score Matching

- Instead of parametrizing p , we can parametrize the score

$$\underbrace{D_F(p_{\text{data}}(x) \parallel p_{\theta}(x))}_{\text{Fisher divergence}} = \mathbb{E}_{p_{\text{data}}(x)} \left[\frac{1}{2} \left\| \underbrace{\nabla_x \log p_{\text{data}}(x)}_{\text{score of data}} - \underbrace{\nabla_x \log p_{\theta}(x)}_{\text{parameterized score}} \right\|^2 \right]$$

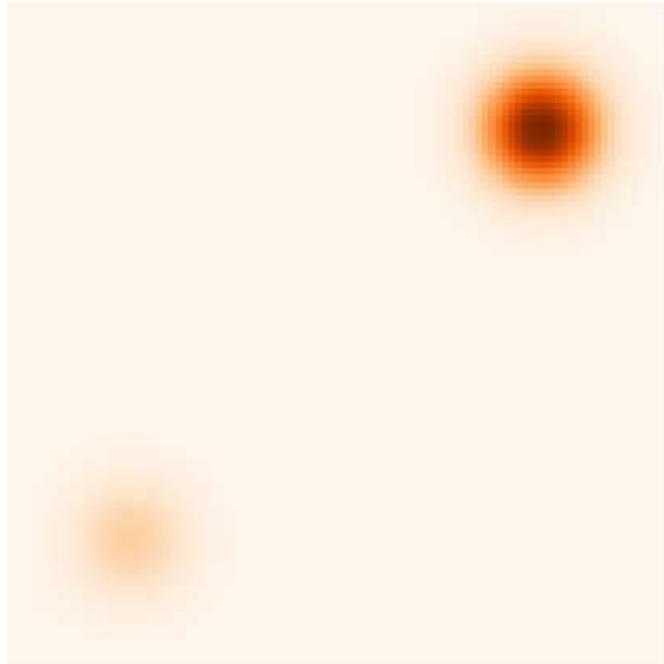


*only visualize

Score estimation by training score-based models

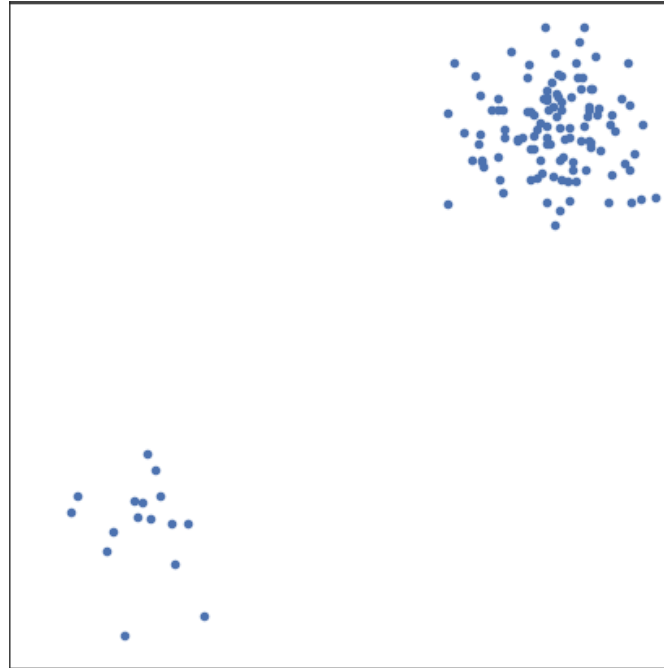
Probability density

$$p_{\text{data}}(\mathbf{x})$$



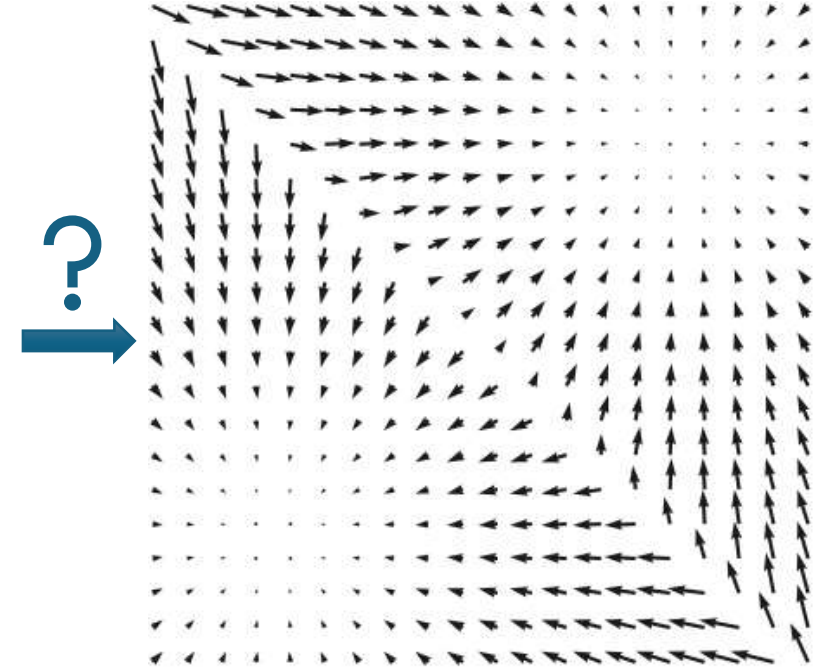
i.i.d. samples

$$\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$$



Score function

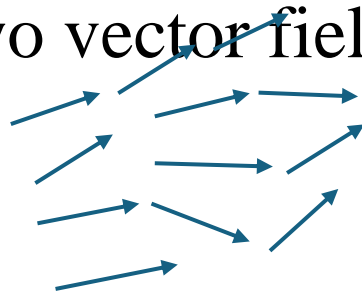
$$\mathbf{s}_\theta(\mathbf{x}) \approx \nabla_{\mathbf{x}} \log p_{\text{data}}(\mathbf{x})$$



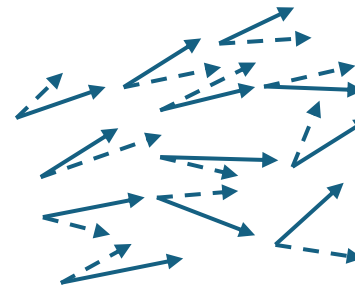
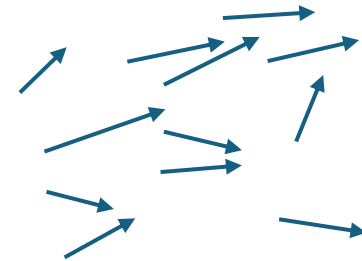
Score estimation by training score-based models

- Given: i.i.d. samples $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\} \sim p_{\text{data}}(\mathbf{x})$
- Task: Estimating the score $\nabla_{\mathbf{x}} \log p_{\text{data}}(\mathbf{x})$
- Score Model: A learnable vector-valued function $\mathbf{s}_{\theta}(\mathbf{x}) : \mathbb{R}^d \rightarrow \mathbb{R}^d$
- Goal: $\mathbf{s}_{\theta}(\mathbf{x}) \approx \nabla_{\mathbf{x}} \log p_{\text{data}}(\mathbf{x})$
- How to compare two vector fields of scores?

$\nabla_{\mathbf{x}} \log p_{\text{data}}(\mathbf{x})$



$\mathbf{s}_{\theta}(\mathbf{x})$



Average
Euclidean distance
over the space

Denoising Score Matching

- with noised data $\tilde{x} := x + \epsilon$, it can be proven: [Vincent, 2011]

$$\underline{\underline{D_F(q(\tilde{x}) \parallel p_\theta(\tilde{x}))}} = \mathbb{E}_{\underline{\underline{q(x, \tilde{x})}}} \left[\frac{1}{2} \left\| \underline{\underline{\nabla_{\tilde{x}} \log q(\tilde{x} | x)}} - \underline{\underline{\nabla_{\tilde{x}} \log p_\theta(\tilde{x})}} \right\|^2 \right] + \text{constant}$$

Fisher divergence of noised data joint distribution score of conditional parameterized score

Denoising Score Matching

- with noised data $\tilde{x} := x + \epsilon$, it can be proven: [Vincent, 2011]

$$D_F(q(\tilde{x}) \parallel p_\theta(\tilde{x})) = \mathbb{E}_{q(x, \tilde{x})} \left[\frac{1}{2} \left\| \underline{\nabla_{\tilde{x}} \log q(\tilde{x} | x)} - \underline{\nabla_{\tilde{x}} \log p_\theta(\tilde{x})} \right\|^2 \right] + \text{constant}$$

Gaussian
noise:

$$= \frac{1}{\sigma^2} (x - \tilde{x})$$

$-\epsilon$


a network to predict
(negative) noise

Langevin Dynamics

- Given a score function, we can sample x from p by iterating:

$$x_t \leftarrow x_{t-1} + \underbrace{\left(\frac{\sigma^2}{2}\right)}_{\text{step size}} \underbrace{\nabla_x \log p_\theta(x_{t-1})}_{\text{score function}} + \sigma \underbrace{z_t}_{\mathcal{N}(0, \mathbf{I})}$$

(don't need to know p)

 (neg) gradient of energy

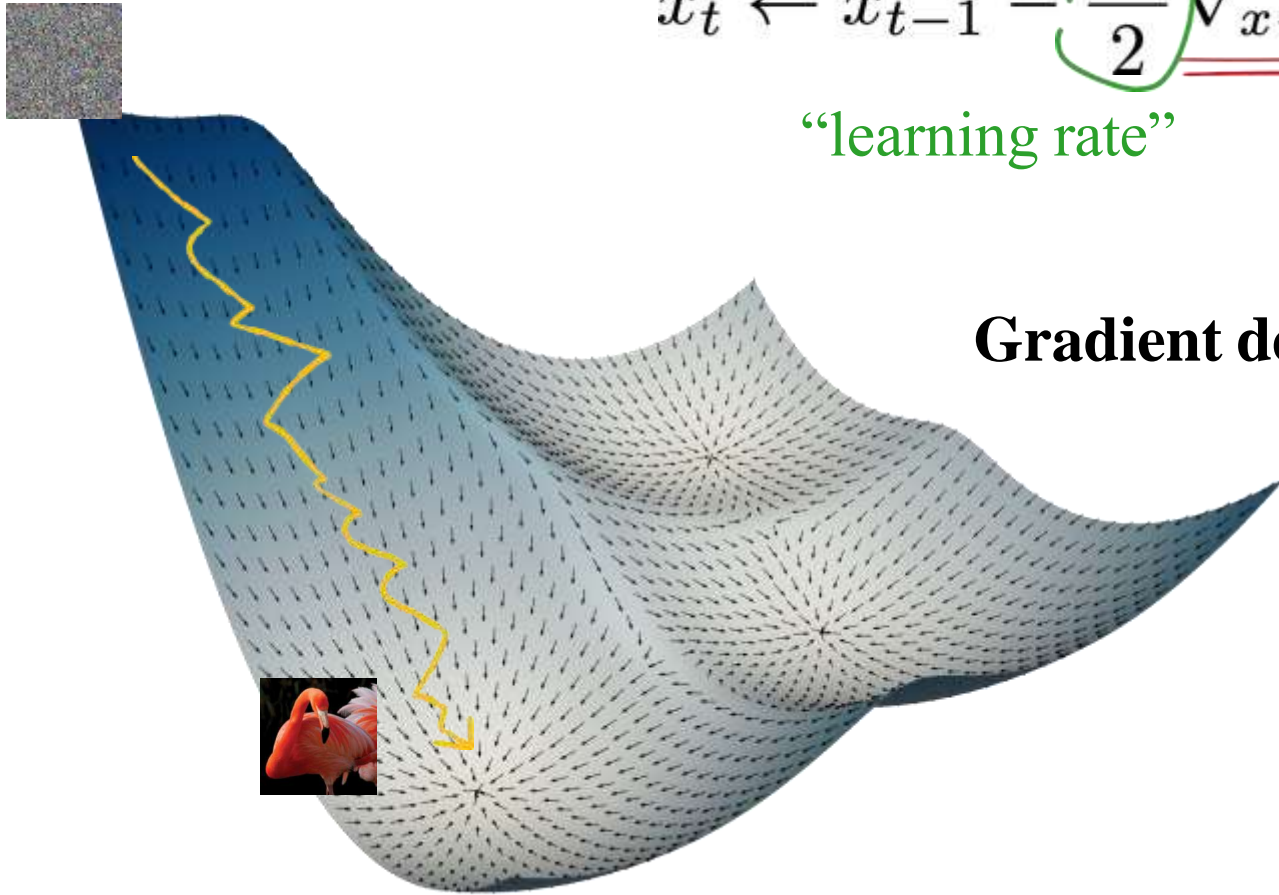
$$-\nabla_x E_\theta(x_{t-1})$$

Langevin Dynamics

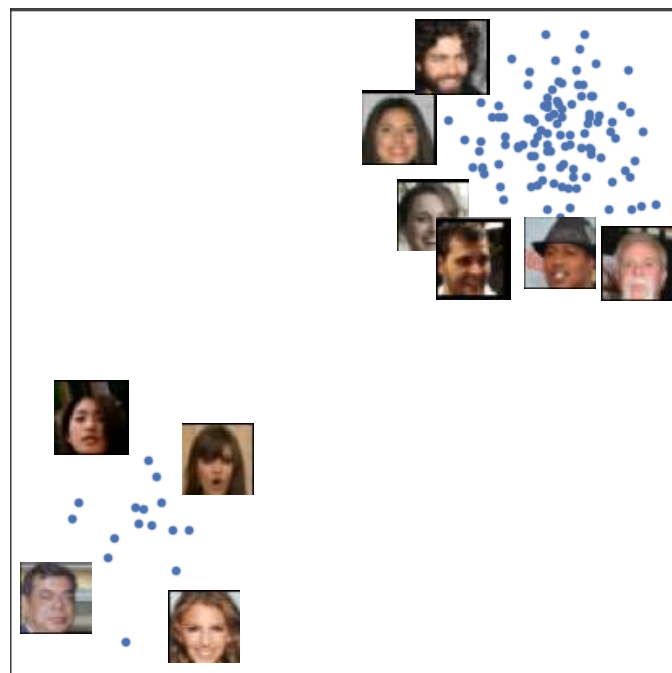
- Given a score function, we can sample x from p by iterating:

$$x_t \leftarrow x_{t-1} - \underbrace{\frac{\sigma^2}{2}}_{\text{“learning rate”}} \underbrace{\nabla_x E_\theta(x_{t-1})}_{\text{gradient}} + \underbrace{(\sigma z_t)}_{\text{perturbation}}$$

Gradient decent in the energy landscape



Score-based generative modeling

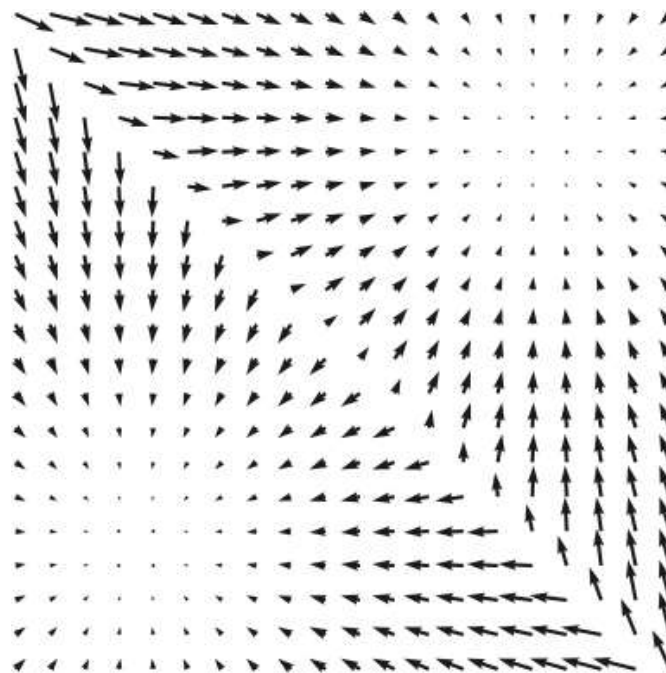


Data samples

$$\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\} \sim p_{\text{data}}(\mathbf{x})$$



Score
Matching



Scores

$$\mathbf{s}_{\theta}(\mathbf{x}) \approx \nabla_{\mathbf{x}} \log p_{\text{data}}(\mathbf{x})$$

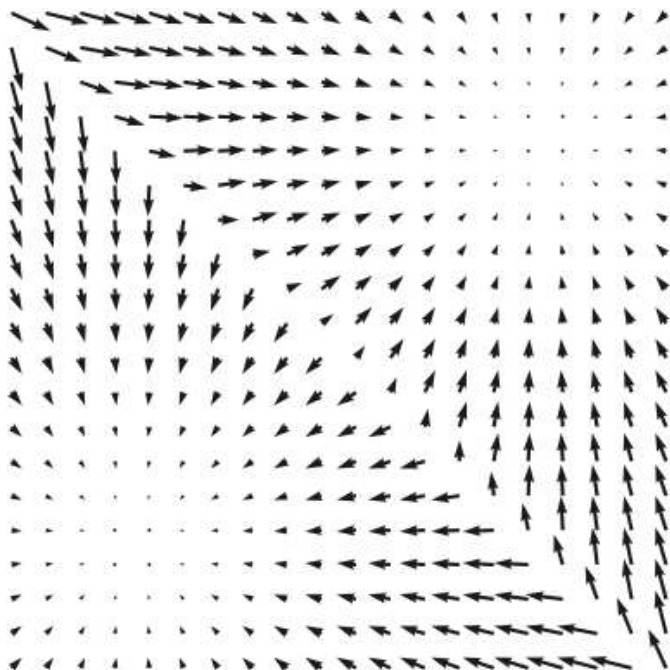


?



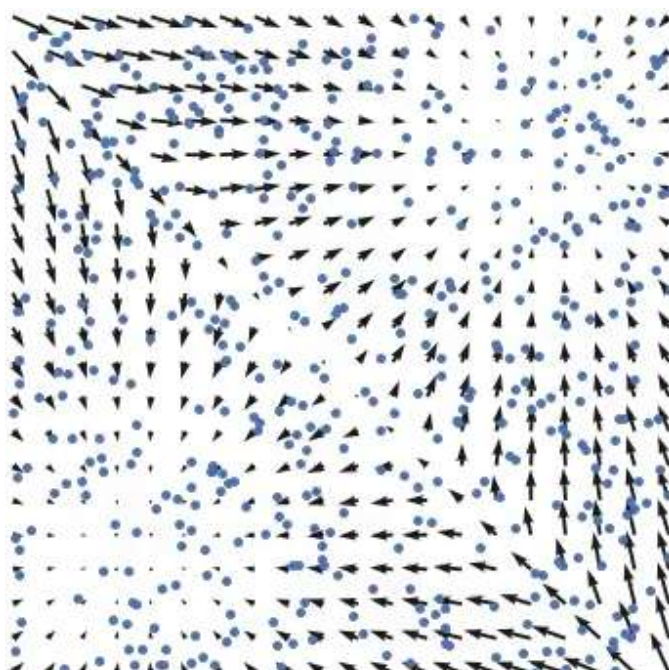
New samples

From scores to samples: Langevin MCMC



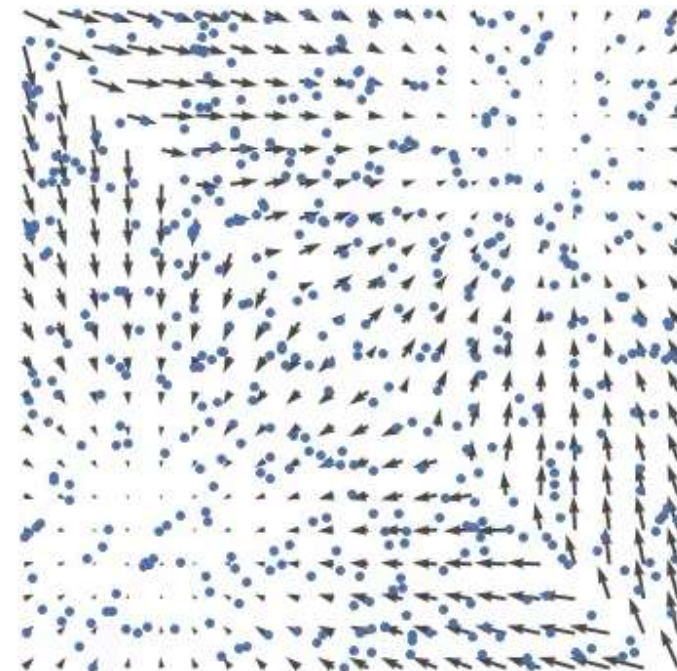
Scores

$$\mathbf{s}_\theta(\mathbf{x})$$



Follow the scores

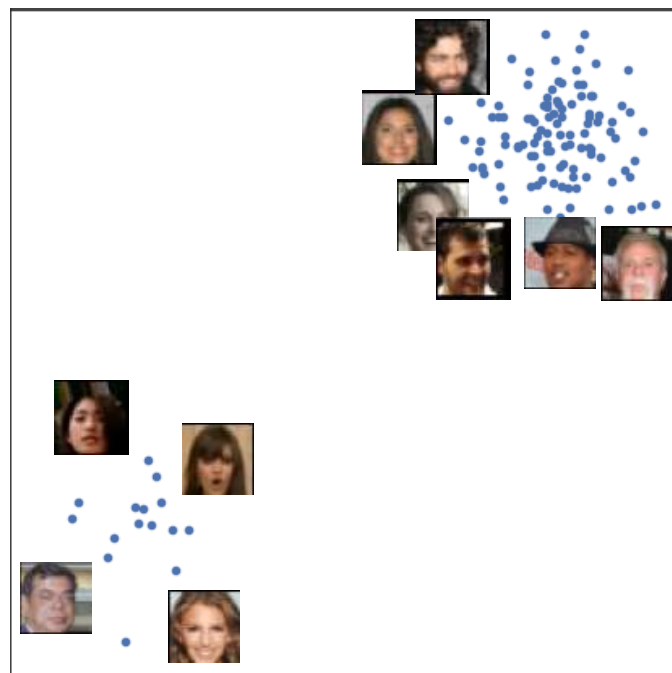
$$\tilde{\mathbf{x}}_{t+1} \leftarrow \tilde{\mathbf{x}}_t + \frac{\epsilon}{2} \mathbf{s}_\theta(\tilde{\mathbf{x}}_t)$$



Follow noisy scores:
Langevin MCMC

$$\mathbf{z}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$
$$\tilde{\mathbf{x}}_{t+1} \leftarrow \tilde{\mathbf{x}}_t + \frac{\epsilon}{2} \mathbf{s}_\theta(\tilde{\mathbf{x}}_t) + \sqrt{\epsilon} \mathbf{z}_t$$

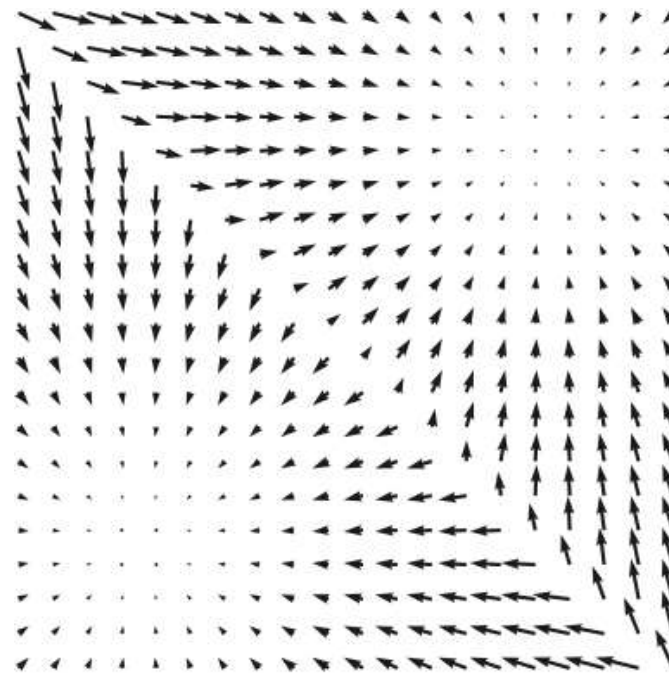
Score-based generative modeling



Data samples

$$\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\} \sim p_{\text{data}}(\mathbf{x})$$

score
matching



Scores

$$\mathbf{s}_{\theta}(\mathbf{x}) \approx \nabla_{\mathbf{x}} \log p_{\text{data}}(\mathbf{x})$$



Langevin
dynamics



New samples

(Recap) Diffusion algorithm

Algorithm 1 Training

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

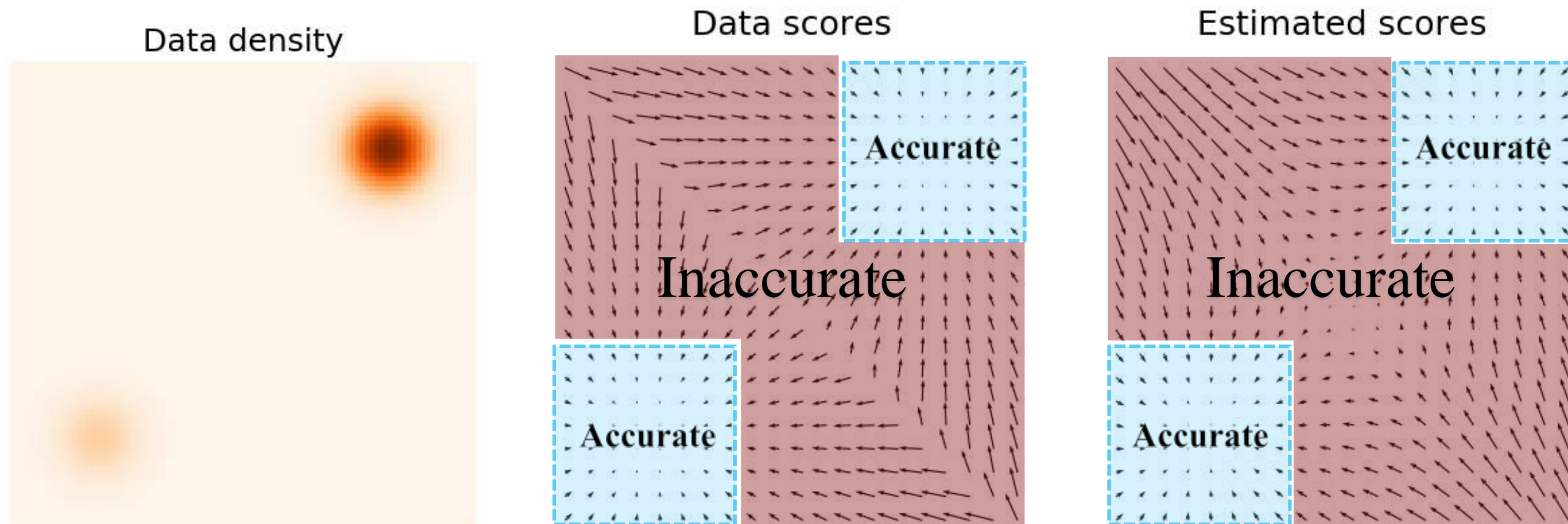
score function

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: $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \underbrace{\boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t)}_{\text{score function}} \right) + \sigma_t \mathbf{z}$
 - 5: **end for**
 - 6: **return** \mathbf{x}_0
-

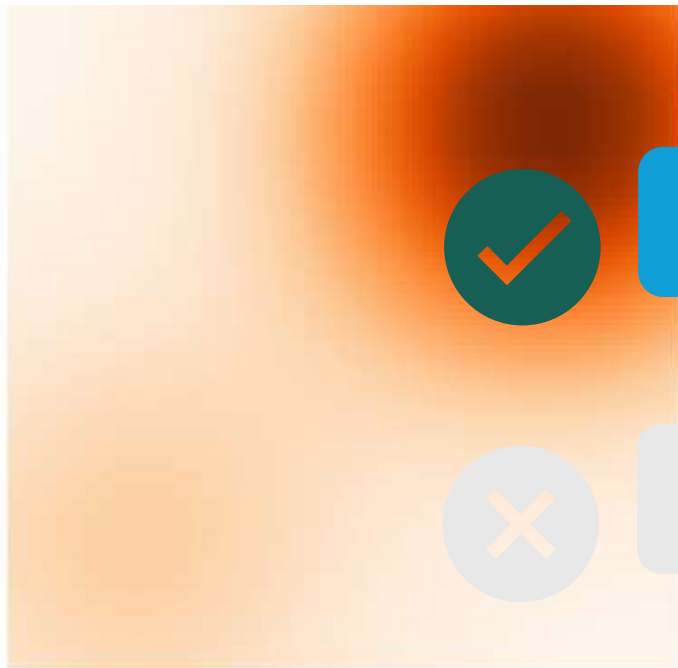
Langevin
Dynamics

Challenge in low data density regions

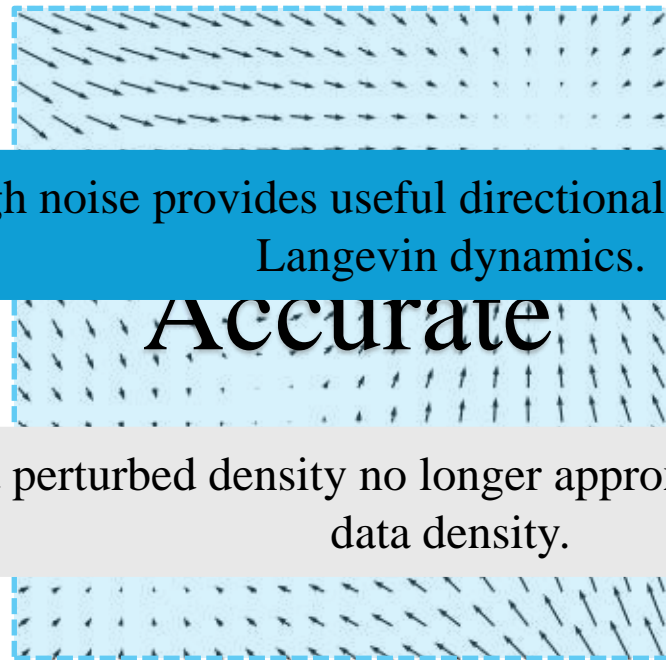


Improving score estimation by adding noise

Perturbed density

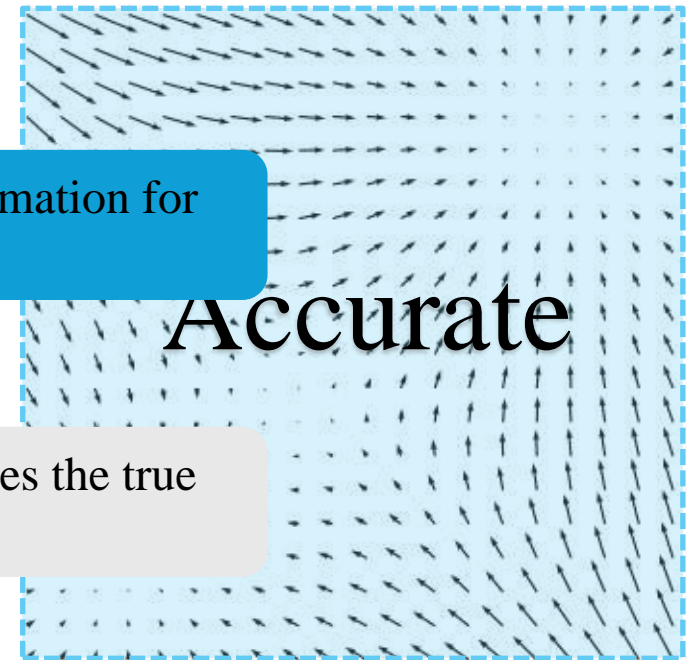


Perturbed scores



Accurate

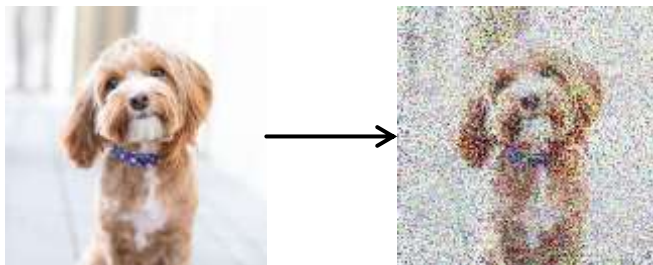
Estimated scores



Accurate

High noise provides useful directional information for Langevin dynamics.

But perturbed density no longer approximates the true data density.

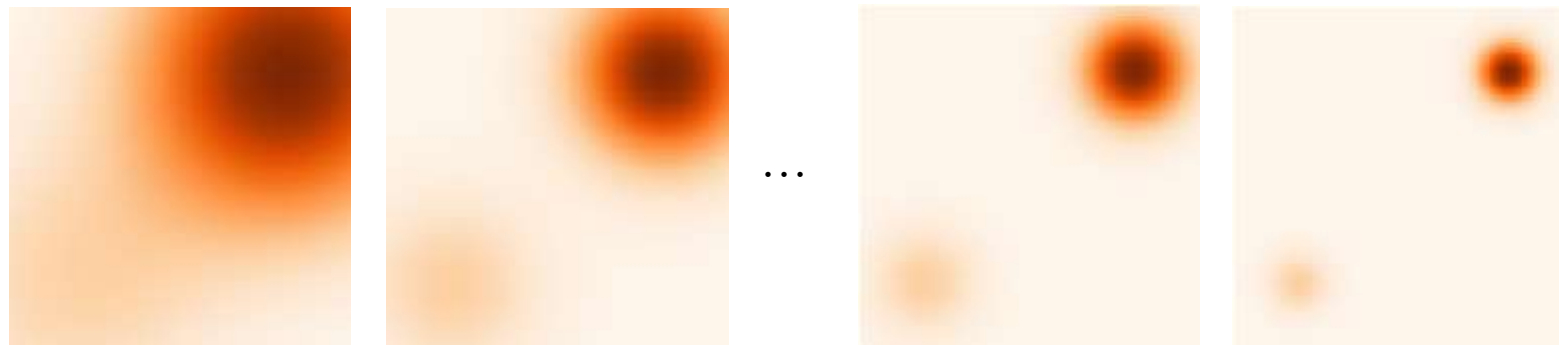


Multi-scale Noise Perturbation

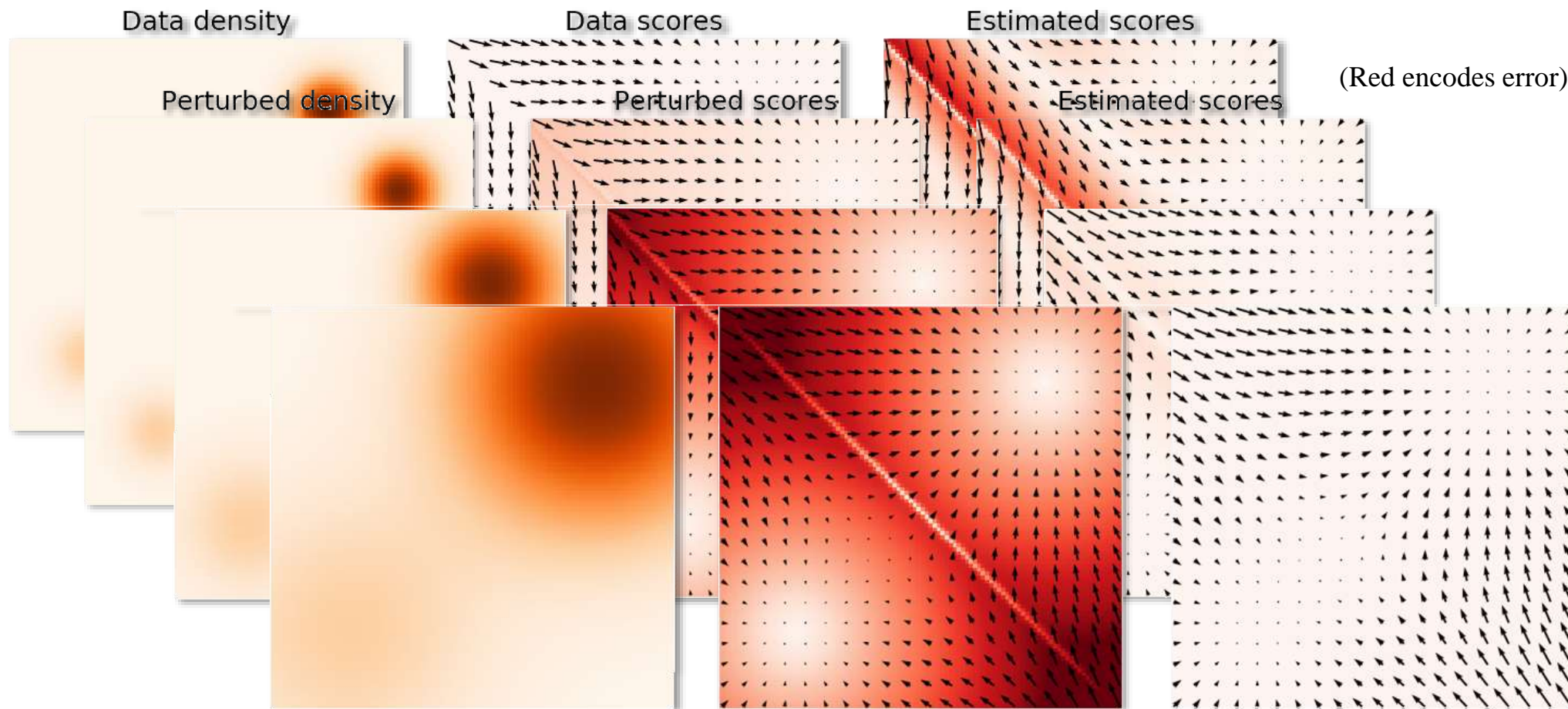
- How much noise to add?



- Multi-scale noise perturbations: $\sigma_1 > \sigma_2 > \dots > \sigma_{L-1} > \sigma_L$



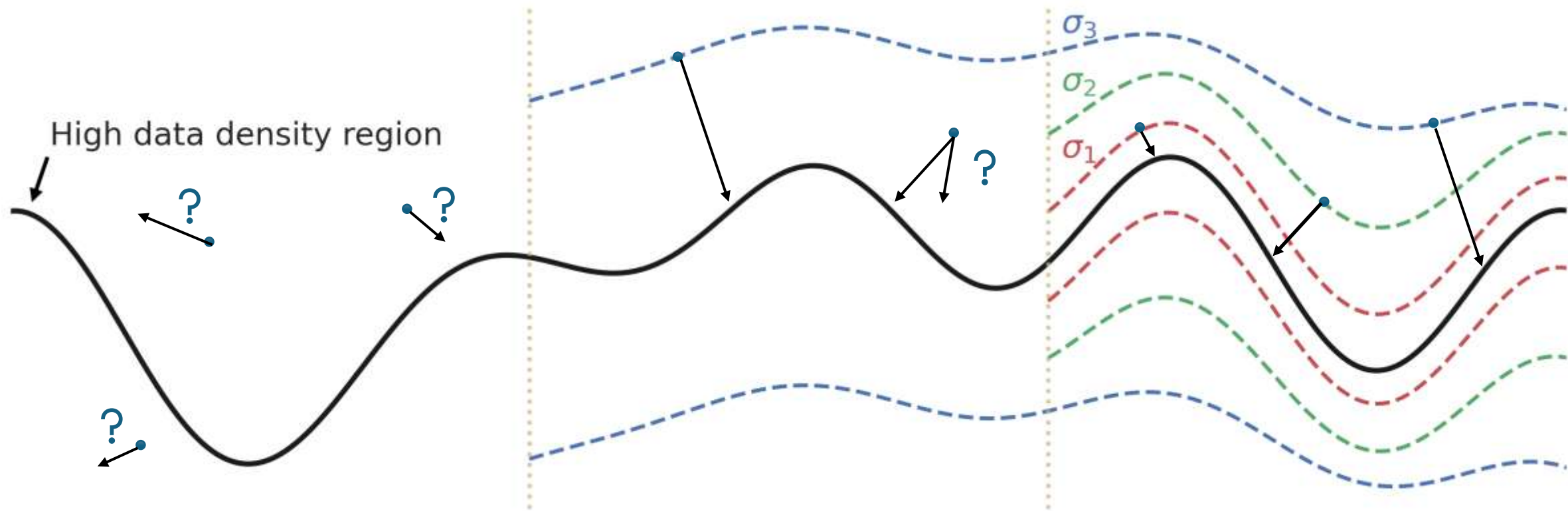
Trading off Data Quality and Estimation Accuracy



Worse data quality!

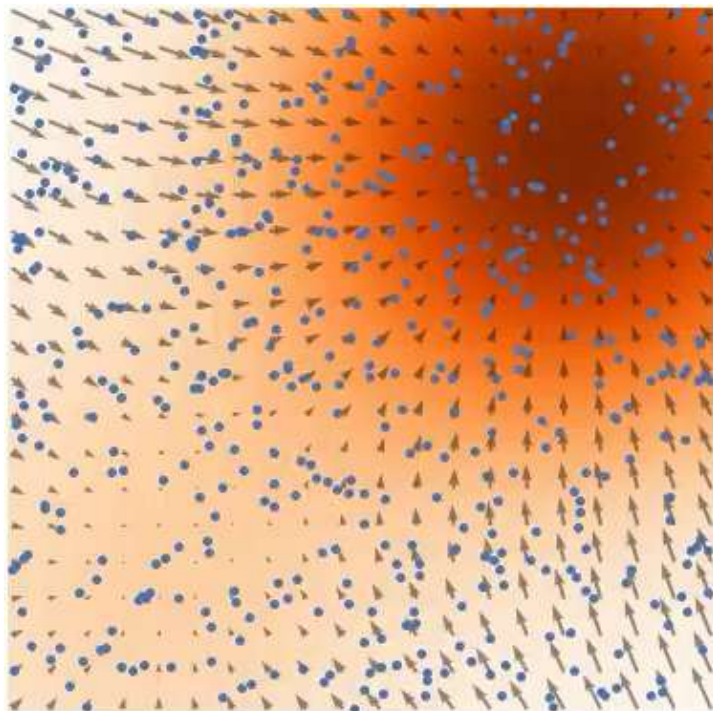
Better score estimation!

Using multiple noise scales

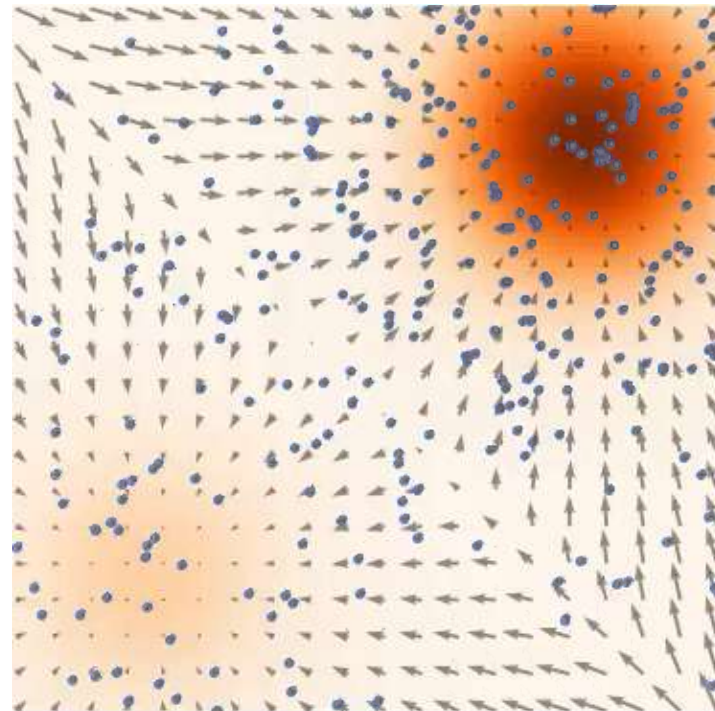


Annealed Langevin Dynamics: Joint Scores to Samples

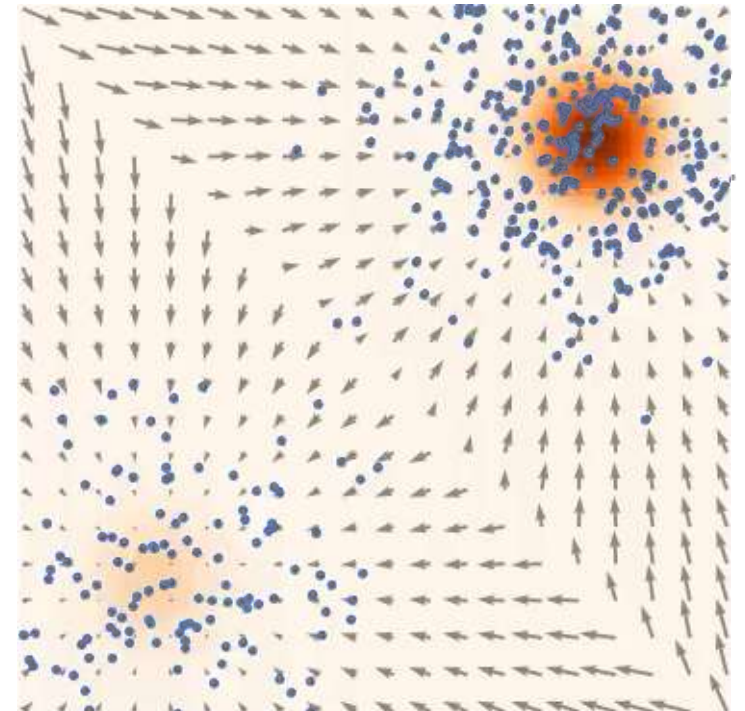
- Sample using $\sigma_1, \sigma_2, \dots, \sigma_L$ sequentially with Langevin dynamics.
- Anneal down the noise level.
- Samples used as initialization for the next level.



σ_1



σ_2



σ_3

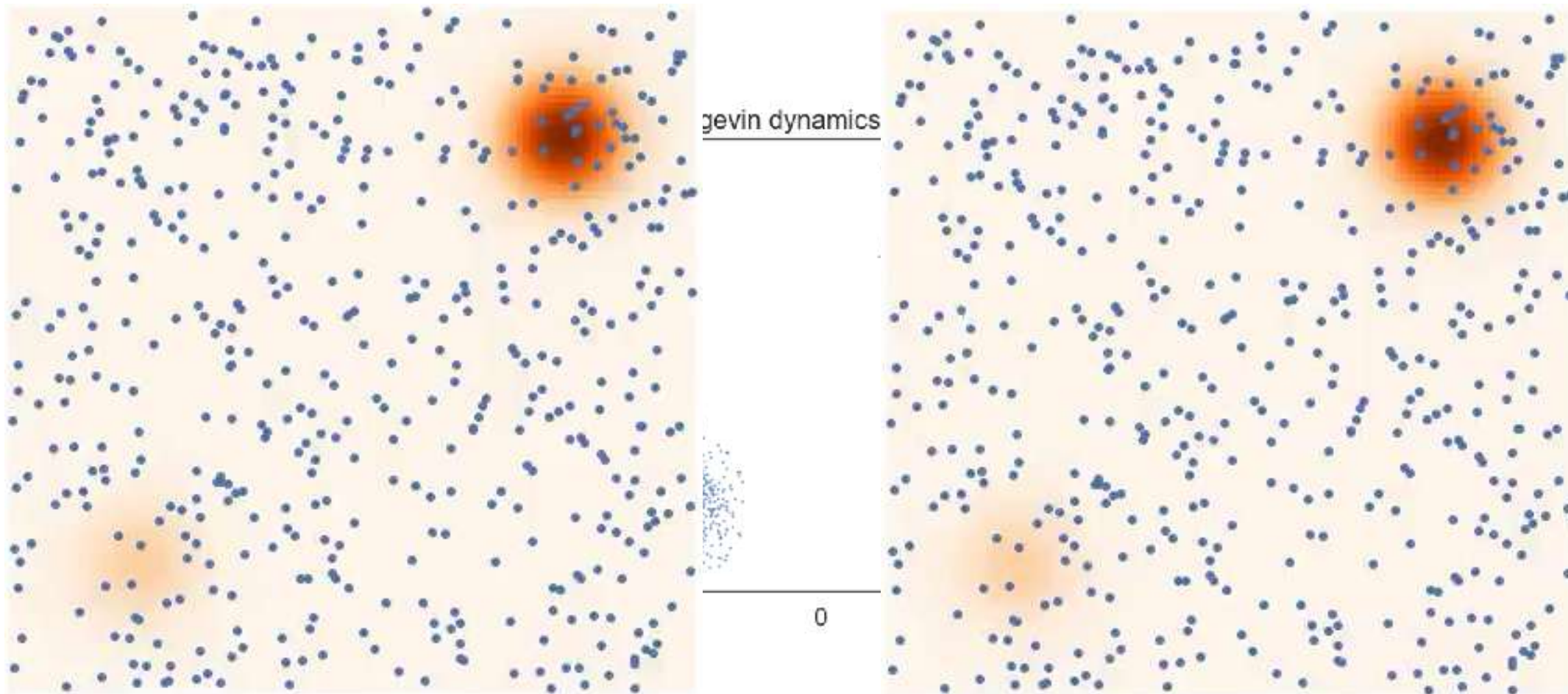
Annealed Langevin dynamics

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** $\tilde{\mathbf{x}}_T$
-

Comparison to the vanilla Langevin dynamics

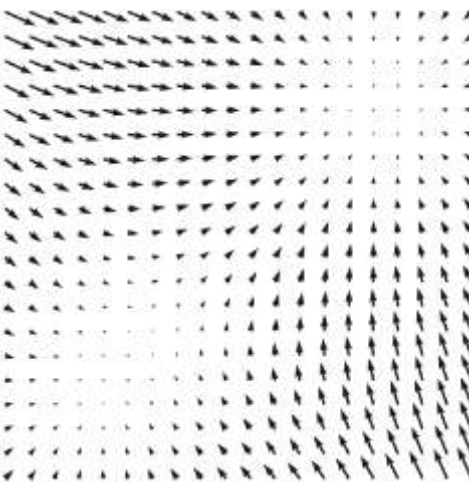
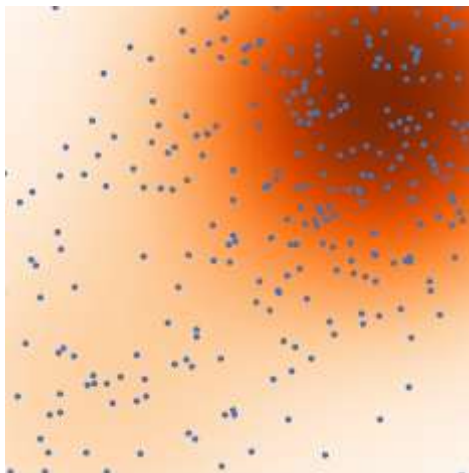


Langevin dynamics

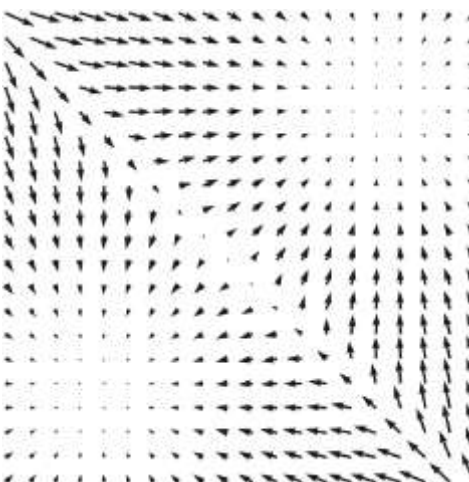
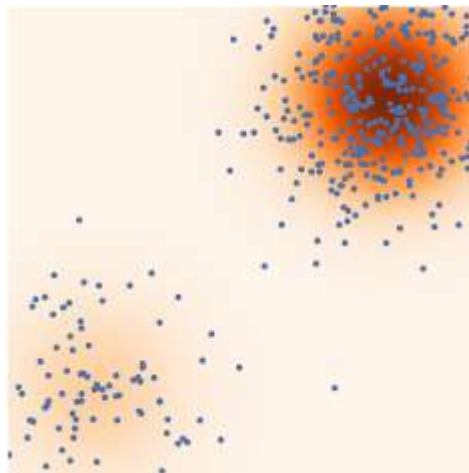
Annealed Langevin dynamics

Joint Score Estimation via Noise Conditional Score Networks

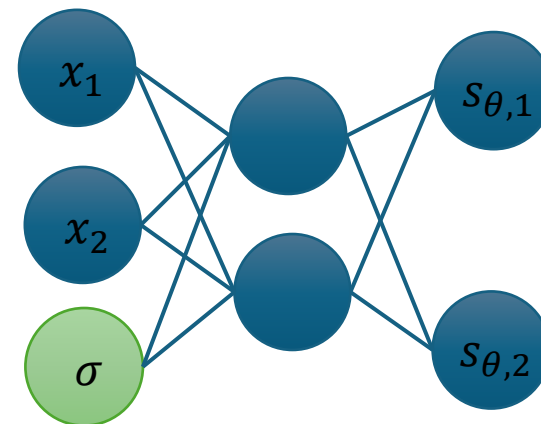
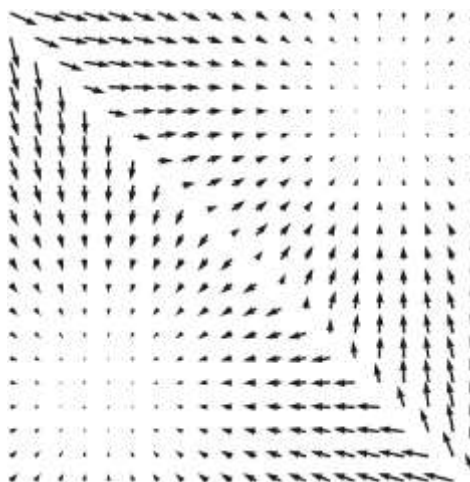
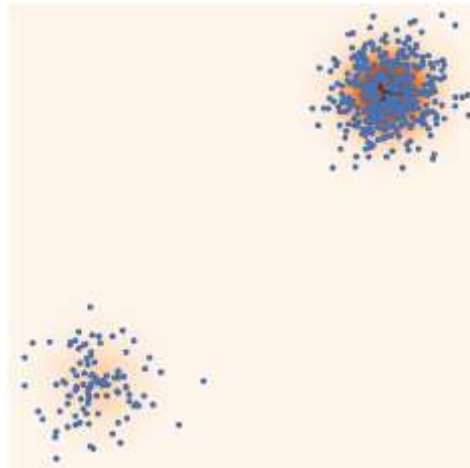
σ_1



σ_2



σ_3

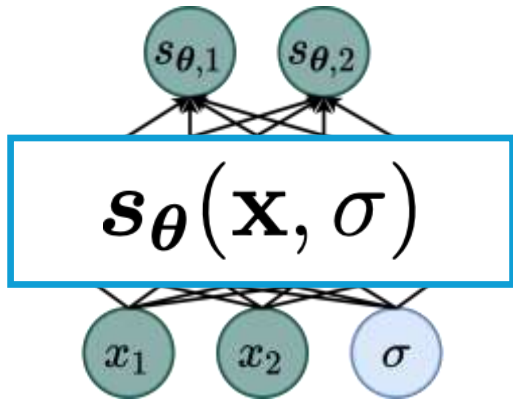
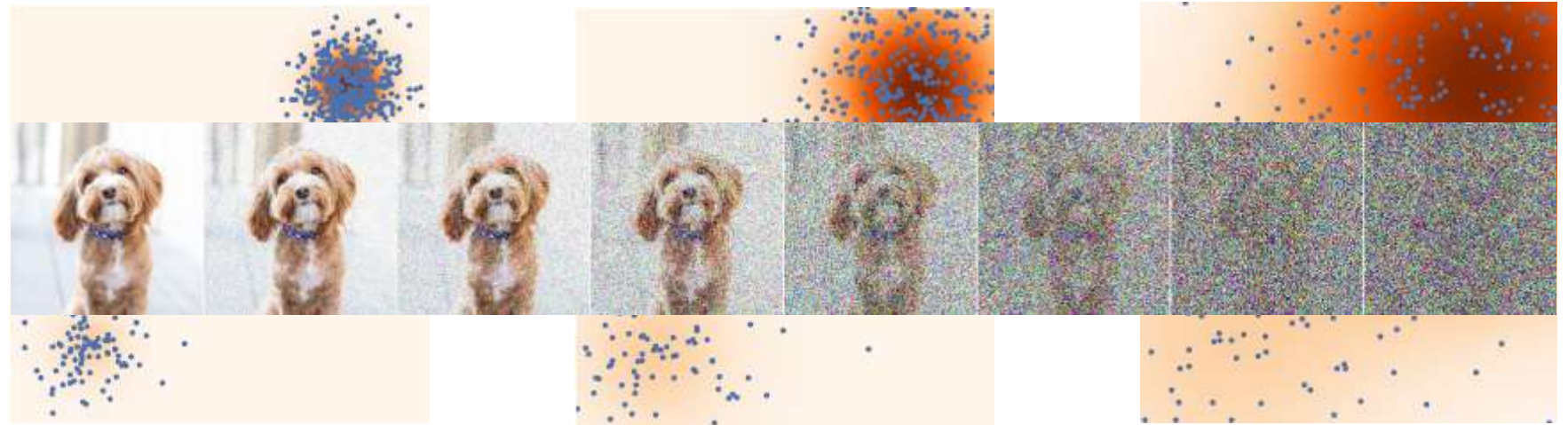


Noise Conditional Score
Network
(NCSN)

Using multiple noise levels

$$p_{\sigma_1}(\mathbf{x}) < p_{\sigma_2}(\mathbf{x}) < p_{\sigma_3}(\mathbf{x})$$

Data



Noise Conditional Score Model

Positive weighting function

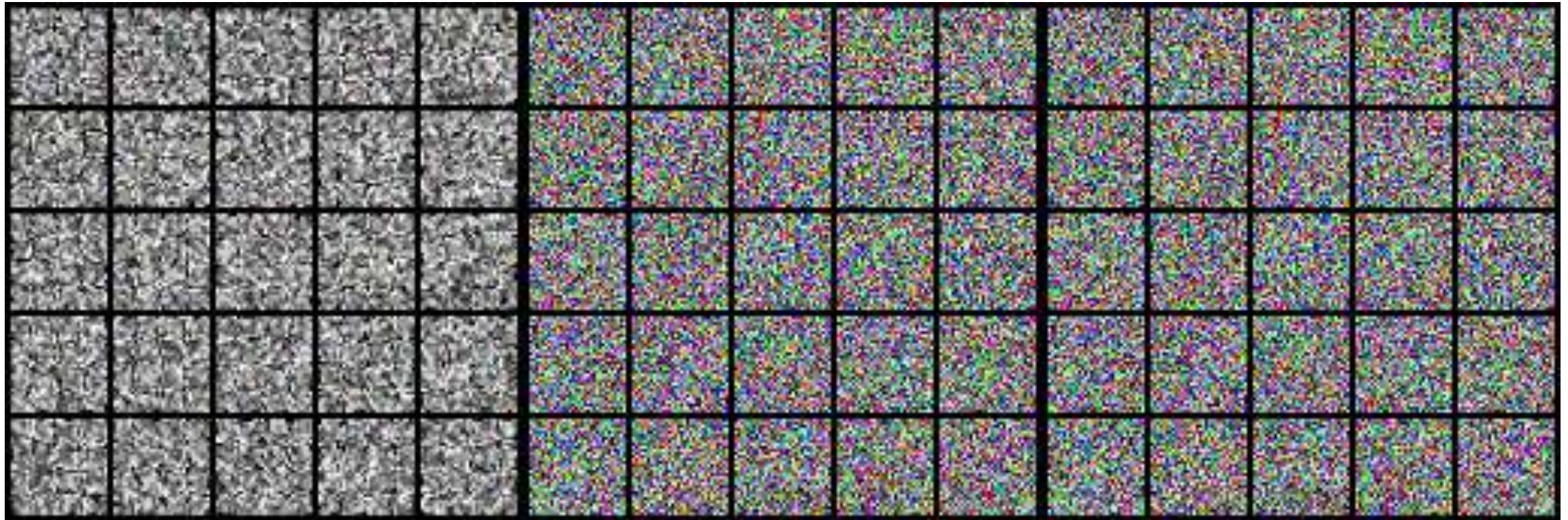
$$\frac{1}{N} \sum_{i=1}^N \lambda(\sigma_i) \mathbb{E}_{p_{\sigma_i}(\mathbf{x})} \left[\left\| \nabla_{\mathbf{x}} \log p_{\sigma_i}(\mathbf{x}) - \mathbf{s}_{\theta}(\mathbf{x}, \sigma_i) \right\|_2^2 \right]$$

Annealed Langevin dynamics

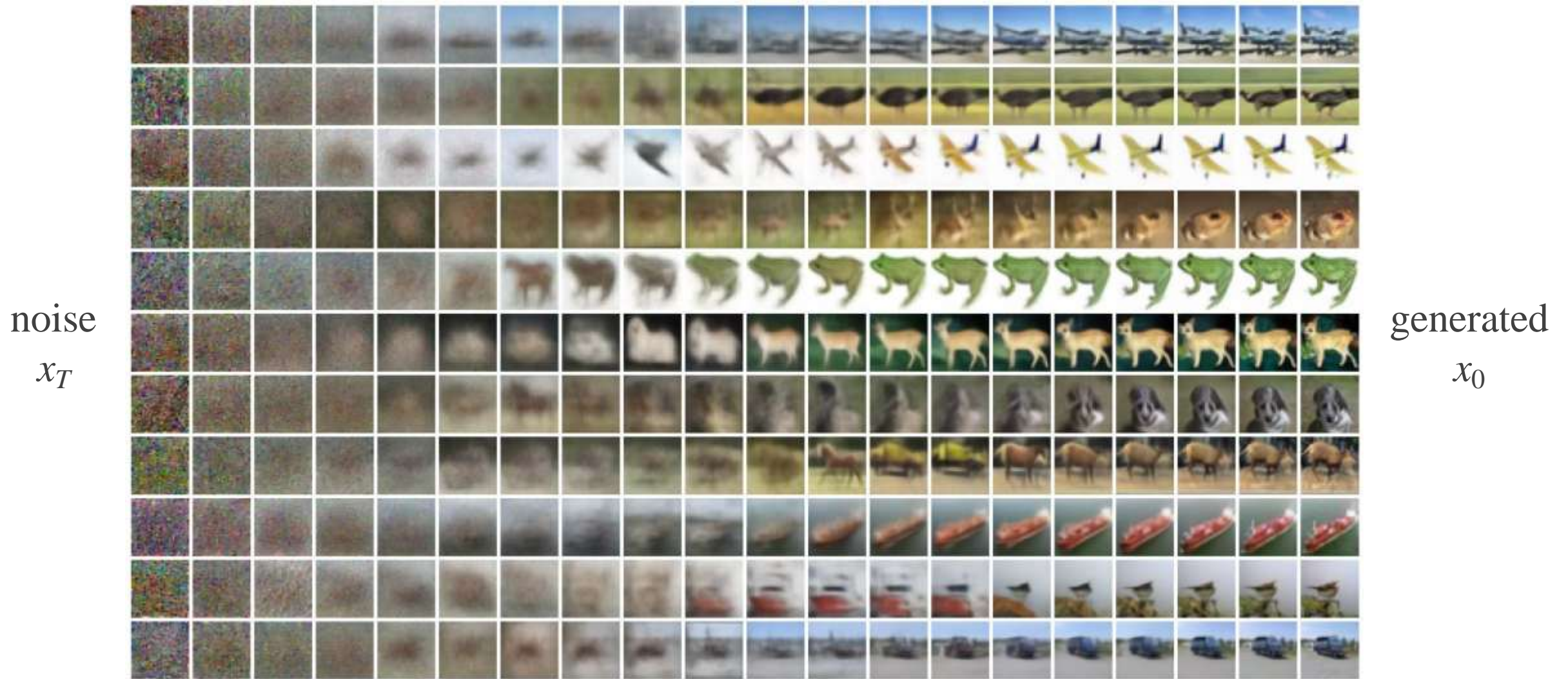
Score matching loss

$$s_{\theta}(\mathbf{x}, \sigma_1) \quad s_{\theta}(\mathbf{x}, \sigma_2) \quad s_{\theta}(\mathbf{x}, \sigma_3)$$

Experiments: Sampling



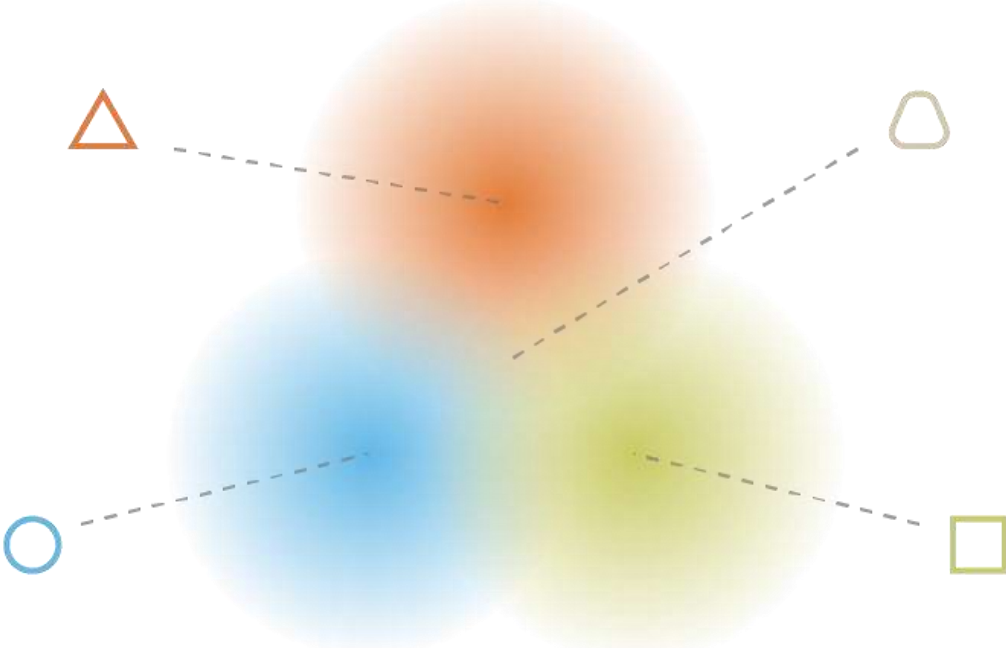
Example: Unconditional Generation on CIFAR-10



Example: shared intermediate latents



Problems

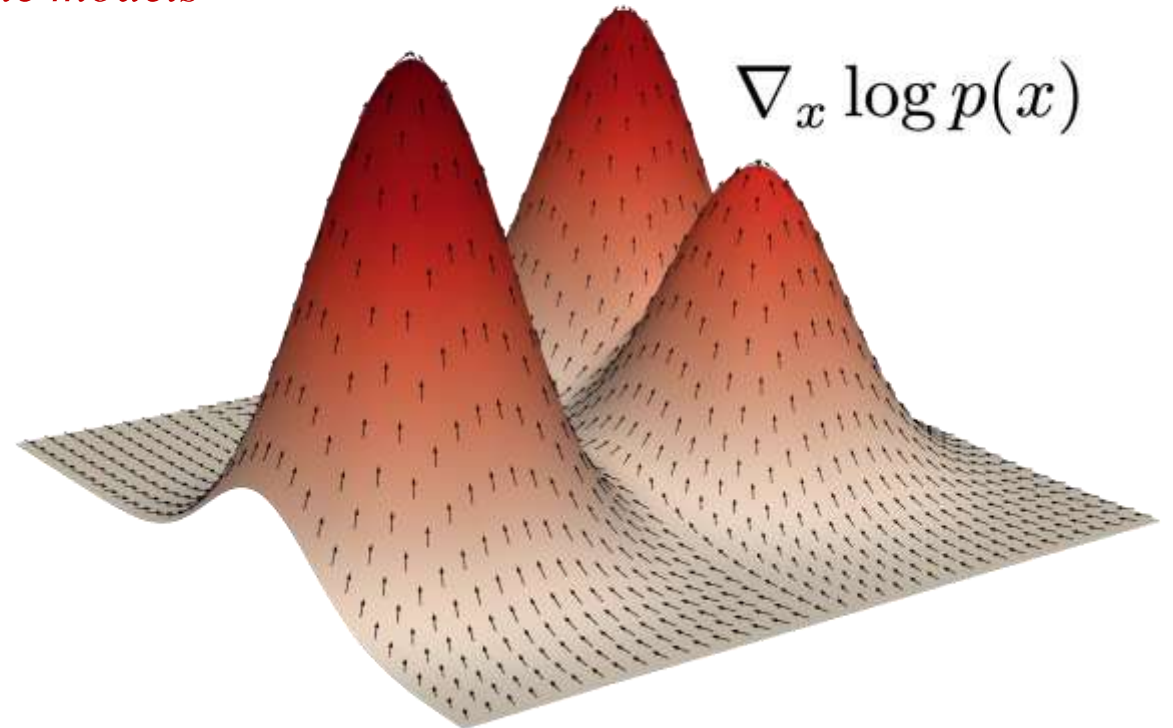
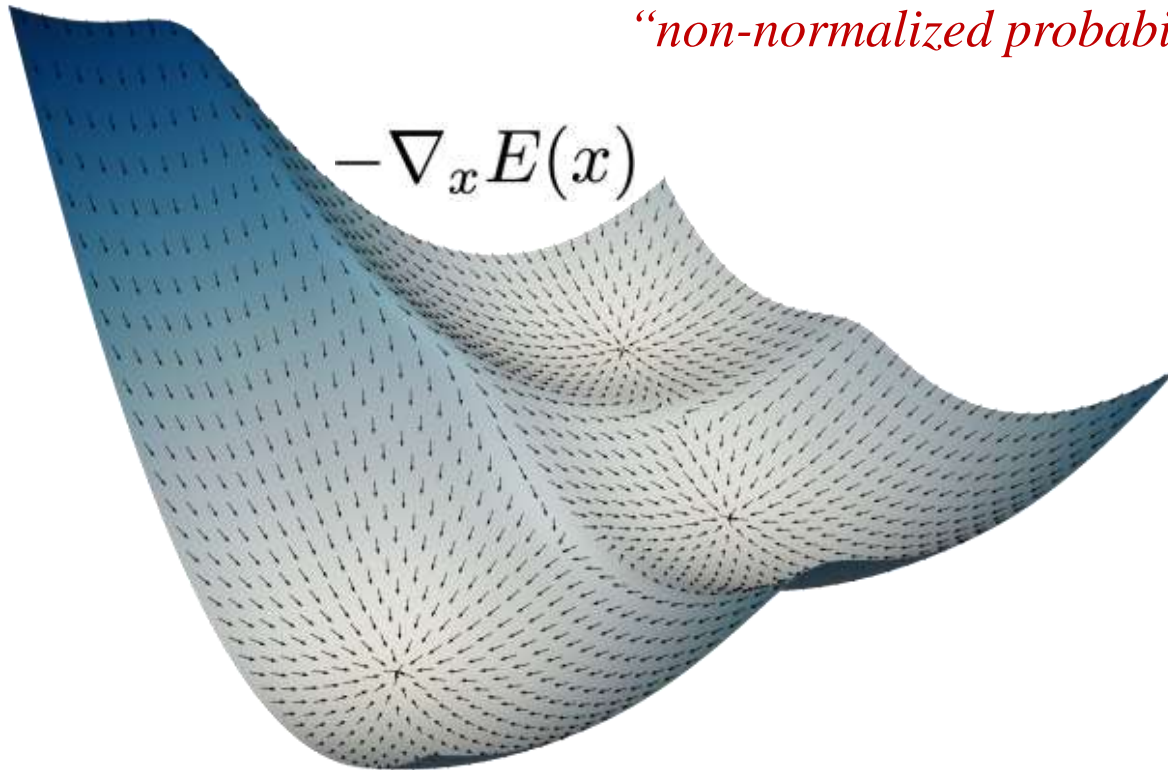


Energy-based Models

- “Score function”: gradient of log-probability

$$\nabla_x \log p(x) = -\nabla_x E(x)$$

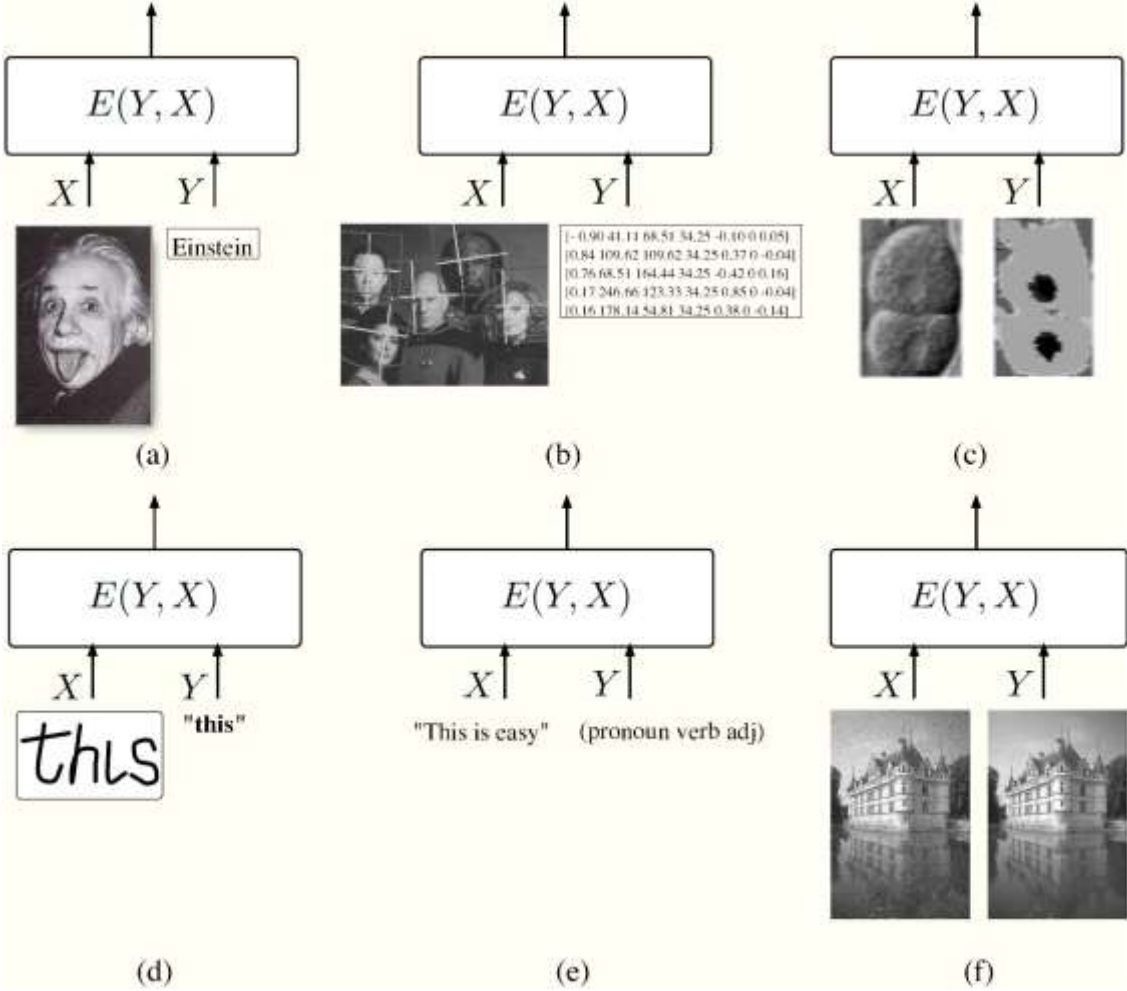
“non-normalized probabilistic models”



*only visualize directions

More about Energy-based Models ...

- At inference time, find a solution that minimizes energy



Various Perspectives on Diffusion Models ...

- Hierarchical VAE
- Energy-based Models and Score Matching
- Autoregressive models

be a fully expressive conditional distribution. With these choices, $D_{\text{KL}}(q(\mathbf{x}_T) \parallel p(\mathbf{x}_T)) = 0$, and minimizing $D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t))$ trains p_{θ} to copy coordinates $t + 1, \dots, T$ unchanged and to predict the t^{th} coordinate given $t + 1, \dots, T$. Thus, training p_{θ} with this particular diffusion is training an autoregressive model.

[Ho et al, 2020]

- SDE and ODE
- Normalizing Flows
- Recurrent Neural Networks