

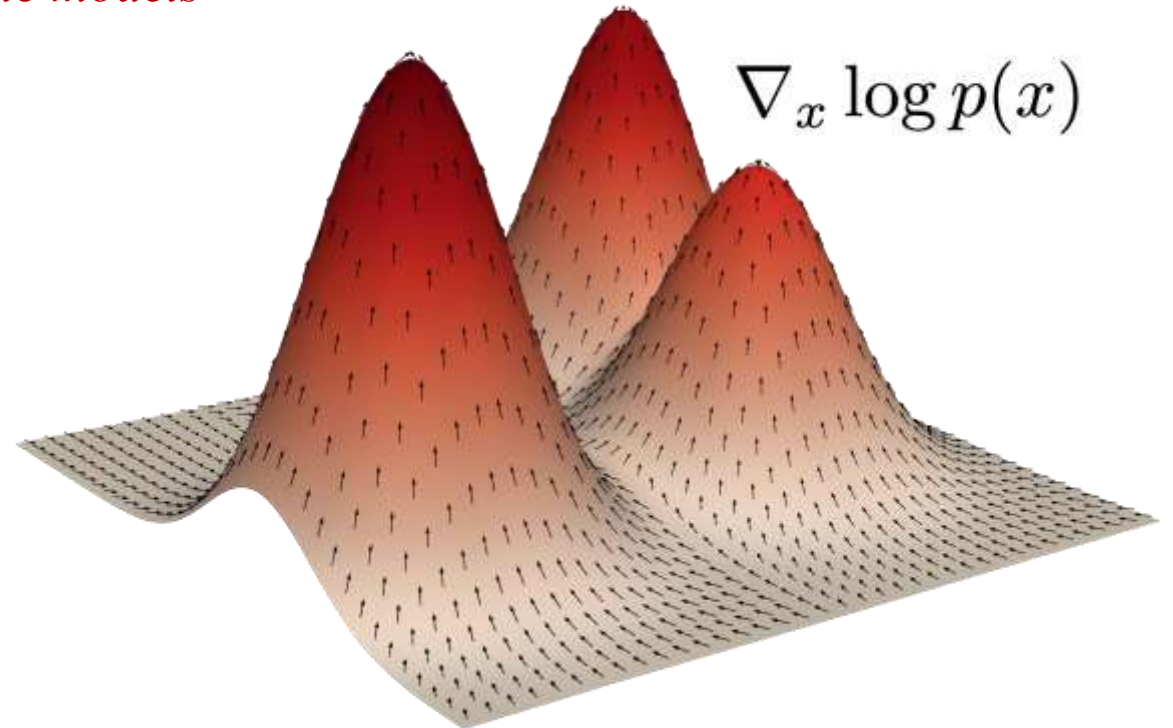
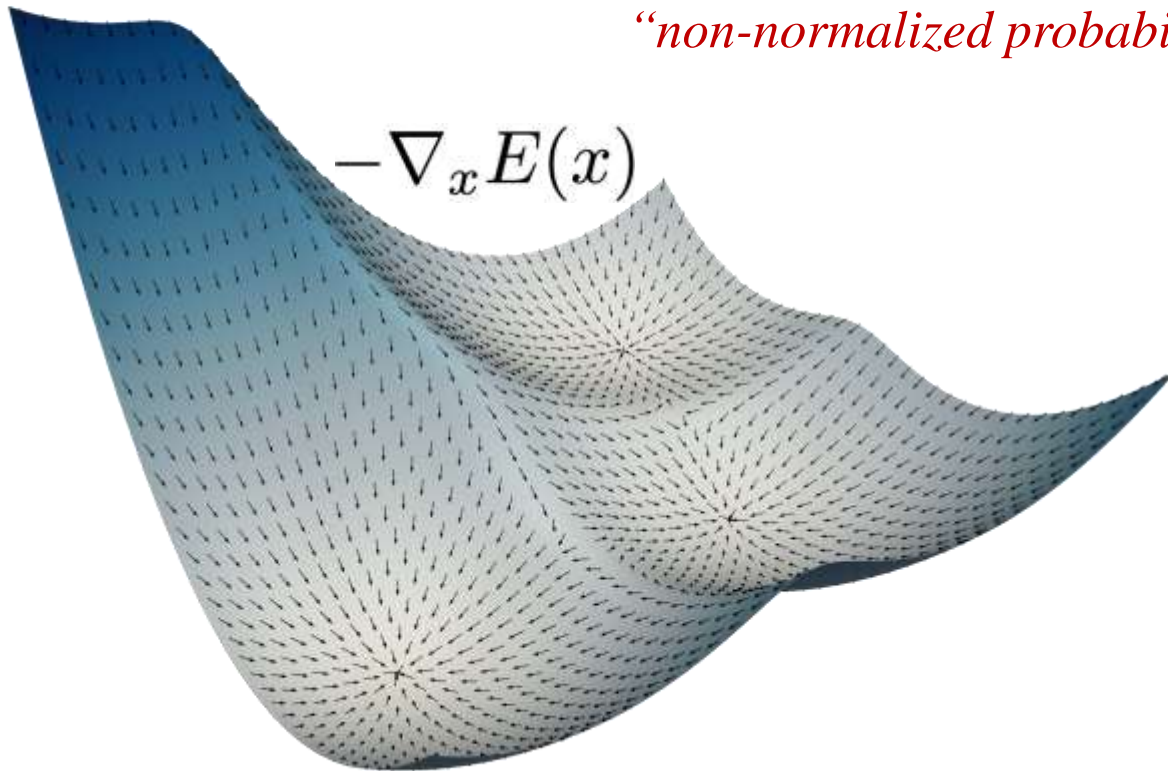
Recap.

Energy-based Models

- “Score function”: gradient of log-probability

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

“non-normalized probabilistic models”

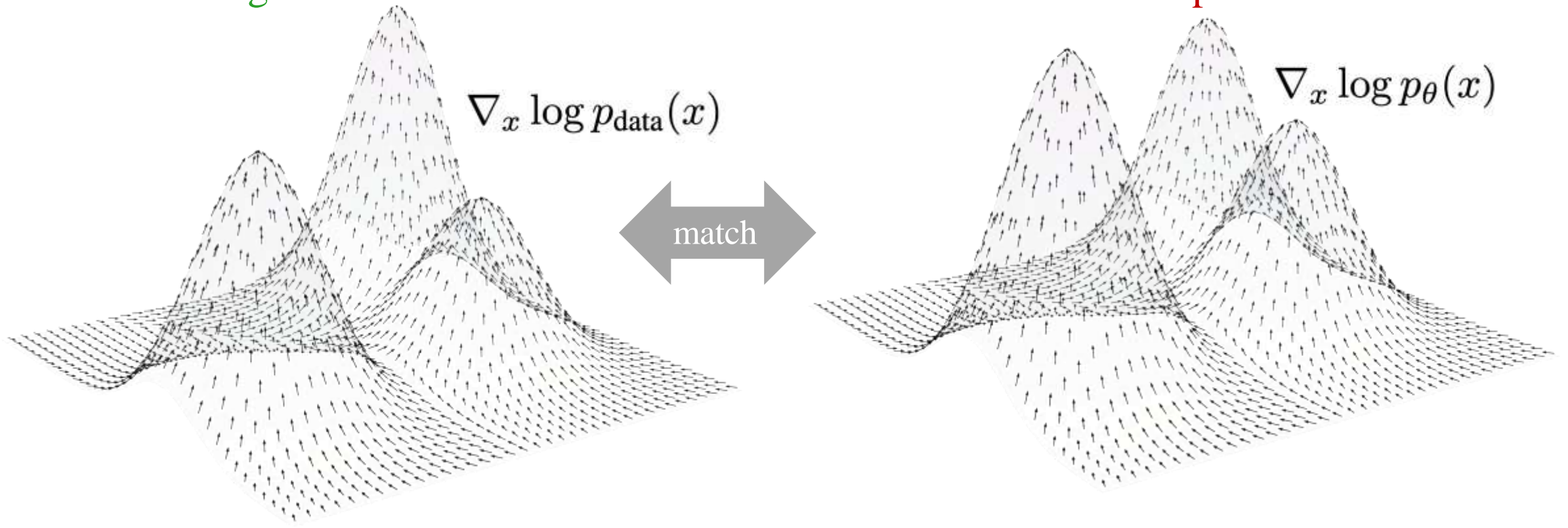


*only visualize directions

Score Matching

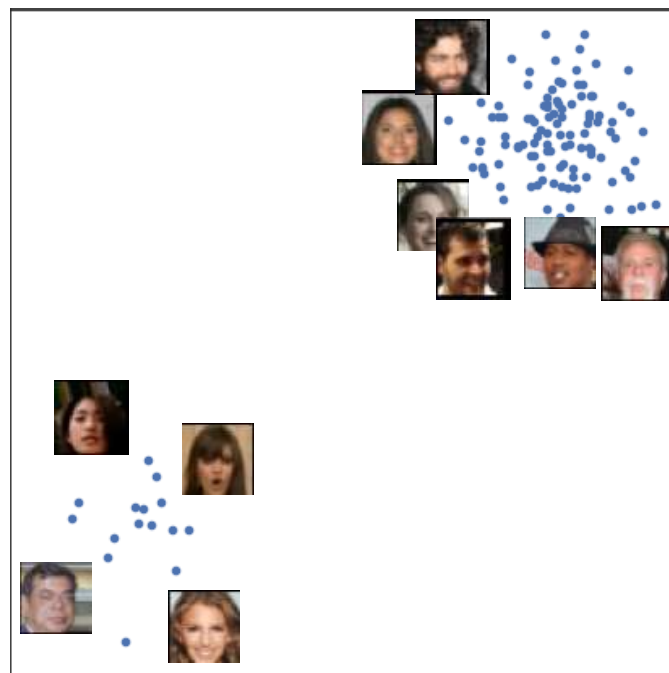
- Instead of parametrizing p , we can parametrize the score

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



*only visualize directions

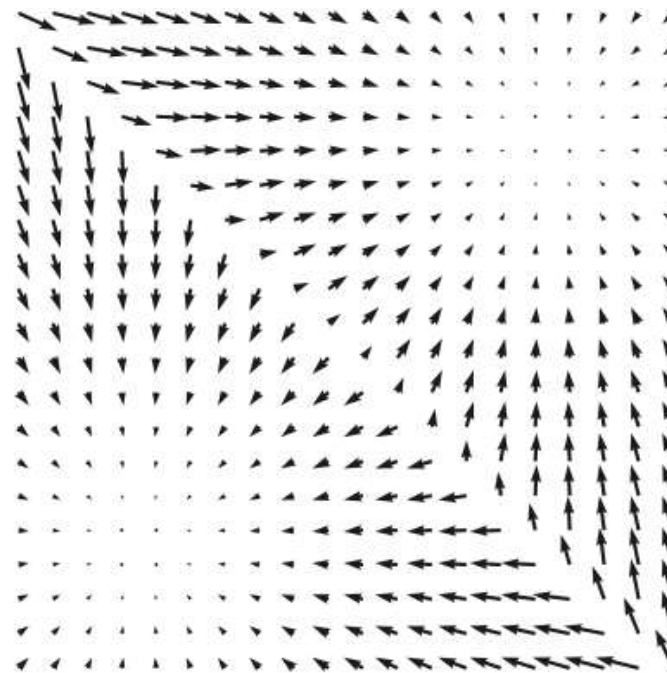
Score-based generative modeling



Data samples

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

score
matching



Scores

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

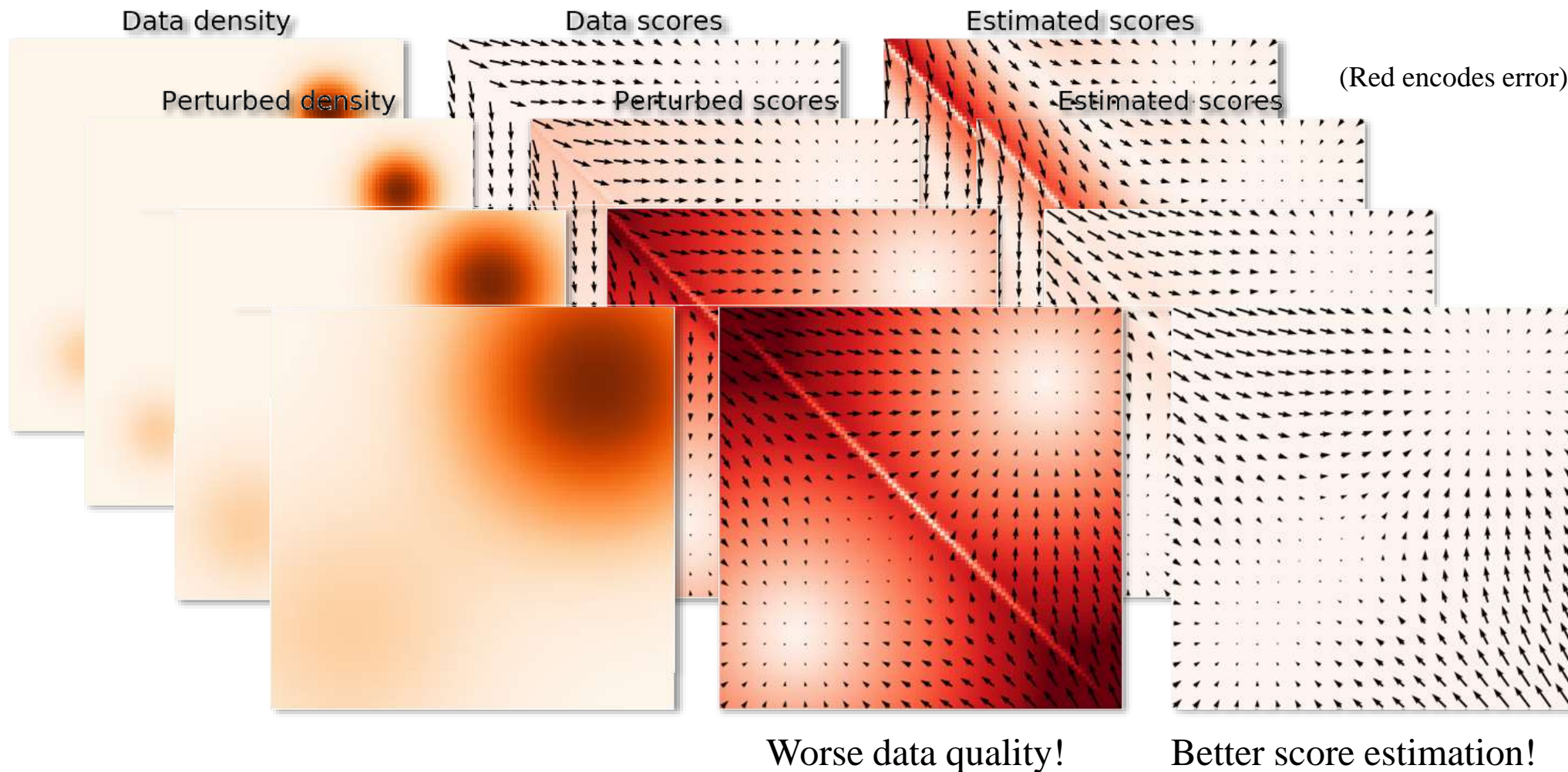


Langevin
dynamics



New samples

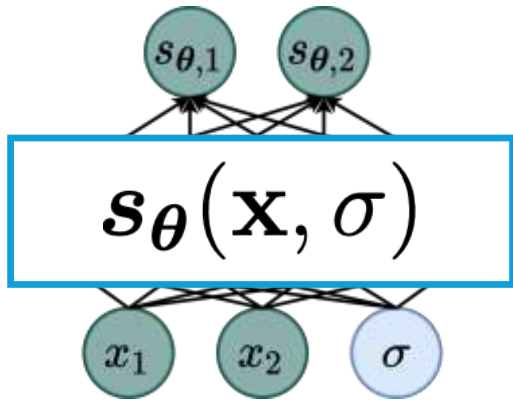
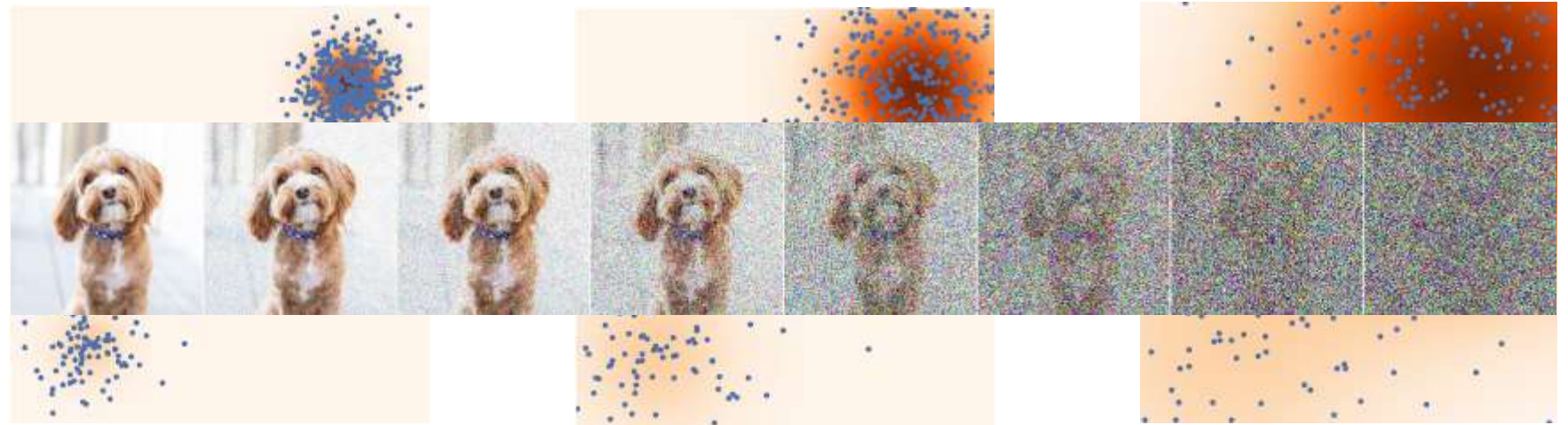
Trading off Data Quality and Estimation Accuracy



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

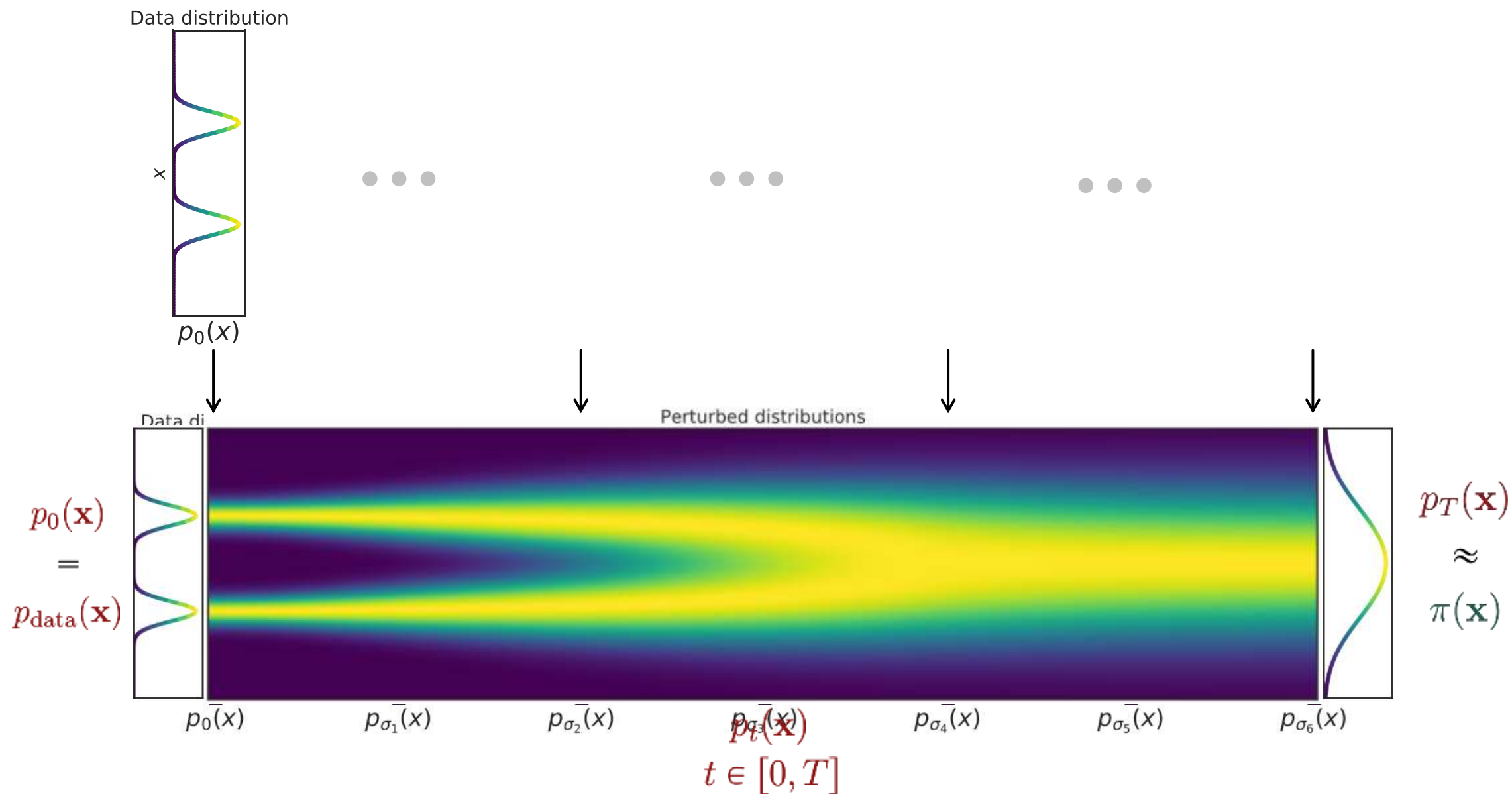
$$\frac{1}{N} \sum_{i=1}^N \lambda(\sigma_i) \mathbb{E}_{p_{\sigma_i}(\mathbf{x})} \left[\left\| \nabla_{\mathbf{x}} \log p_{\sigma_i}(\mathbf{x}) - \mathbf{s}_{\theta}(\mathbf{x}, \sigma_i) \right\|_2^2 \right]$$

Annealed Langevin dynamics

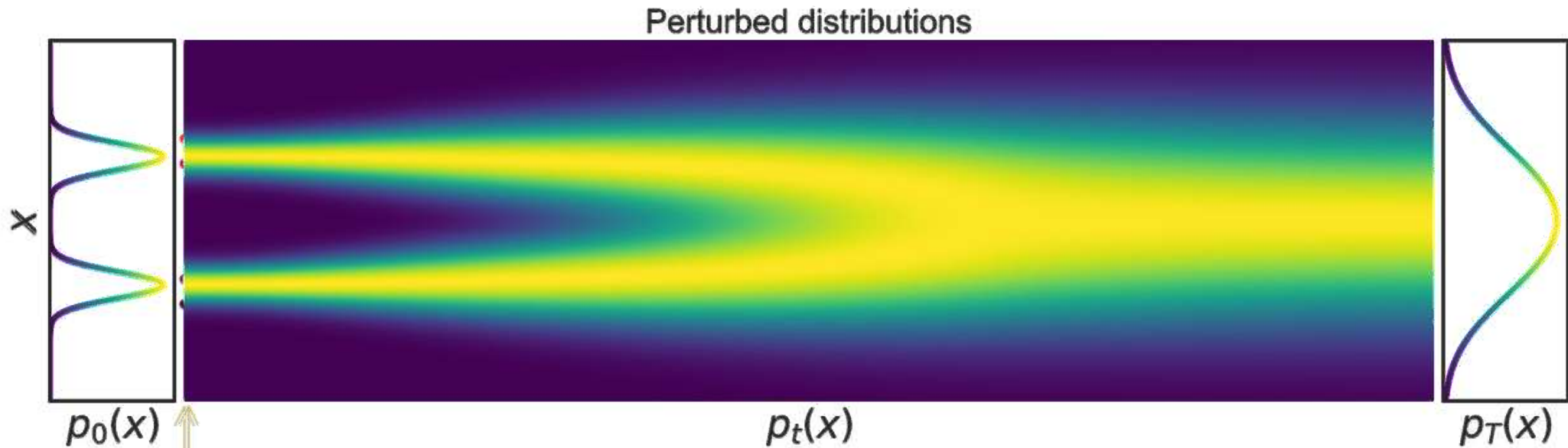
Score matching loss

$\mathbf{s}_{\theta}(\mathbf{x}, \sigma_1)$ $\mathbf{s}_{\theta}(\mathbf{x}, \sigma_2)$ $\mathbf{s}_{\theta}(\mathbf{x}, \sigma_3)$

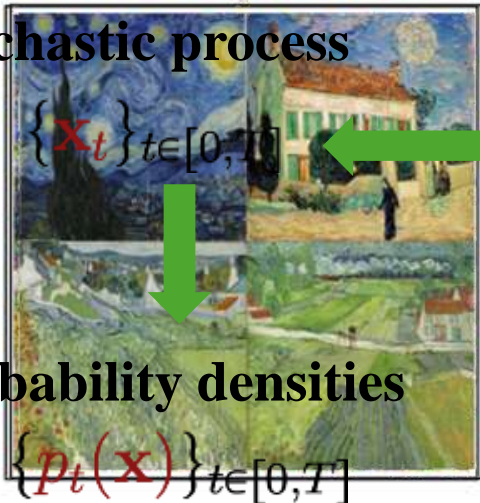
Infinite noise levels



Perturbing data with stochastic processes



Stochastic process



Probability densities

Stochastic differential equation (SDE)

$$dx_t = \underbrace{f(x_t, t)}_{\text{Deterministic drift}} dt + g(t) \underbrace{dw_t}_{\text{Infinitesimal noise}}$$

Deterministic drift

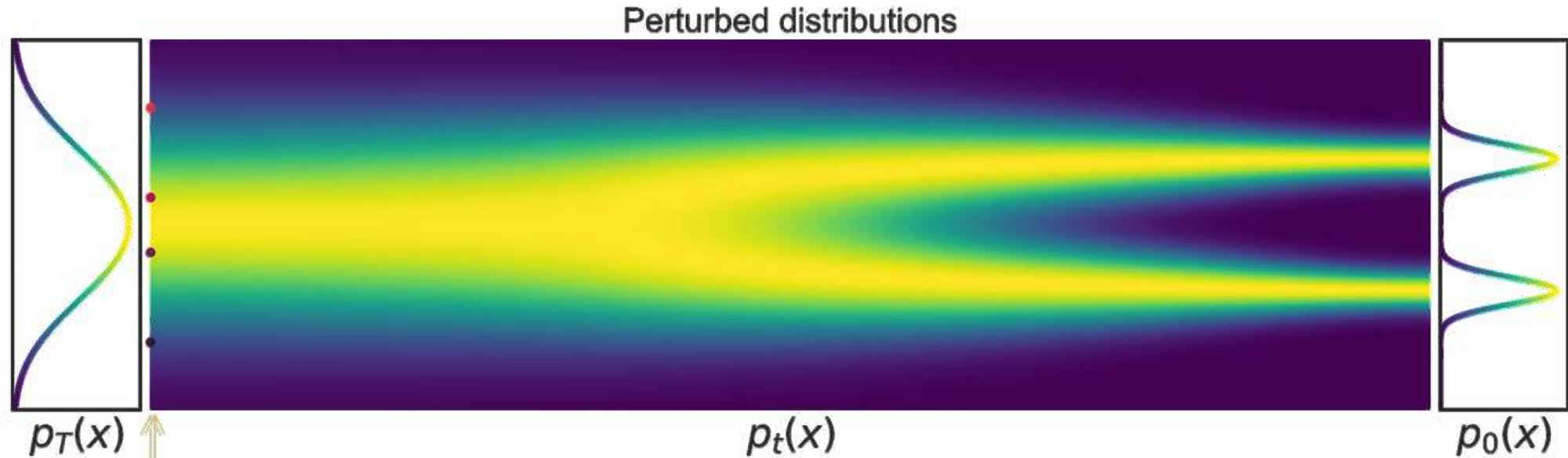
Infinitesimal noise

WLOG: Toy SDE

$$dx_t = \sigma(t) dw_t$$

$$p_T(\mathbf{x}) \approx \pi(\mathbf{x})$$

Generation via reverse stochastic processes



$\pi(\mathbf{x})$
 \approx
 $p_T(\mathbf{x})$

Forward SDE (t: 0 → T)

$$d\mathbf{x}_t = \sigma(t) d\mathbf{w}_t$$
Reverse SDE (t: T → 0)

$$d\mathbf{x}_t = -\sigma(t)^2 \nabla_{\mathbf{x}} \log p_t(\mathbf{x}_t) dt + \sigma(t) d\bar{\mathbf{w}}_t$$

Infinitesimal noise in
the reverse time
direction

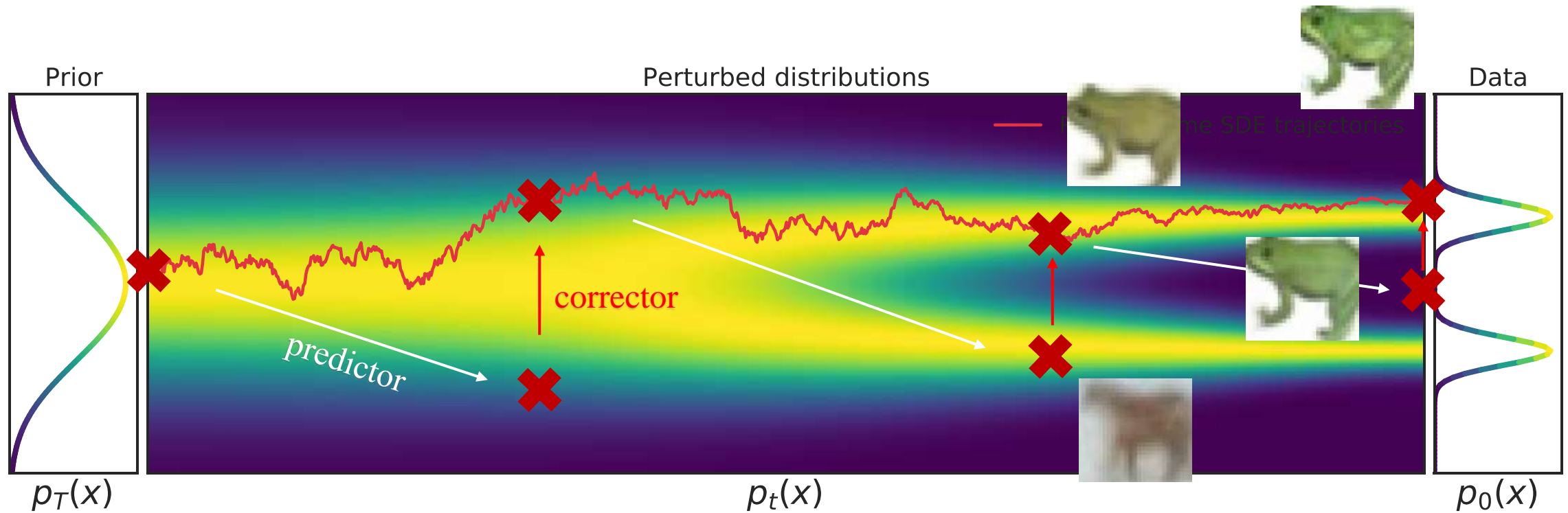
↑

↓

Score function!

Predictor-Corrector sampling methods

- Predictor-Corrector sampling.
 - **Predictor:** Numerical SDE solver
 - **Corrector:** Score-based MCMC



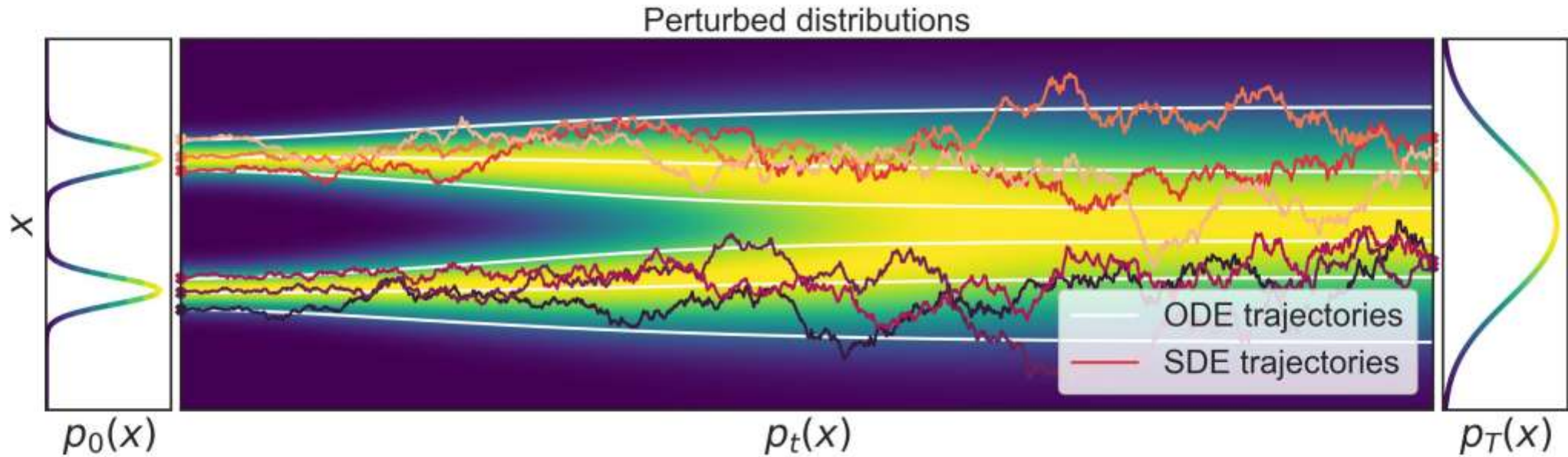
Predictor-Corrector sampling methods

- Predictor-Corrector sampling.
 - **Predictor:** Numerical SDE solver
 - **Corrector:** Score-based MCMC

Algorithm 3 PC sampling (VP SDE)

```
1:  $\mathbf{x}_N \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
2: for  $i = N - 1$  to 0 do
3:    $\mathbf{x}'_i \leftarrow (2 - \sqrt{1 - \beta_{i+1}})\mathbf{x}_{i+1} + \beta_{i+1}\mathbf{s}_{\theta^*}(\mathbf{x}_{i+1}, i + 1)$ 
4:    $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
5:    $\mathbf{x}_i \leftarrow \mathbf{x}'_i + \sqrt{\beta_{i+1}}\mathbf{z}$  Predictor
6:   for  $j = 1$  to  $M$  do Corrector
7:      $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
8:      $\mathbf{x}_i \leftarrow \mathbf{x}_i + \epsilon_i\mathbf{s}_{\theta^*}(\mathbf{x}_i, i) + \sqrt{2\epsilon_i}\mathbf{z}$ 
9: return  $\mathbf{x}_0$ 
```

Converting the SDE to an ODE



SDE

$$d\mathbf{x}_t = \sigma(t) d\mathbf{w}_t$$

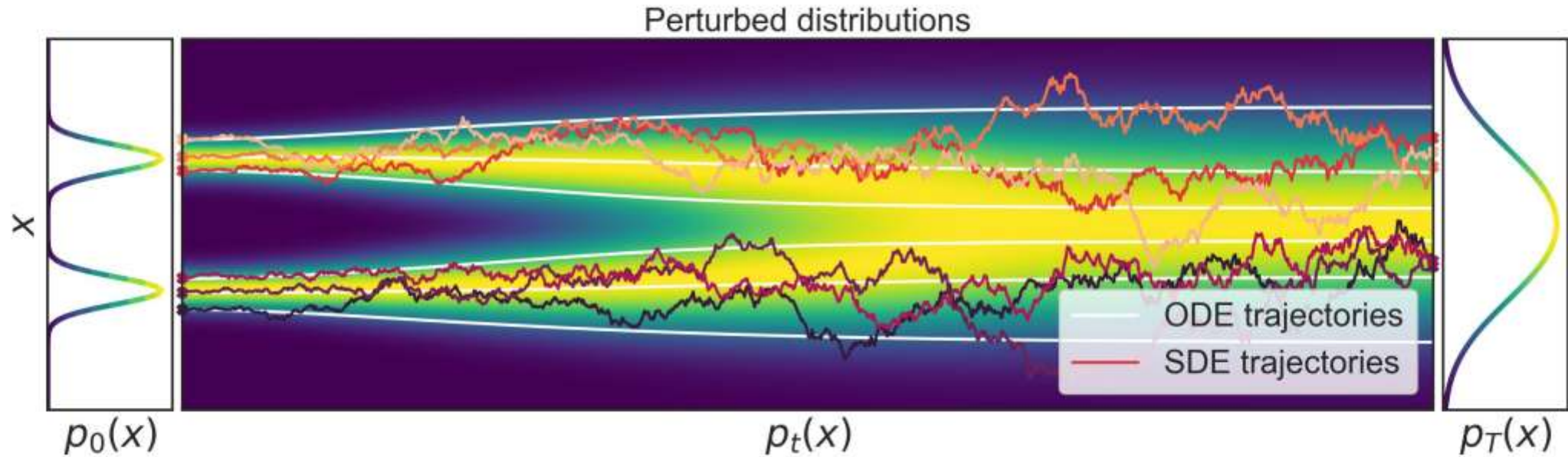


Ordinary differential equation (ODE)

$$\frac{d\mathbf{x}_t}{dt} = -\frac{1}{2}\sigma(t)^2 \nabla_{\mathbf{x}} \log p_t(\mathbf{x}_t)$$

Score function
 $\approx s_{\theta}(\mathbf{x}, t)$

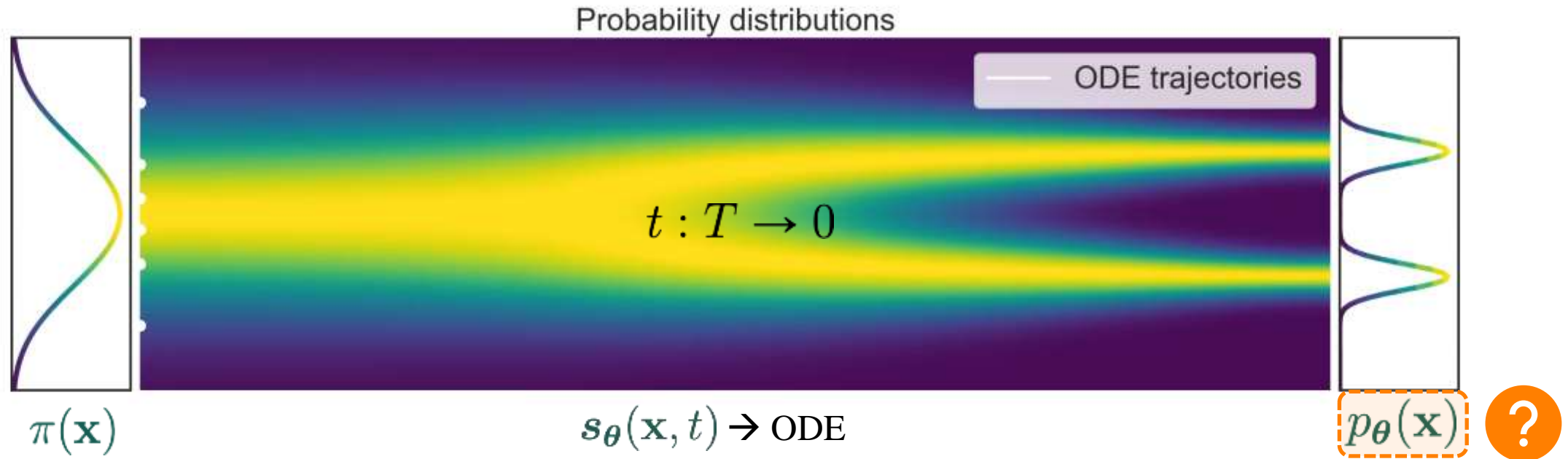
Converting the SDE to an ODE



We can think of this as a (continuous time, infinite depth) normalizing flow

1. Unique ODE solution \rightarrow Invertible mapping
2. To invert, solve ODE backwards from T to 0

Evaluating the probabilities with ODEs



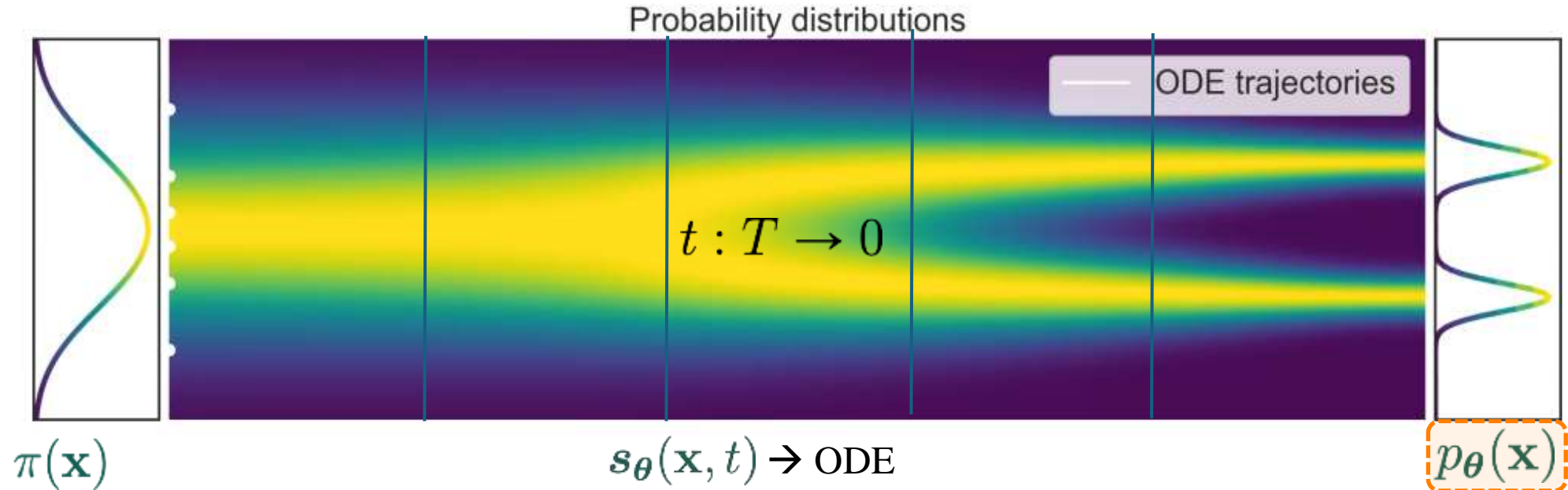
Computing the probability density function (change of variables formula)

$$\log p_{\theta}(\mathbf{x}_0) = \log \pi(\mathbf{x}_T) - \frac{1}{2} \int_0^T \sigma(t)^2 \text{trace}(\nabla_{\mathbf{x}} s_{\theta}(\mathbf{x}, t)) dt$$

ODE solver
Computed in polynomial time

- **It is a (continuous-time) normalizing flow model!**

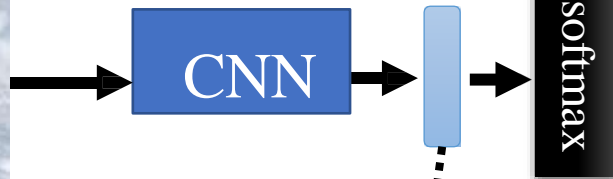
Accelerated sampling



- Numerical methods + ODE formulation to accelerate sampling
- DDIM [Song and Ermon, 2021]:
 - Coarsely discretize the time axis, take big steps
 - Corresponds to exponential integrator (semi-linear ODE) [Lu et al, 2022; Zhang and Chen, 2022]
 - 10x-50x speedups, comparable sample quality

Fréchet Inception Distance (FID)

<https://arxiv.org/abs/1706.08500>



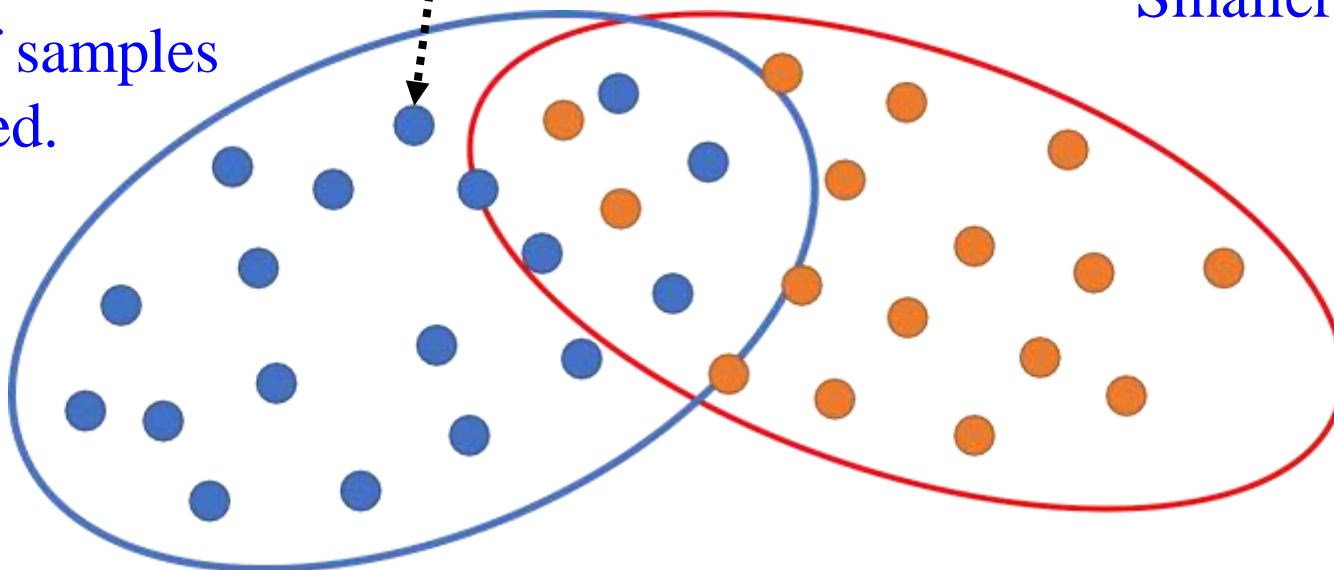
red points: real images

blue points: generated images

FID = Fréchet distance
between the two **Gaussians**

Smaller is better

A lot of samples
is needed.



Accelerated sampling

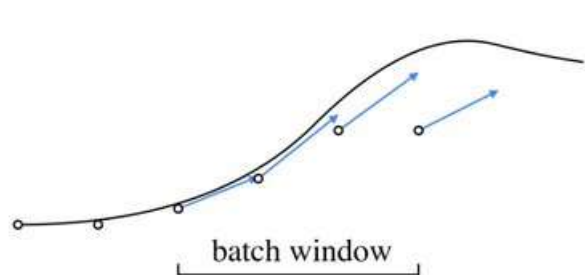
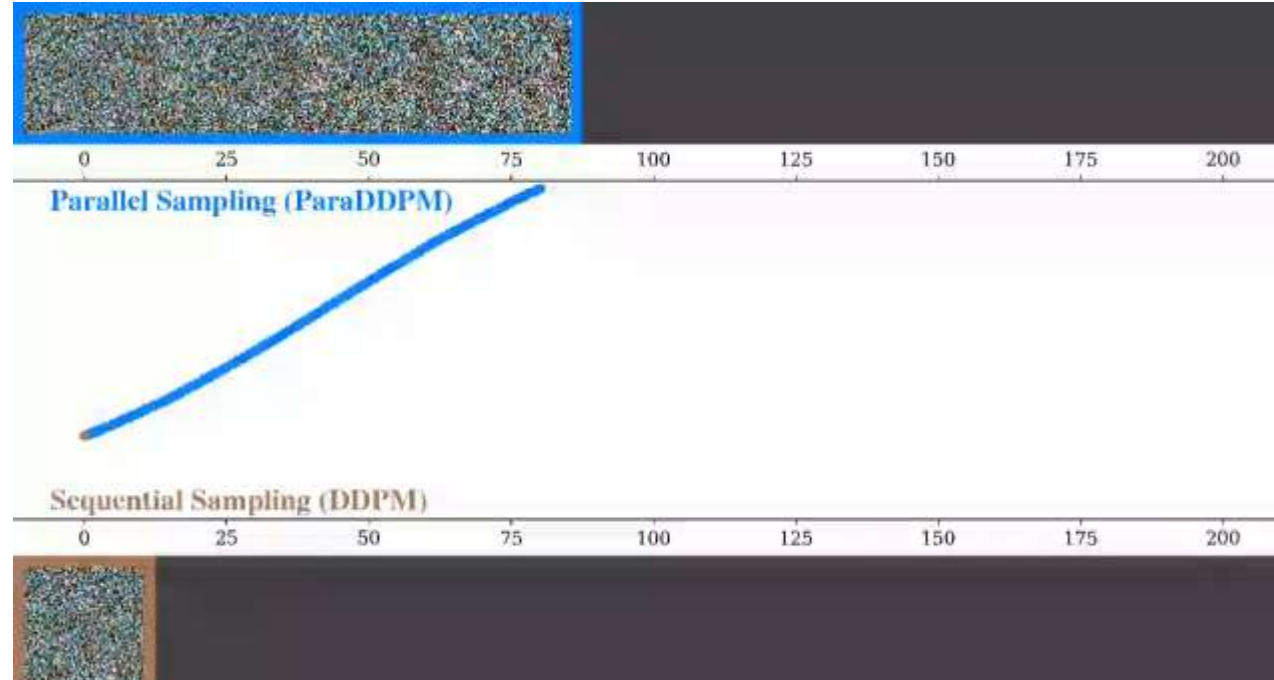
S	CIFAR10 (32×32)					CelebA (64×64)					
	10	20	50	100	1000	10	20	50	100	1000	
η	0.0	13.36	6.84	4.67	4.16	4.04	17.33	13.73	9.17	6.53	3.51
	0.2	14.04	7.11	4.77	4.25	4.09	17.66	14.11	9.51	6.79	3.64
	0.5	16.66	8.35	5.25	4.46	4.29	19.86	16.06	11.01	8.09	4.28
	1.0	41.07	18.36	8.01	5.78	4.73	33.12	26.03	18.48	13.93	5.98
$\hat{\sigma}$	367.43	133.37	32.72	9.99	3.17	299.71	183.83	71.71	45.20	3.26	

- Numerical methods + ODE formulation to accelerate sampling
- DDIM [Song and Ermon, 2021]:
 - Coarsely discretize the time axis, take big steps
 - Corresponds to exponential integrator (semi-linear ODE) [Lu et al, 2022; Zhang and Chen, 2022]
 - 10x-50x speedups, comparable sample quality

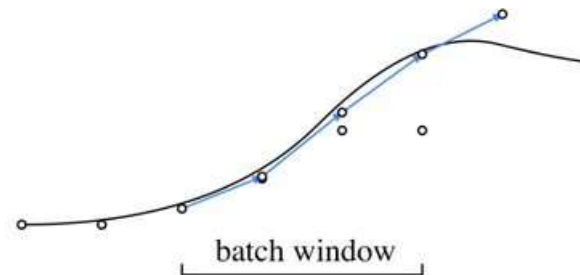
Compared to DDPM, DDIM is able to:

- Generate higher-quality samples using a much fewer number of steps.
- Have “consistency” property since the generative process is deterministic, meaning that multiple samples conditioned on the same latent variable should have similar high-level features.
- Because of the consistency, DDIM can do semantically meaningful interpolation in the latent variable.

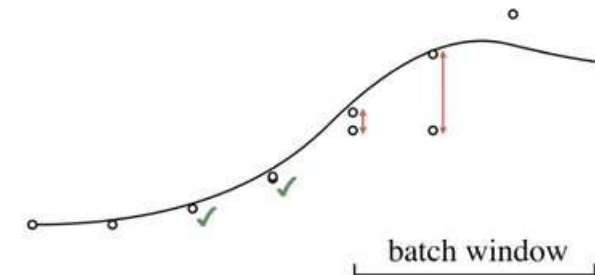
Parallel ODE solving



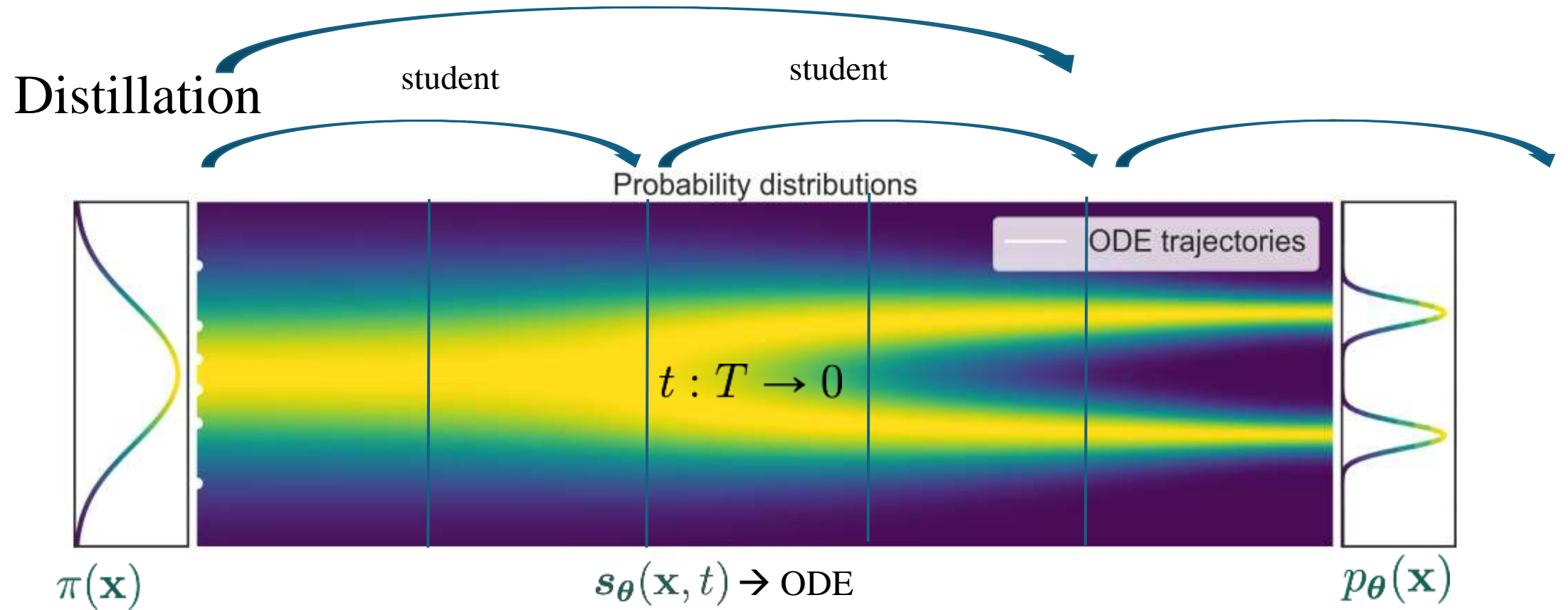
(a) Compute the drift of $x_{t:t+p}^k$ on a batch window of size $p = 4$, in parallel



(b) Update the values to $x_{t:t+p}^{k+1}$ using the cumulative drift of points in the window



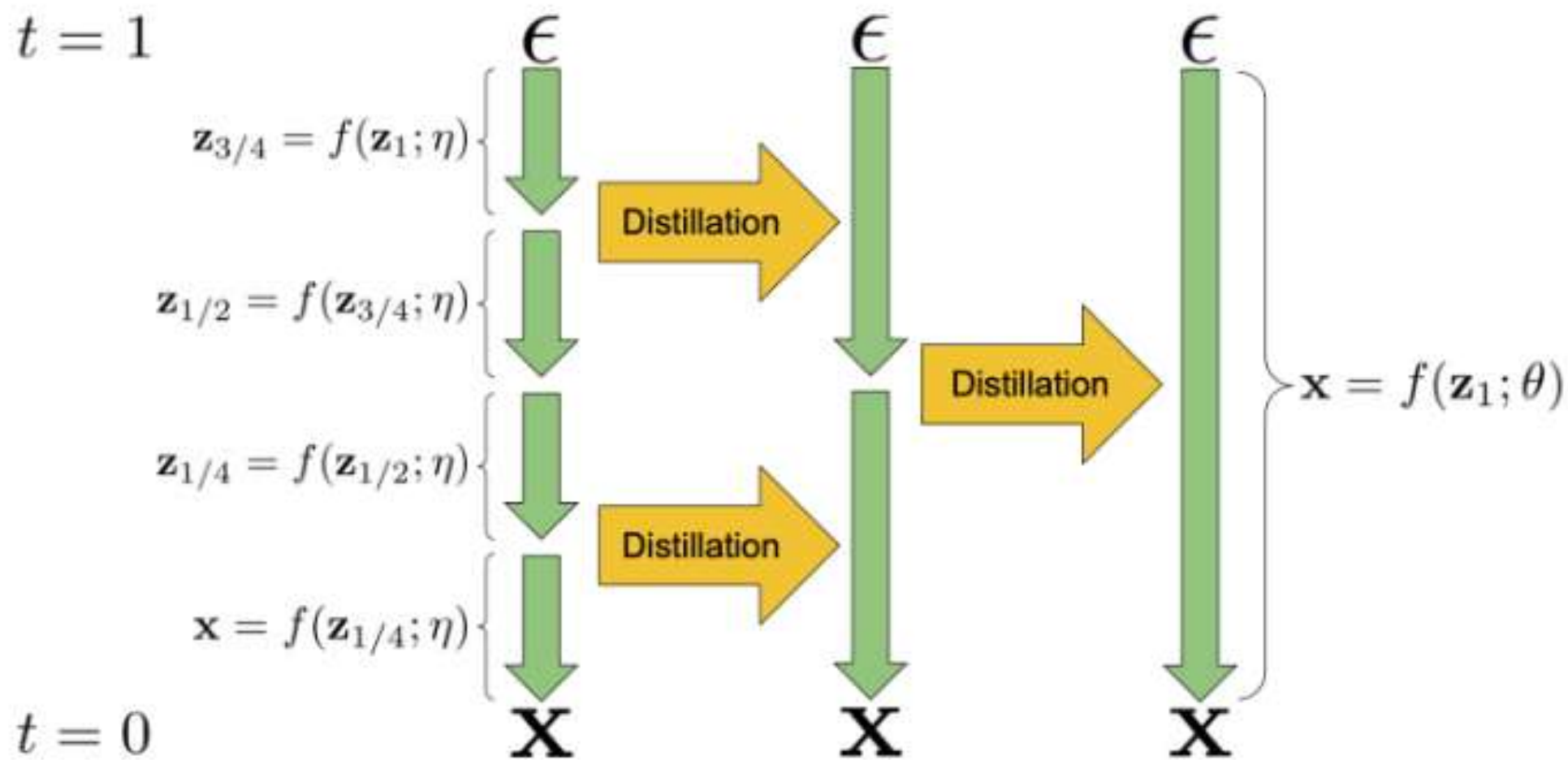
(c) Determine how far to slide the window forward, based on the error $\|x_i^{k+1} - x_i^k\|^2$.



- Progressive distillation [Salimans, Ho 2022]

- DDIM sampler as a teacher model
- Student model trained to do in 1 step what DDIM achieves in 2 steps
- Applied recursively to drastically reduce the number of steps required

Distillation



Algorithm 1 Standard diffusion training

Require: Model $\hat{\mathbf{x}}_\theta(\mathbf{z}_t)$ to be trained

Require: Data set \mathcal{D}

Require: Loss weight function $w()$

while not converged **do**

$\mathbf{x} \sim \mathcal{D}$ \triangleright Sample data

$t \sim U[0, 1]$ \triangleright Sample time

$\epsilon \sim N(0, I)$ \triangleright Sample noise

$\mathbf{z}_t = \alpha_t \mathbf{x} + \sigma_t \epsilon$ \triangleright Add noise to data

$\tilde{\mathbf{x}} = \mathbf{x}$ \triangleright Clean data is target for $\hat{\mathbf{x}}$

$\lambda_t = \log[\alpha_t^2/\sigma_t^2]$ \triangleright log-SNR

$L_\theta = w(\lambda_t) \|\tilde{\mathbf{x}} - \hat{\mathbf{x}}_\theta(\mathbf{z}_t)\|_2^2$ \triangleright Loss

$\theta \leftarrow \theta - \gamma \nabla_\theta L_\theta$ \triangleright Optimization

end while

Algorithm 2 Progressive distillation

Require: Trained teacher model $\hat{\mathbf{x}}_\eta(\mathbf{z}_t)$

Require: Data set \mathcal{D}

Require: Loss weight function $w()$

Require: Student sampling steps N

for K iterations **do**

$\theta \leftarrow \eta$ \triangleright Init student from teacher

while not converged **do**

$\mathbf{x} \sim \mathcal{D}$

$t = i/N, i \sim \text{Cat}[1, 2, \dots, N]$

$\epsilon \sim N(0, I)$

$\mathbf{z}_t = \alpha_t \mathbf{x} + \sigma_t \epsilon$

2 steps of DDIM with teacher

$t' = t - 0.5/N, t'' = t - 1/N$

$\mathbf{z}_{t'} = \alpha_{t'} \hat{\mathbf{x}}_\eta(\mathbf{z}_t) + \frac{\sigma_{t'}}{\sigma_t} (\mathbf{z}_t - \alpha_t \hat{\mathbf{x}}_\eta(\mathbf{z}_t))$

$\mathbf{z}_{t''} = \alpha_{t''} \hat{\mathbf{x}}_\eta(\mathbf{z}_{t'}) + \frac{\sigma_{t''}}{\sigma_{t'}} (\mathbf{z}_{t'} - \alpha_{t'} \hat{\mathbf{x}}_\eta(\mathbf{z}_{t'}))$

$\tilde{\mathbf{x}} = \frac{\mathbf{z}_{t''} - (\sigma_{t''}/\sigma_t) \mathbf{z}_t}{\alpha_{t''} - (\sigma_{t''}/\sigma_t) \alpha_t}$ \triangleright Teacher $\hat{\mathbf{x}}$ target

$\lambda_t = \log[\alpha_t^2/\sigma_t^2]$

$L_\theta = w(\lambda_t) \|\tilde{\mathbf{x}} - \hat{\mathbf{x}}_\theta(\mathbf{z}_t)\|_2^2$

$\theta \leftarrow \theta - \gamma \nabla_\theta L_\theta$

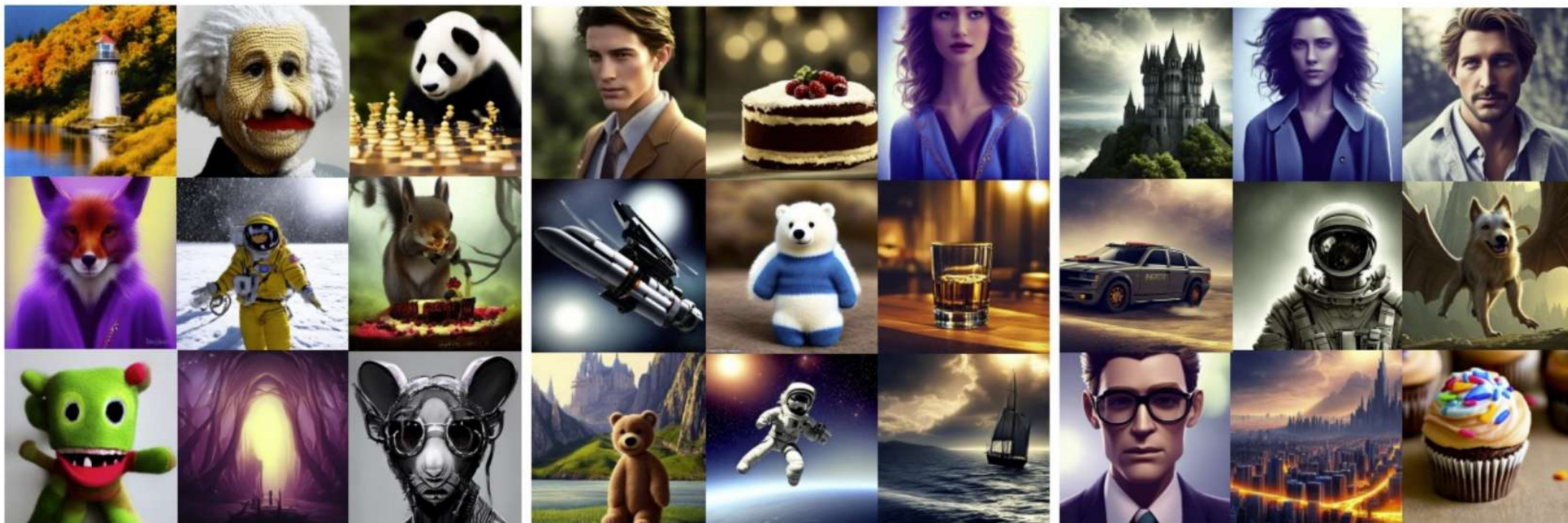
end while

$\eta \leftarrow \theta$ \triangleright Student becomes next teacher

$N \leftarrow N/2$ \triangleright Halve number of sampling steps

end for

On distillation of guided diffusion models

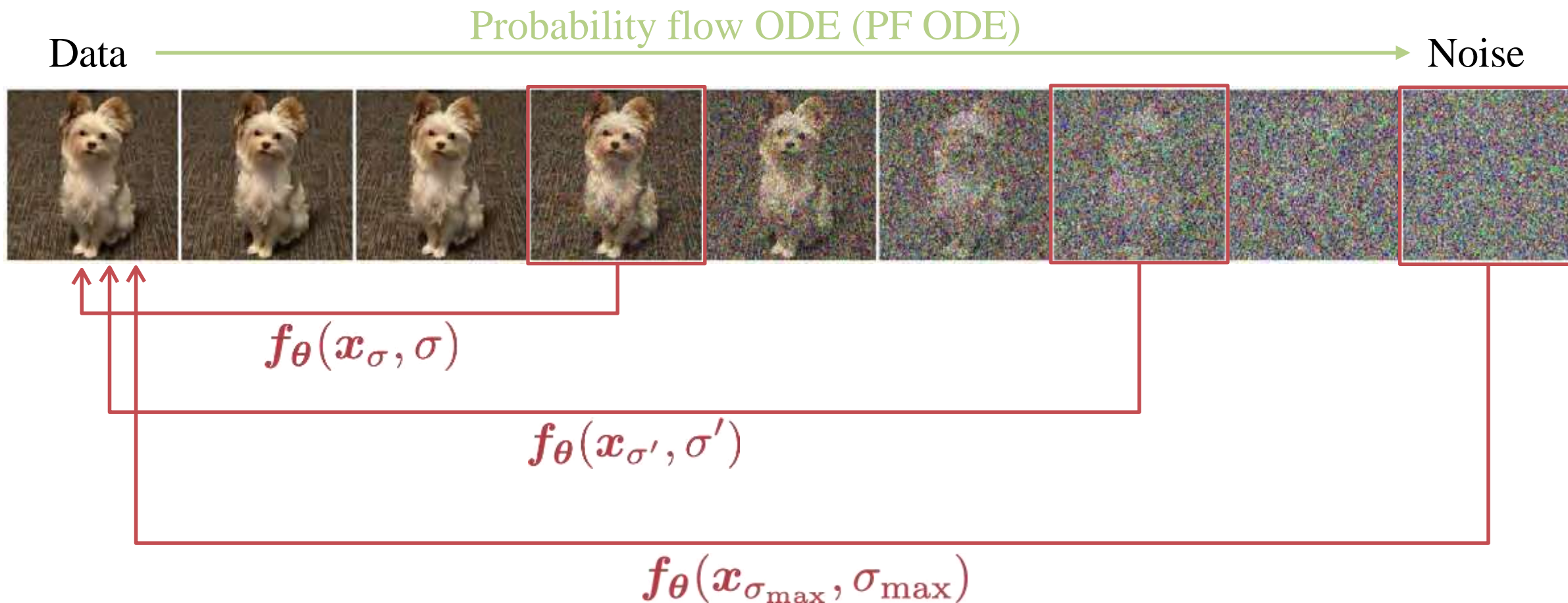


(a) 2 denoising steps

(b) 4 denoising steps

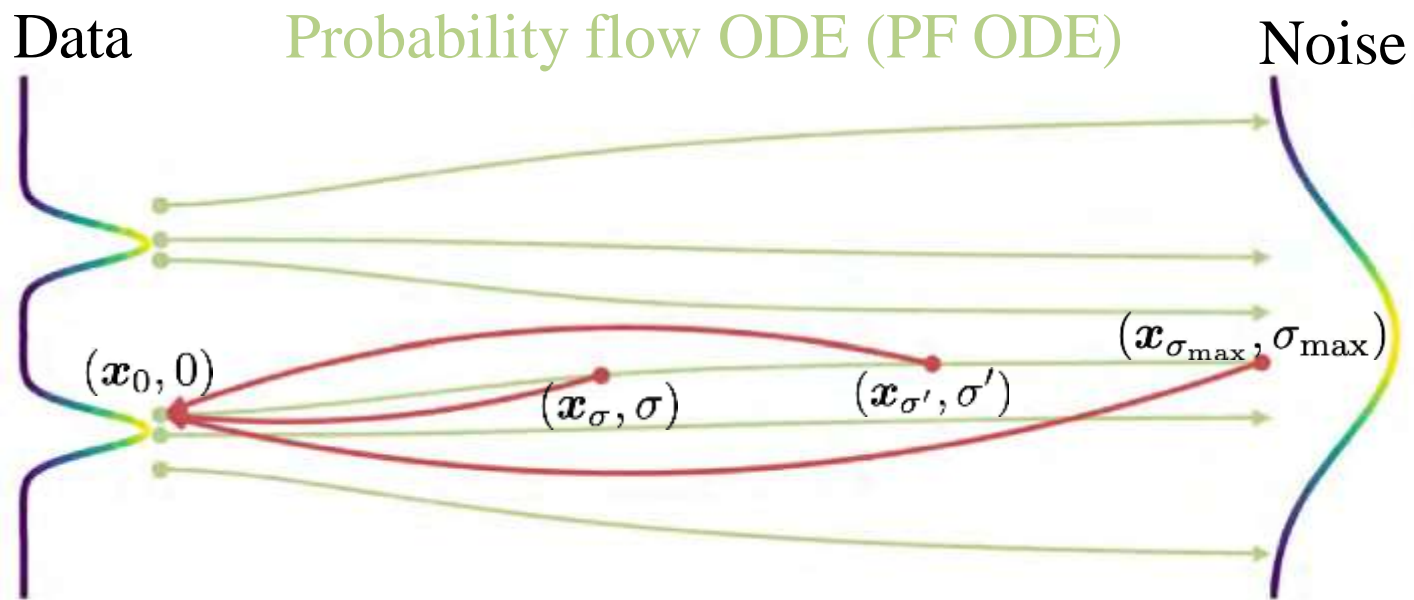
(c) 8 denoising steps

Consistency models are designed for one-step generation



How does this differ from a denoiser?

Consistency models are designed for one-step generation



Consistency models are trained to map points on any trajectory of the PF ODE to the trajectory's origin in one step.

$$f_{\theta}(\mathbf{x}_{\sigma}, \sigma) = \mathbf{x}_0$$

Boundary condition

$$f_{\theta}(\mathbf{x}_0, 0) = \mathbf{x}_0$$

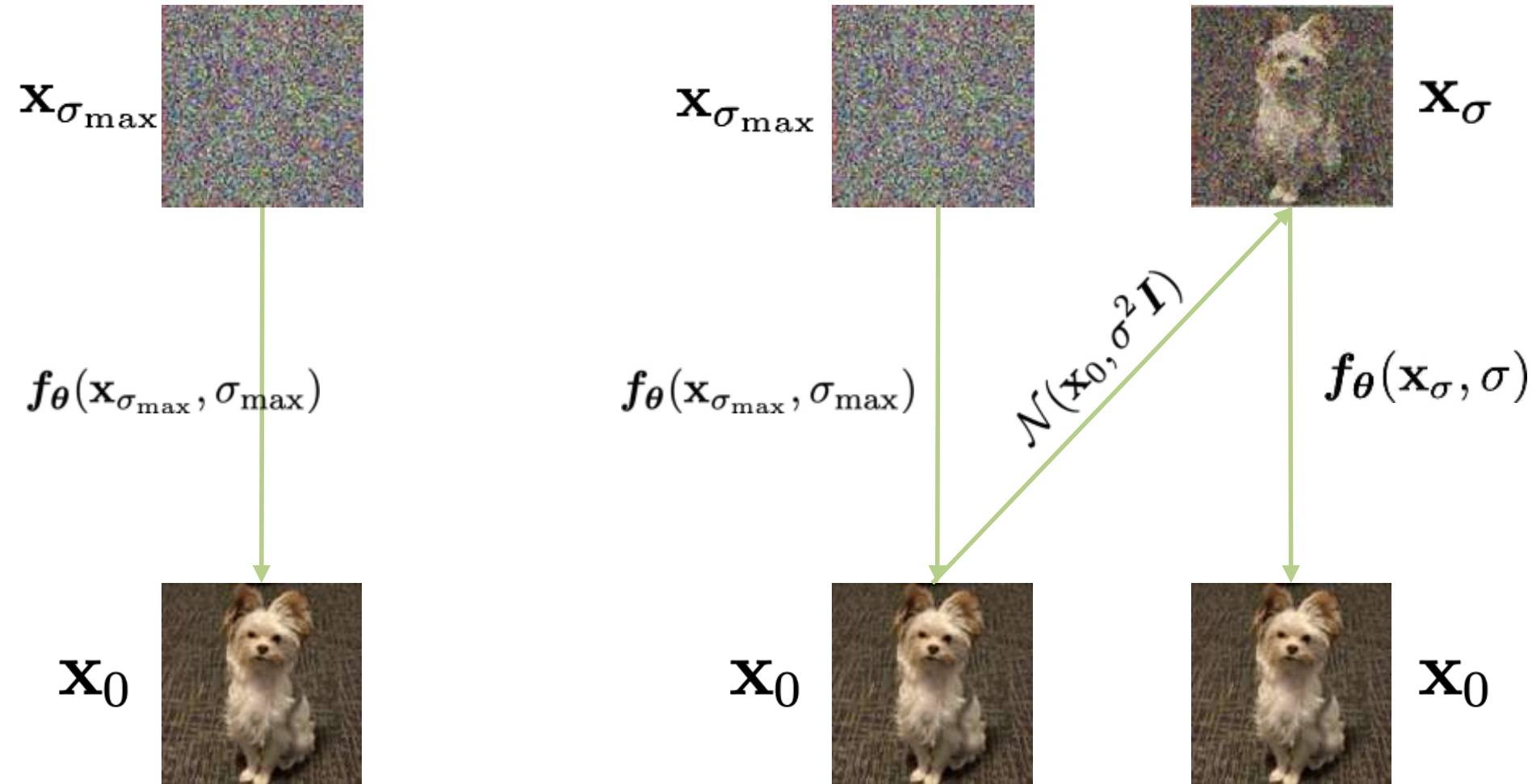
Enforced via network parameterization

Enforced via learning

Self-consistency

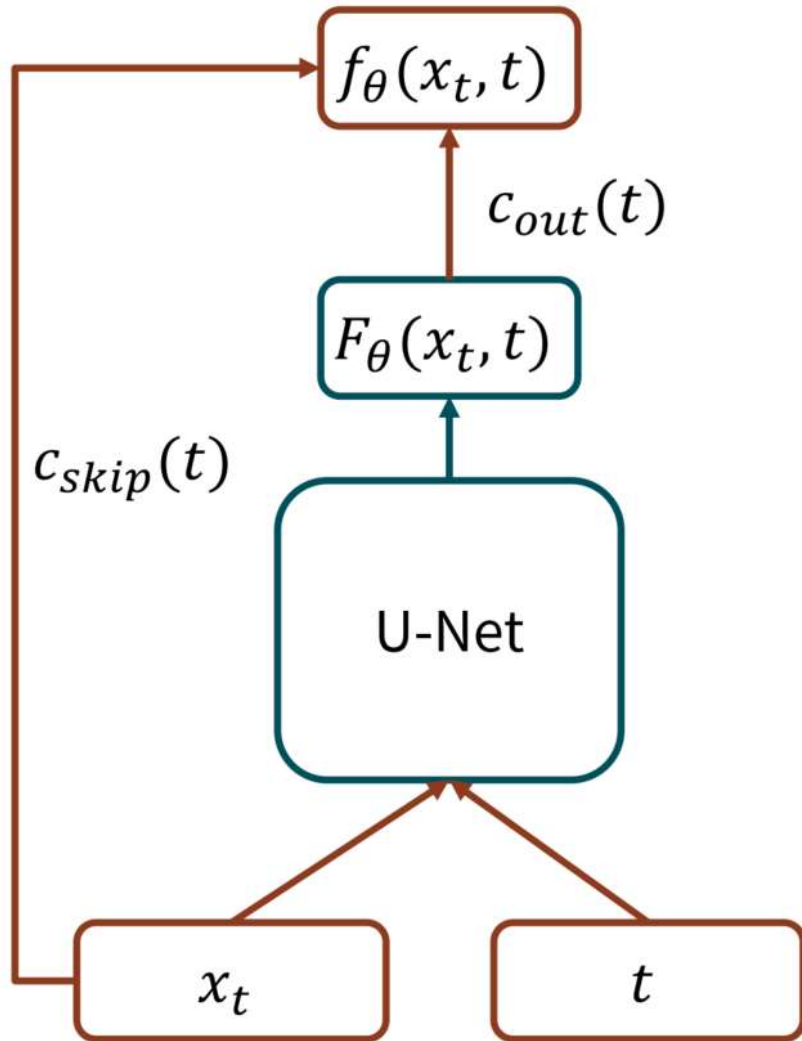
$$\forall \sigma, \sigma' \in [0, \sigma_{\max}] : f_{\theta}(\mathbf{x}_{\sigma}, \sigma) = f_{\theta}(\mathbf{x}_{\sigma'}, \sigma')$$

Sampling from consistency models



- Trading compute for quality
- Zero-shot image editing

Enforcing the boundary condition



- Skip connections for enforcing the boundary condition:

$$f_\theta(\mathbf{x}, t) = c_{\text{skip}}(t)\mathbf{x} + c_{\text{out}}(t)F_\theta(\mathbf{x}, t)$$

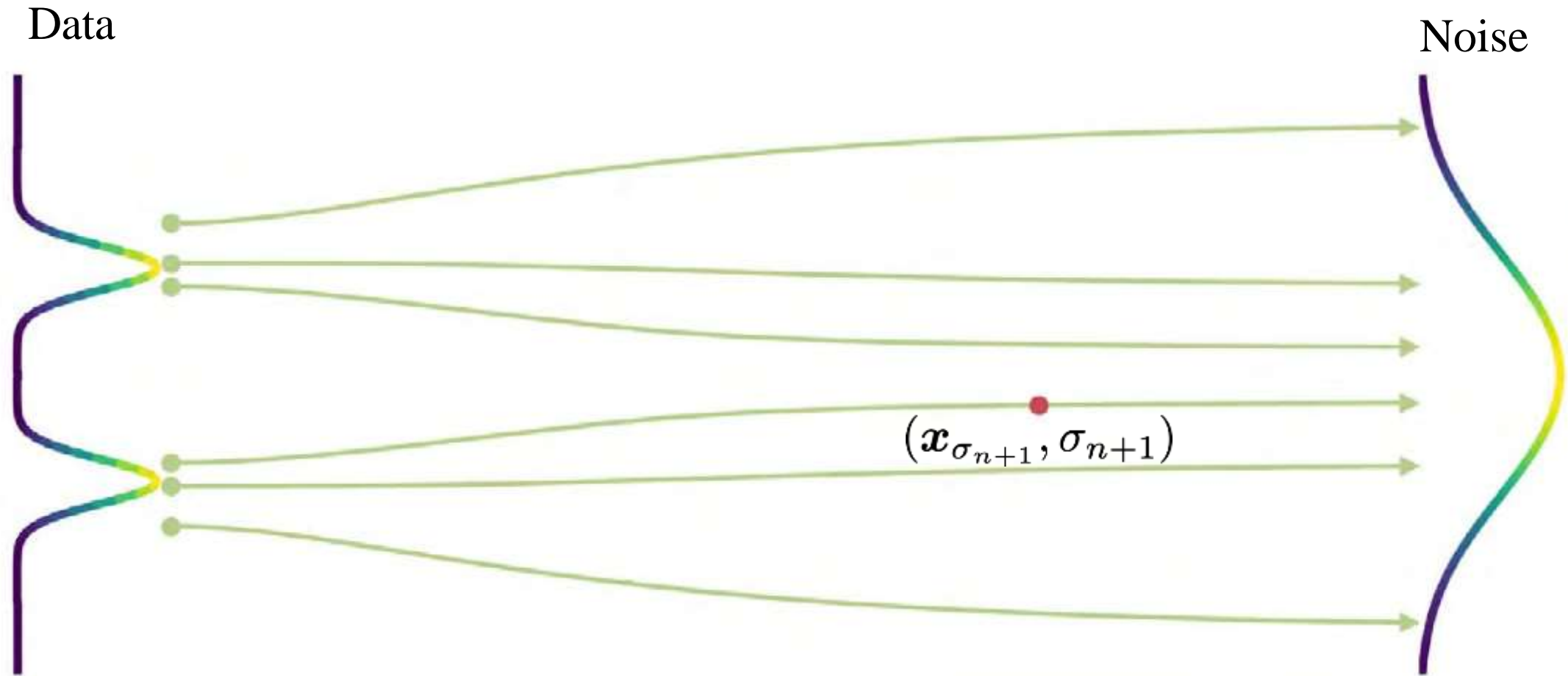
$$c_{\text{skip}}(0) = 1$$

$$c_{\text{out}}(0) = 0$$

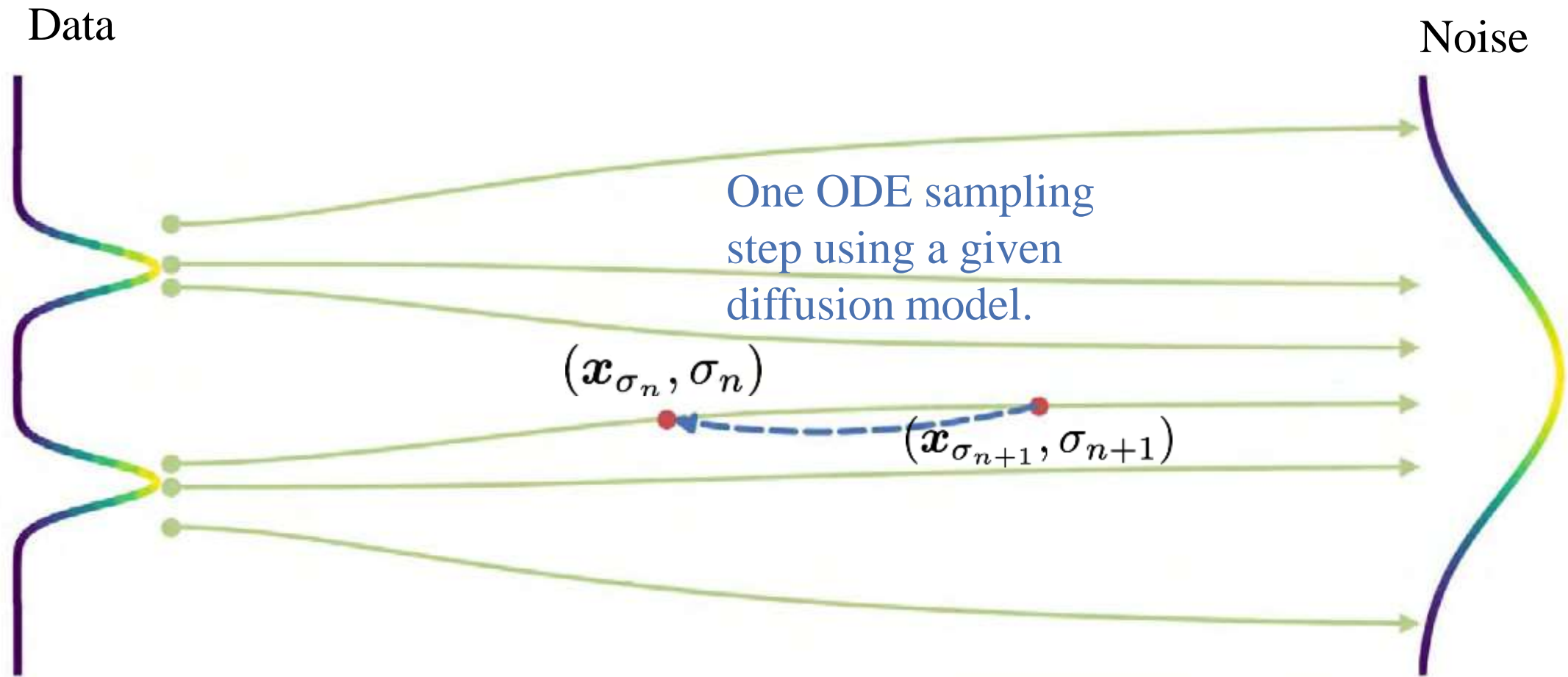
- The denoising/score network in diffusion models often has a similar parameterization (cf., EDM, v-prediction, etc.)

Karras, Tero, et al. "Elucidating the design space of diffusion-based generative models." *arXiv preprint arXiv:2206.00364* (2022).

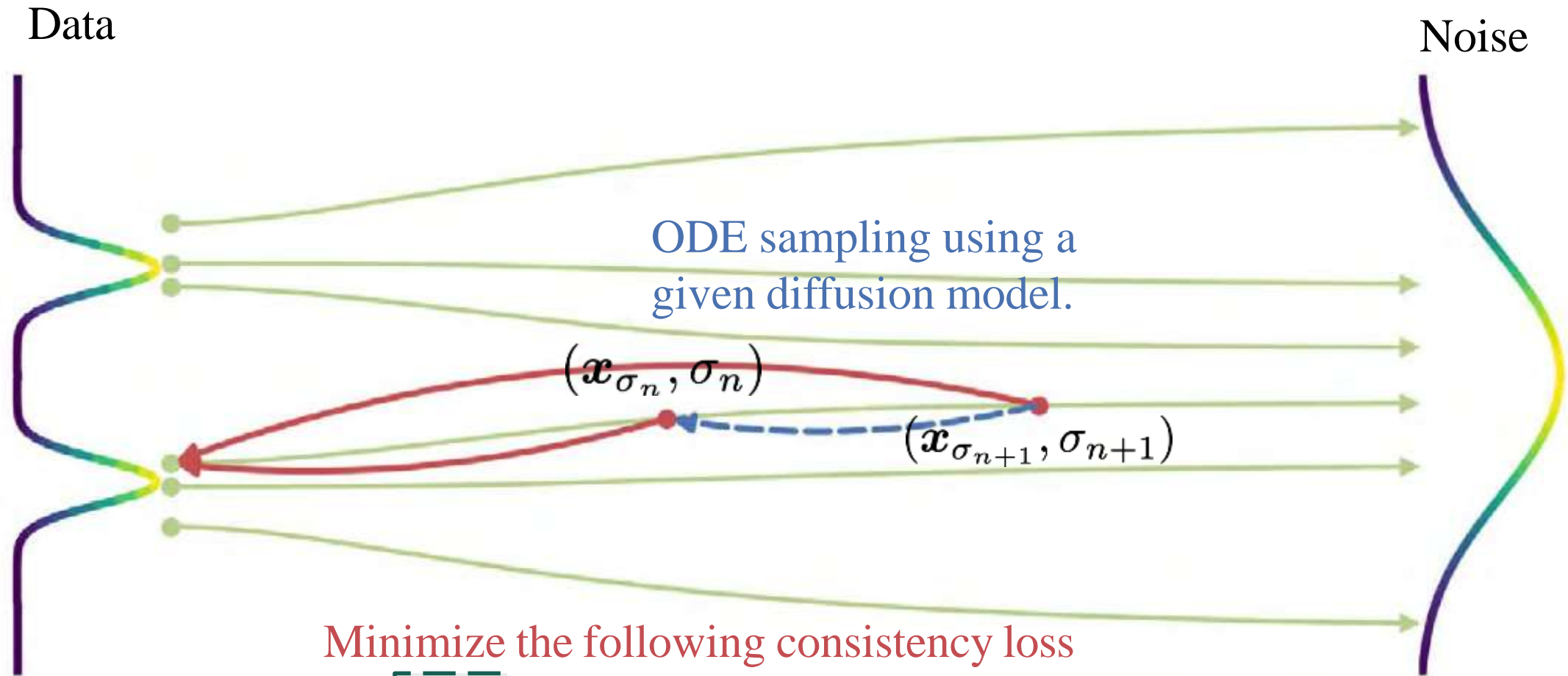
Training consistency models via distillation



Training consistency models via distillation



Training consistency models via distillation



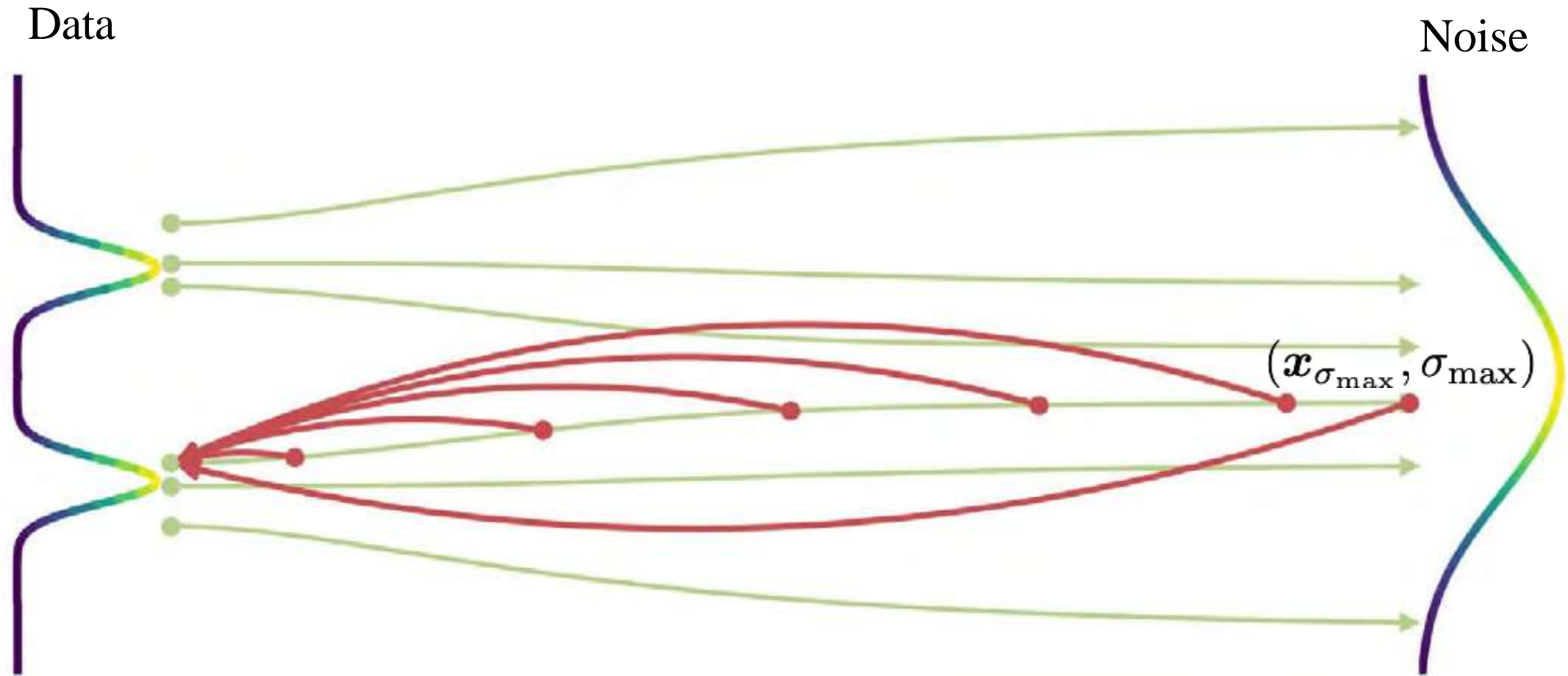
Minimize the following consistency loss

$$\min_{\theta} \underbrace{\lambda(\sigma_n)}_{\text{Weighting function}} d(\mathbf{f}_{\theta}(\mathbf{x}_{\sigma_{n+1}}, \sigma_{n+1}), \mathbf{f}_{\theta^-}(\mathbf{x}_{\sigma_n}, \sigma_n))$$

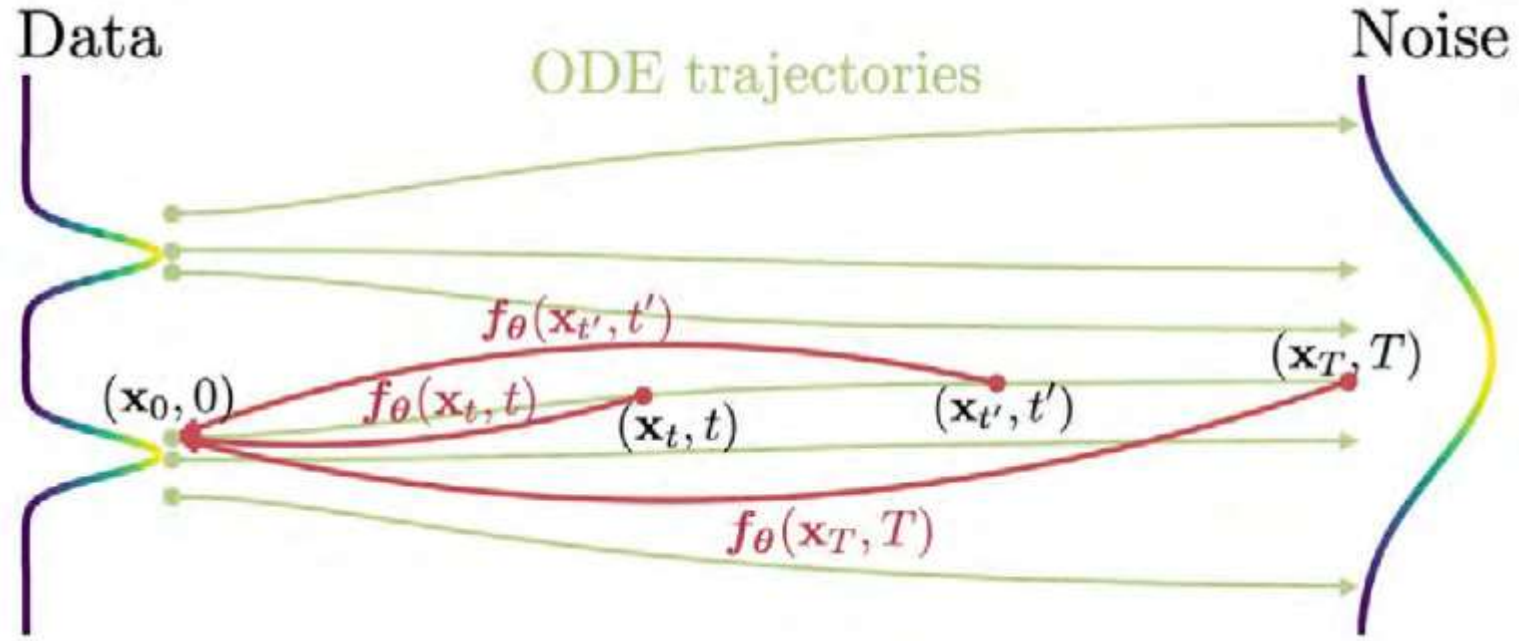
Weighting function

$$\theta^- = \text{EMA}_{\mu}(\theta)$$

Training consistency models via distillation



Enforcing self-consistency via distillation



$$\min_{\theta} \|\mathbf{f}_\theta(\mathbf{x}_{t'}, t') - \mathbf{f}_\theta(\mathbf{x}_t, t)\|_2^2$$

$\mathbf{x}_0 \leftarrow \mathbf{x}_t \leftarrow \mathbf{x}_{t'} \leftarrow \mathbf{x}_T$

ODE solver
+ pretrained
score function

ODE solver
+ pretrained
score function

ODE solver
+ pretrained
score function

$$\min_{\theta} \|\mathbf{f}_\theta(\mathbf{x}_t, t) - \mathbf{f}_\theta(\mathbf{x}_0, 0)\|_2^2$$

$$\min_{\theta} \|\mathbf{f}_\theta(\mathbf{x}_T, T) - \mathbf{f}_\theta(\mathbf{x}_{t'}, t')\|_2^2$$

Consistency Model

Algorithm 2 Consistency Distillation (CD)

Input: dataset \mathcal{D} , initial model parameter θ , learning rate η , ODE solver $\Phi(\cdot, \cdot; \phi)$, $d(\cdot, \cdot)$, $\lambda(\cdot)$, and μ

$\theta^- \leftarrow \theta$

repeat

 Sample $\mathbf{x} \sim \mathcal{D}$ and $n \sim \mathcal{U}[[1, N - 1]]$

 Sample $\mathbf{x}_{t_{n+1}} \sim \mathcal{N}(\mathbf{x}; t_{n+1}^2 \mathbf{I})$

$\hat{\mathbf{x}}_{t_n}^\phi \leftarrow \mathbf{x}_{t_{n+1}} + (t_n - t_{n+1})\Phi(\mathbf{x}_{t_{n+1}}, t_{n+1}; \phi)$

$\mathcal{L}(\theta, \theta^-; \phi) \leftarrow$

$\lambda(t_n)d(\mathbf{f}_\theta(\mathbf{x}_{t_{n+1}}, t_{n+1}), \mathbf{f}_{\theta^-}(\hat{\mathbf{x}}_{t_n}^\phi, t_n))$

$\theta \leftarrow \theta - \eta \nabla_\theta \mathcal{L}(\theta, \theta^-; \phi)$

$\theta^- \leftarrow \text{stopgrad}(\mu \theta^- + (1 - \mu)\theta)$

until convergence

Distill a diffusion model into a consistency model by minimizing the difference between model outputs for pairs generated out of the same trajectory. This enables a much cheaper sampling evaluation.

Consistency Model

Algorithm 3 Consistency Training (CT)

Input: dataset \mathcal{D} , initial model parameter θ , learning rate η , step schedule $N(\cdot)$, EMA decay rate schedule $\mu(\cdot)$, $d(\cdot, \cdot)$, and $\lambda(\cdot)$

$\theta^- \leftarrow \theta$ and $k \leftarrow 0$

repeat

 Sample $\mathbf{x} \sim \mathcal{D}$, and $n \sim \mathcal{U}[[1, N(k) - 1]]$

 Sample $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$

$\mathcal{L}(\theta, \theta^-) \leftarrow$

$\lambda(t_n)d(\mathbf{f}_\theta(\mathbf{x} + t_{n+1}\mathbf{z}, t_{n+1}), \mathbf{f}_{\theta^-}(\mathbf{x} + t_n\mathbf{z}, t_n))$

$\theta \leftarrow \theta - \eta \nabla_{\theta} \mathcal{L}(\theta, \theta^-)$

$\theta^- \leftarrow \text{stopgrad}(\mu(k)\theta^- + (1 - \mu(k))\theta)$

$k \leftarrow k + 1$

until convergence

Train a consistency model independently.

Consistency models distilled from diffusion models



EDM
FID = 2.44
NFE = 79



One step
FID = 6.20
NFE = 1



Two steps
FID = 4.70
NFE = 2

Consistency models distilled from diffusion models



EDM
FID = 3.57
NFE = 79



One step
FID = 7.80
NFE = 1

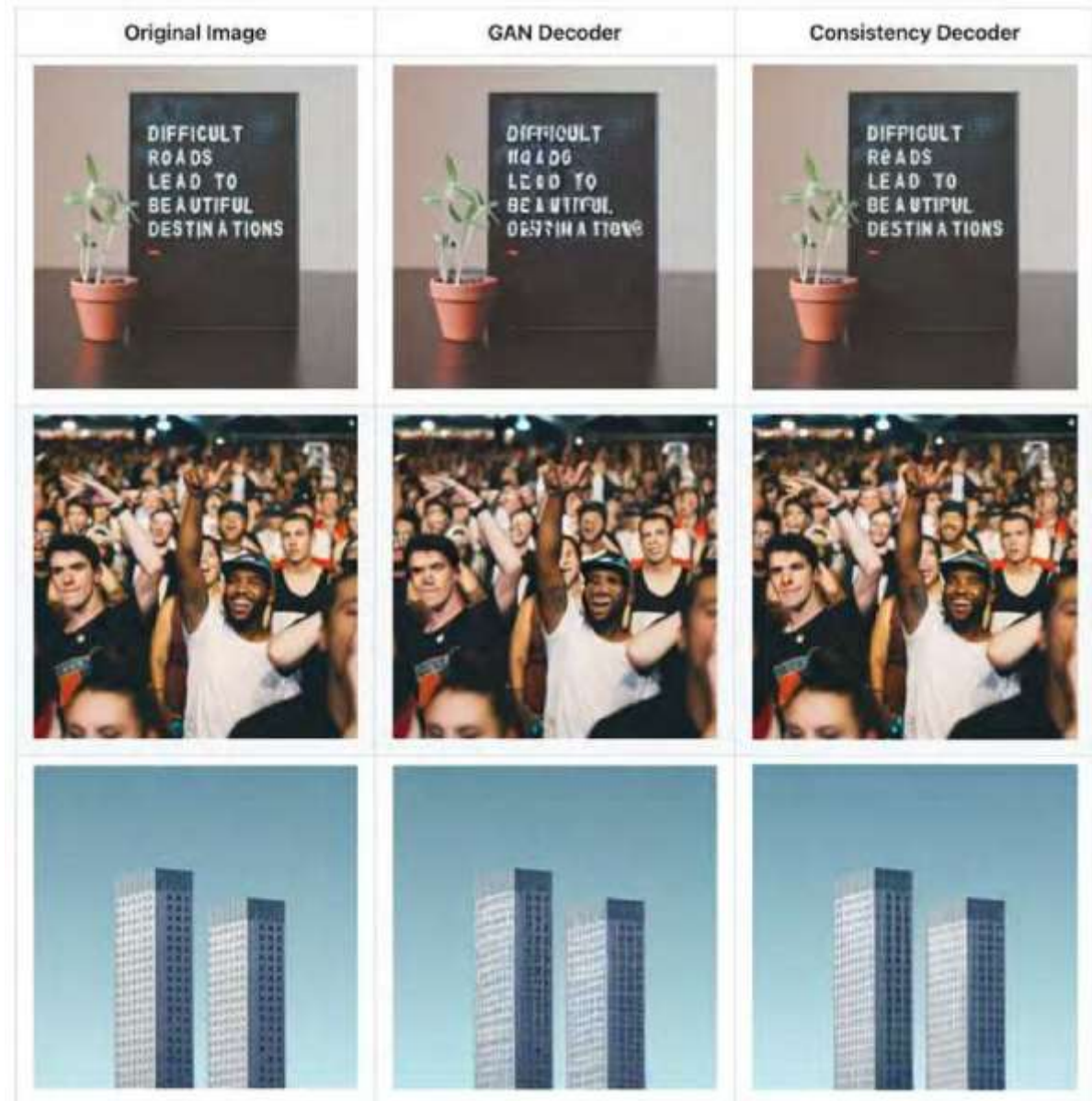


Two steps
FID = 5.22
NFE = 2

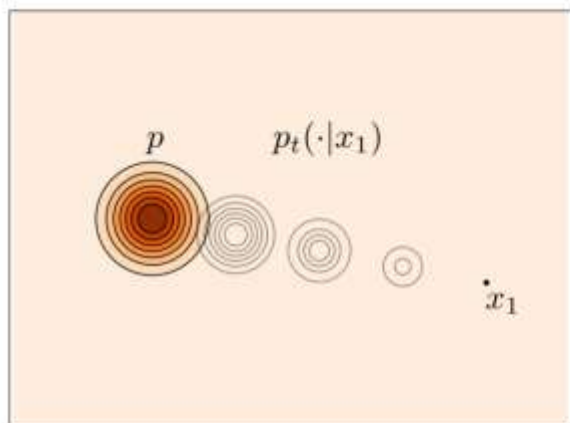
Consistency decoder in DALLE 3

 Gabriel Goh 
@gabeeeegoooh

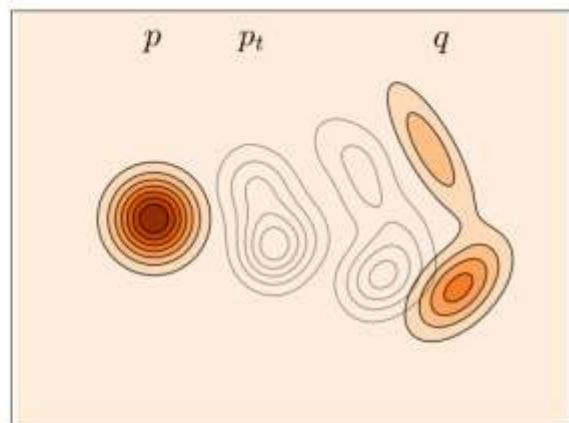
if you get the chance, do try decoding your Stable Diffusion generations with this decoder. You should see some improvements in text, small faces and straight lines. Made with @tim_brooks, @DrYangSong, @model_mechanic, @txhf, @neonbjb
github.com/openai/consist...



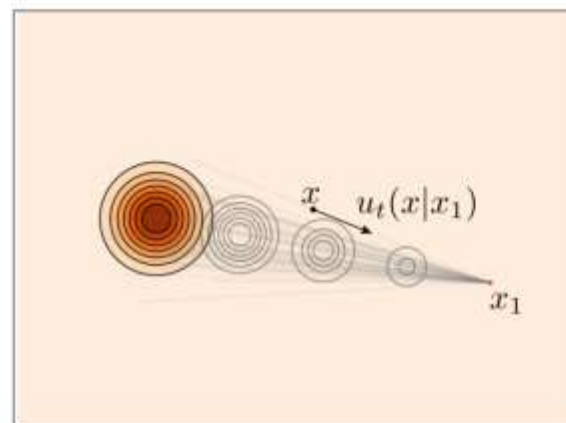
Flow Matching



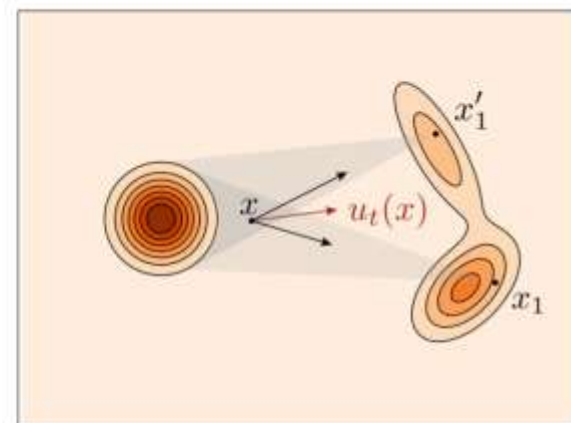
(a) Conditional probability path $p_t(x|x_1)$.



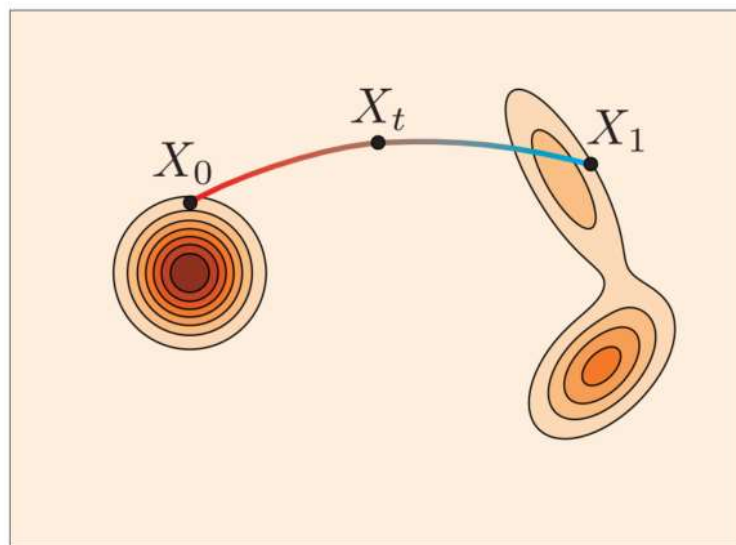
(b) (Marginal) Probability path $p_t(x)$.



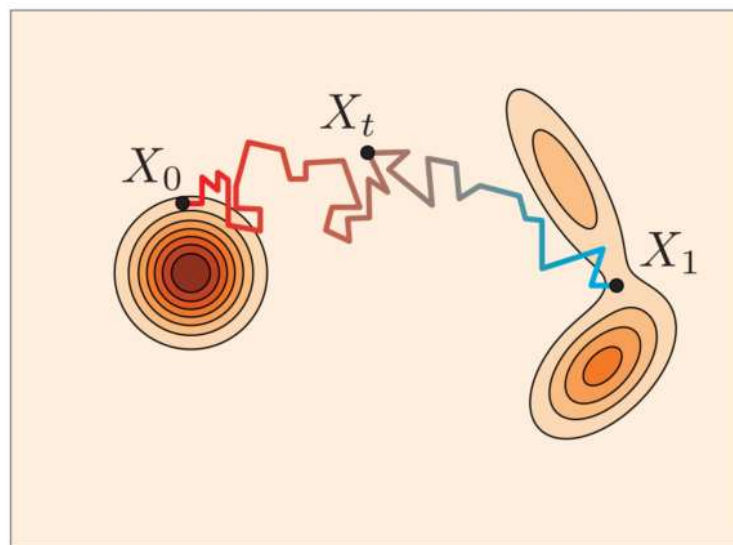
(c) Conditional velocity field $u_t(x|x_1)$.



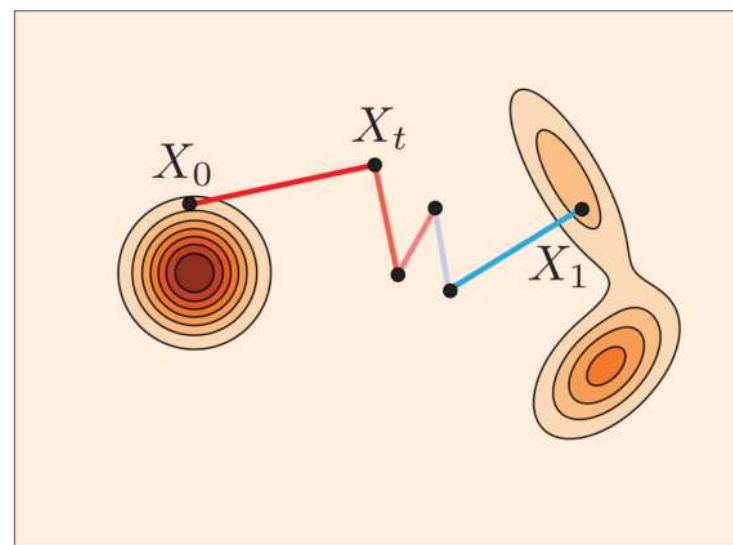
(d) (Marginal) Velocity field $u_t(x)$.



(a) Flow

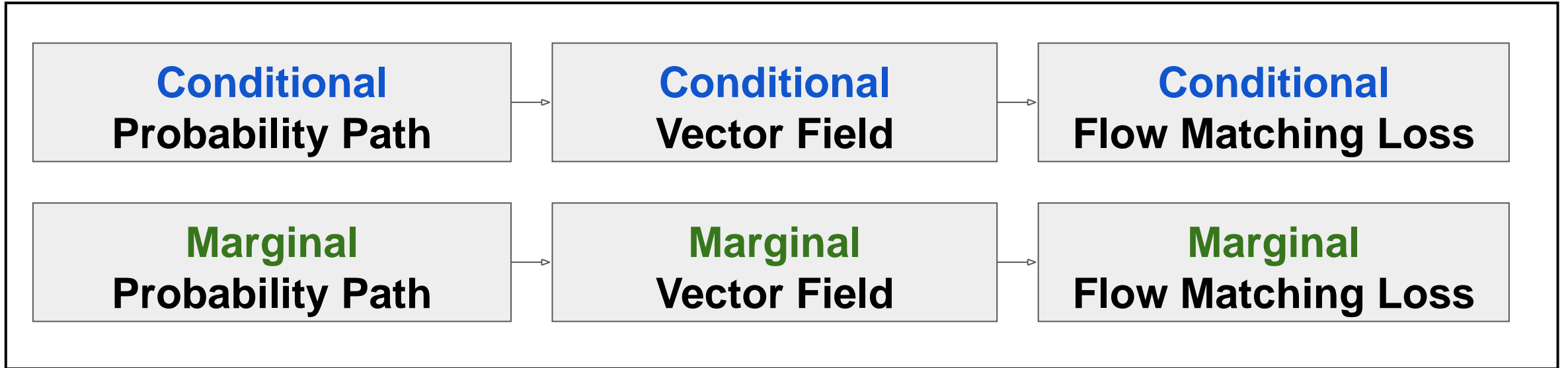


(b) Diffusion



(c) Jump

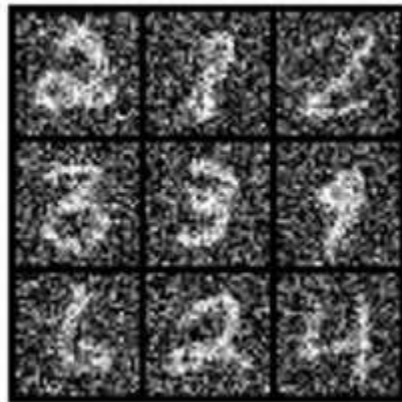
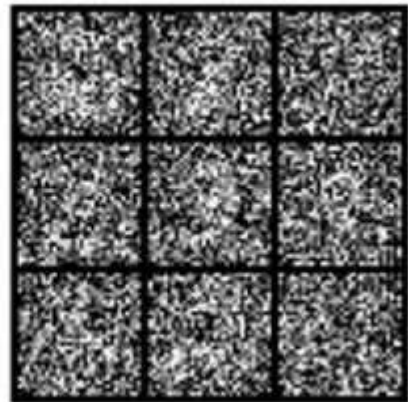
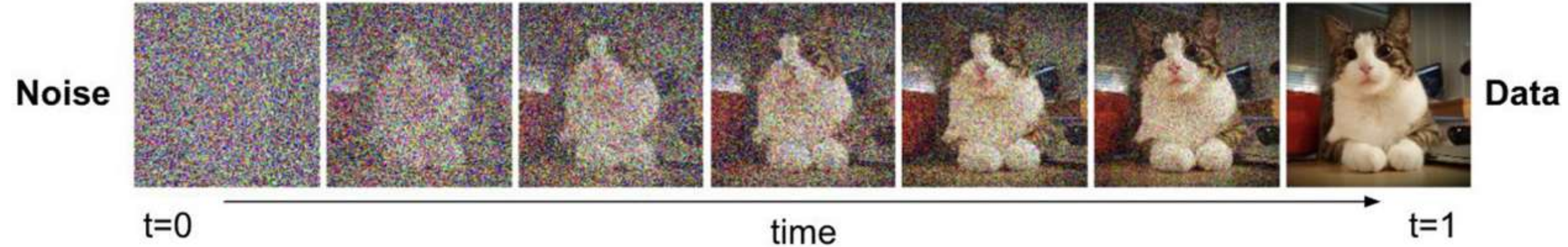
The Flow Matching Matrix



“Conditional” = “Per single data point”

“Marginal” = “Across distribution of data points”

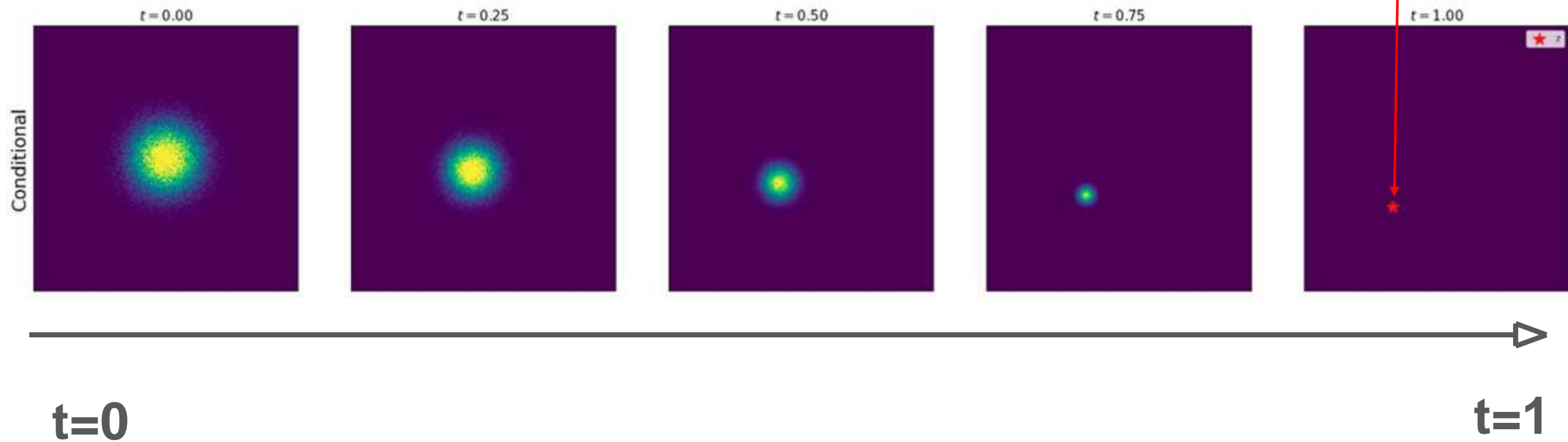
Probability Paths: The Path from Noise to Data



p_{init}

Conditional Probability Path

$$p_t(\cdot|z)$$



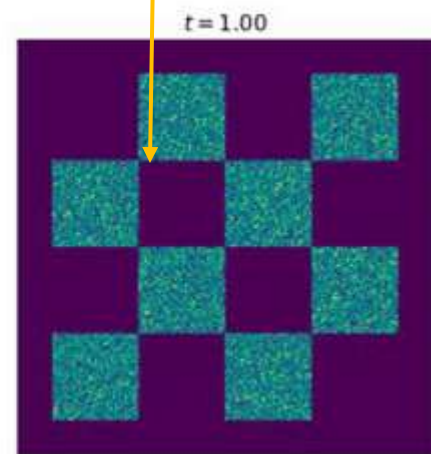
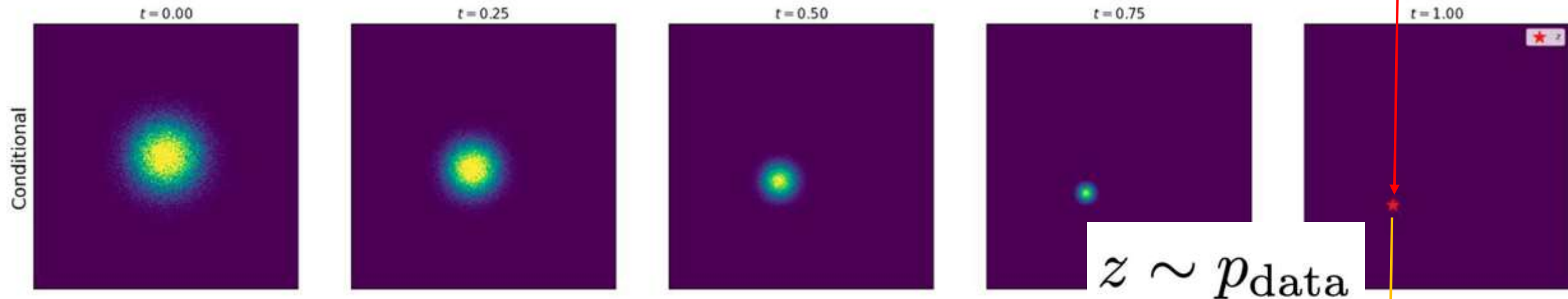
Marginal Probability Path p_t

p_{init}

Conditional Probability Path

$p_t(\cdot|z)$

z

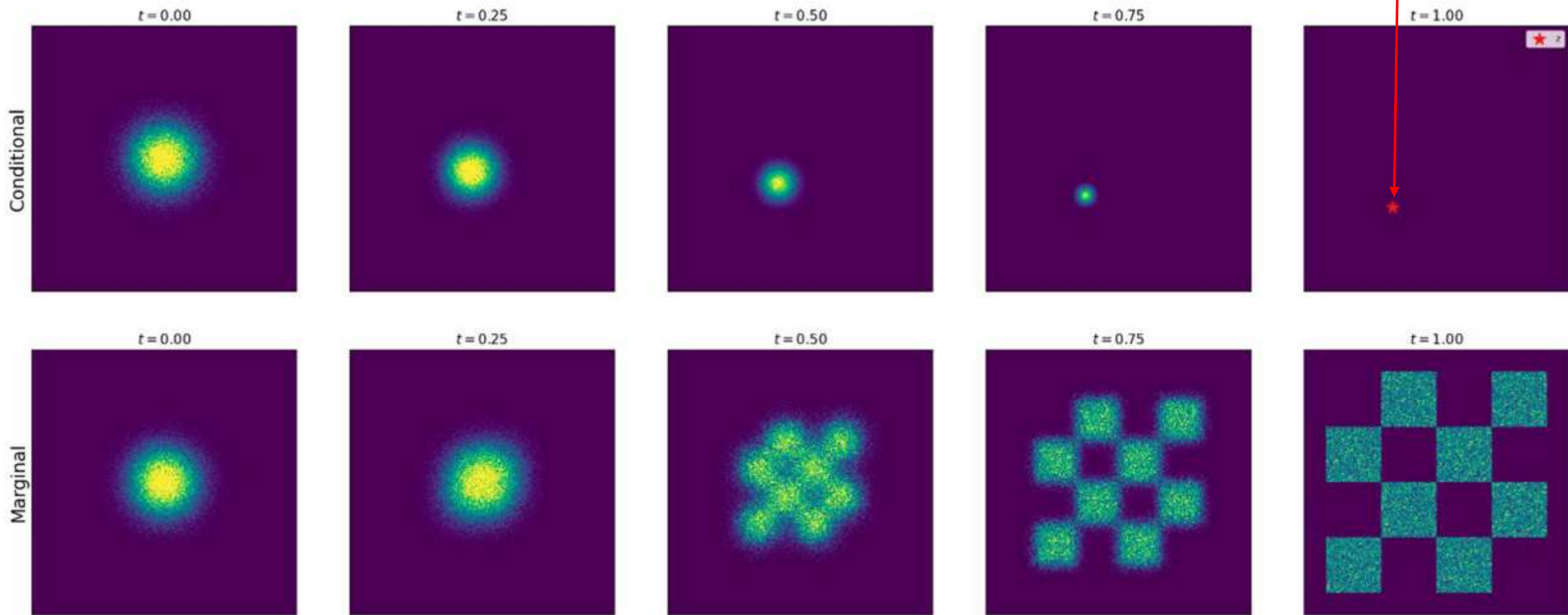


p_{data}

p_{init}

Conditional Probability Path

$$p_t(\cdot|z)$$



p_{init}

Marginal Probability Path

$$p_t$$

p_{data}

Conditional Prob. Path

Notation

Key property

Gaussian example

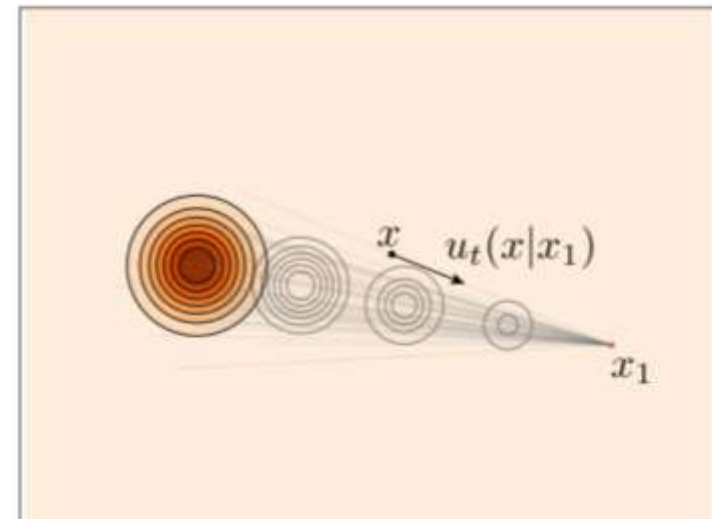
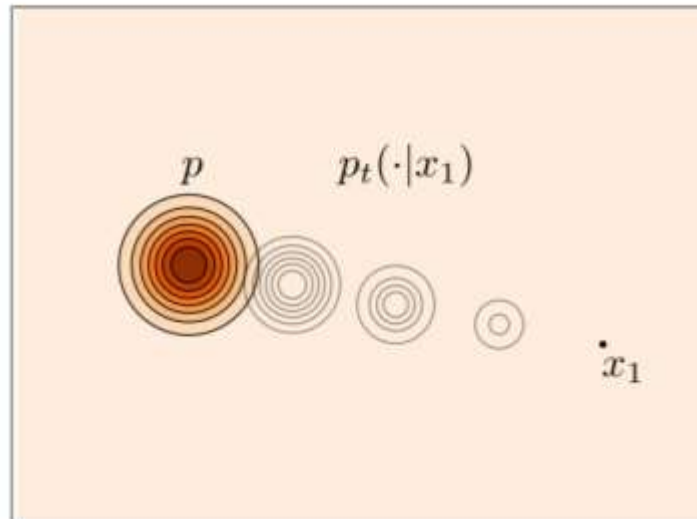
Conditional
Probability Path

$$p_t(\cdot|z)$$

Interpolates p_{init}
and a data point z

$$p_t(\cdot|z) = \mathcal{N}(\alpha_t z, \beta_t^2 I_d)$$

Conditional
Vector Field



Conditional Prob. Path

	Notation	Key property	Gaussian example
Conditional Probability Path	$p_t(\cdot z)$	Interpolates p_{init} and a data point z	$p_t(\cdot z) = \mathcal{N}(\alpha_t z, \beta_t^2 I_d)$
Conditional Vector Field	$u_t^{\text{target}}(x z)$	ODE follows conditional path	$\left(\dot{\alpha}_t - \frac{\dot{\beta}_t}{\beta_t} \alpha_t \right) z + \frac{\dot{\beta}_t}{\beta_t} x$
	$\nabla \log p_t(x z)$	Gradient of log-likelihood	$-\frac{x - \alpha_t z}{\beta_t^2}$

Marginal Prob. Path

Notation

Key property

Formula

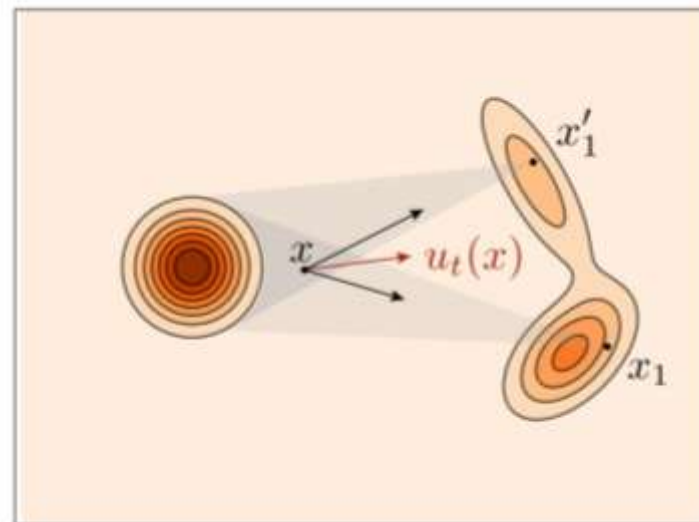
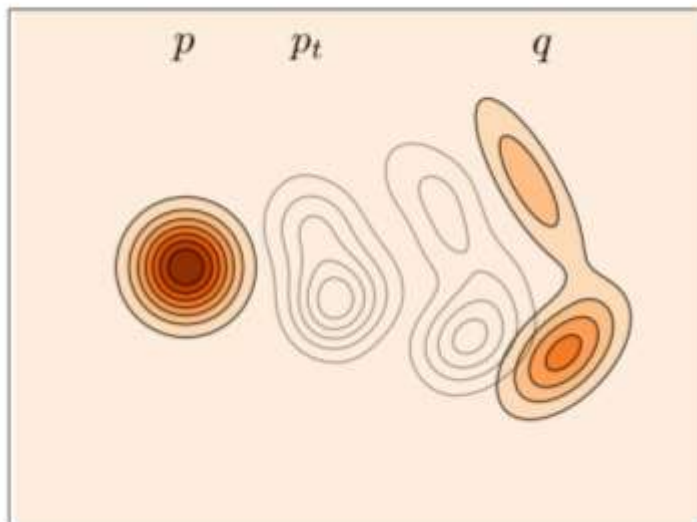
Marginal
Probability
Path

p_t

Interpolates p_{init}
and p_{data}

$$\int p_t(x|z)p_{\text{data}}(z)dz$$

Marginal
Vector
Field



Marginal Prob. Path

	Notation	Key property	Formula
Marginal Probability Path	p_t	Interpolates p_{init} and p_{data}	$\int p_t(x z)p_{\text{data}}(z)dz$
Marginal Vector Field	$u_t^{\text{target}}(x)$	ODE follows marginal path	$\int u_t^{\text{target}}(x z)\frac{p_t(x z)p_{\text{data}}(z)}{p_t(x)}dz$
	$\nabla \log p_t(x)$	Gradient of log-likelihood of marginal path	$\int \nabla \log p_t(x z)\frac{p_t(x z)p_{\text{data}}(z)}{p_t(x)}dz$

Example - Conditional Vector Field for Gaussian

$$u_t^{\text{target}}(x|z) = \left(\dot{\alpha}_t - \frac{\dot{\beta}_t}{\beta_t} \alpha_t \right) z + \frac{\dot{\beta}_t}{\beta_t} x$$

Proof Sketch:

Step 1: By checking ODE, show that the flow of the vector field is given by

$$\psi_t^{\text{target}}(x_0|z) = \alpha_t z + \beta_t x_0$$

Step 2: If $X_0 = x_0 \sim \mathcal{N}(0, I_d)$ is random, then we know that then:

$$X_t = \psi_t(X_0|z) = \alpha_t z + \beta_t X_0 \sim \mathcal{N}(\alpha_t z, \beta_t^2 I_d) = p_t(\cdot|z)$$

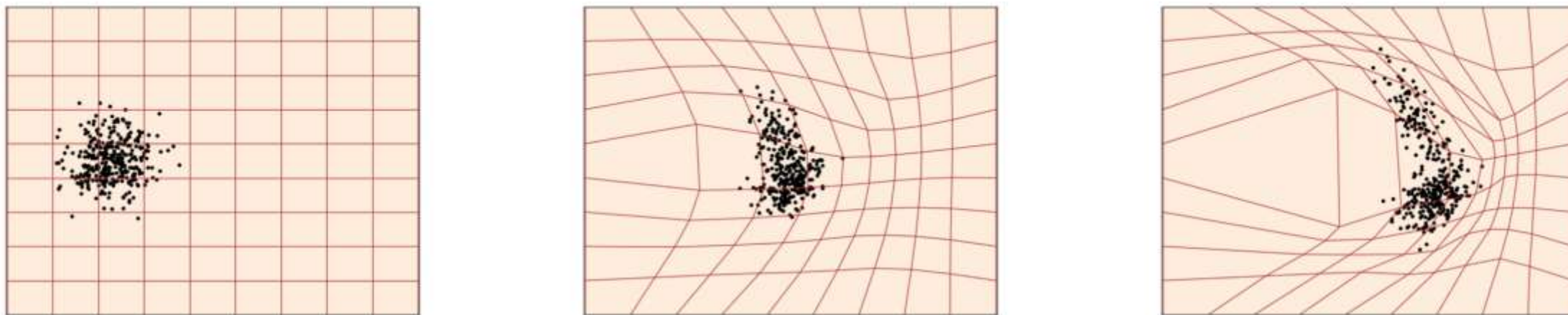


Figure 5 A flow model $X_t = \psi_t(X_0)$ is defined by a diffeomorphism $\psi_t : \mathbb{R}^d \rightarrow \mathbb{R}^d$ (visualized with a brown square grid) pushing samples from a source RV X_0 (left, black points) toward some target distribution q (right). We show three different times t .

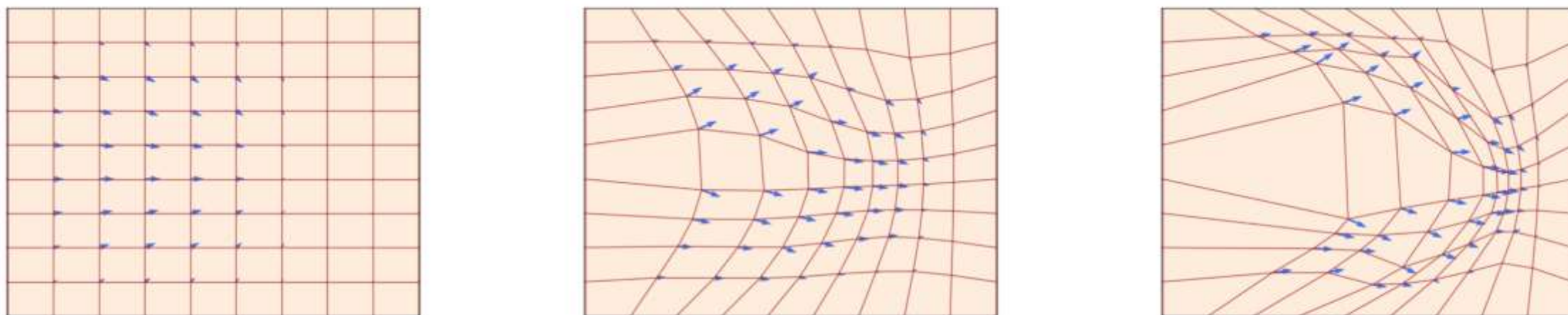


Figure 6 A flow $\psi_t : \mathbb{R}^d \rightarrow \mathbb{R}^d$ (square grid) is defined by a velocity field $u_t : \mathbb{R}^d \rightarrow \mathbb{R}^d$ (visualized with blue arrows) that prescribes its instantaneous movements at all locations. We show three different times t .

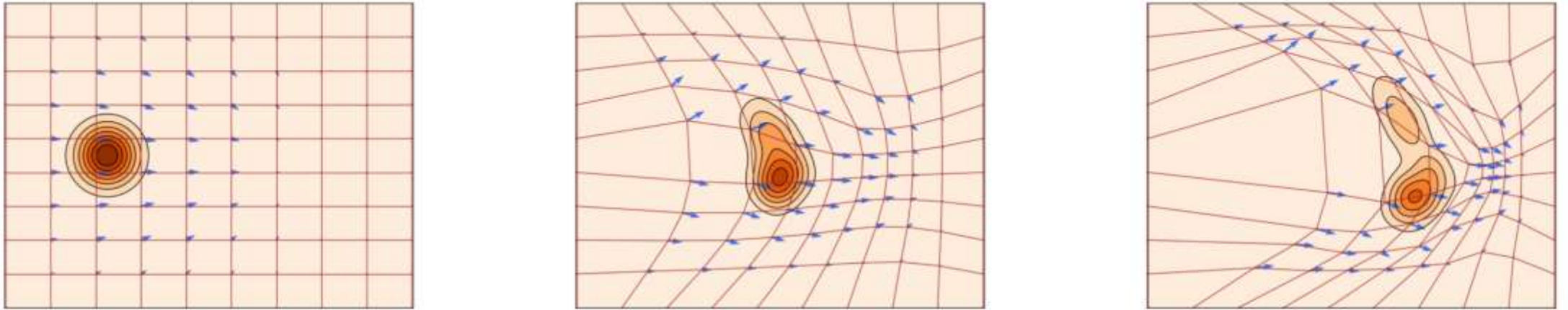
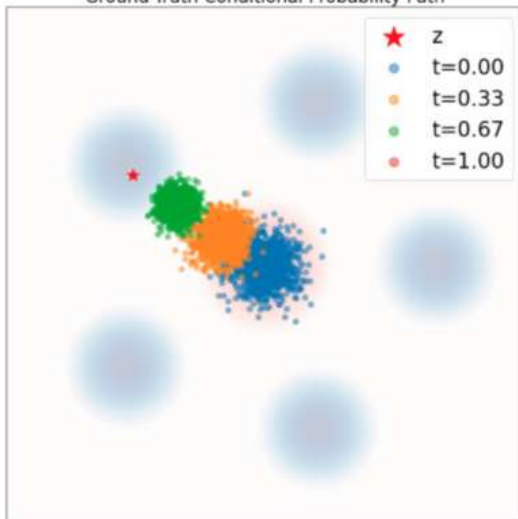


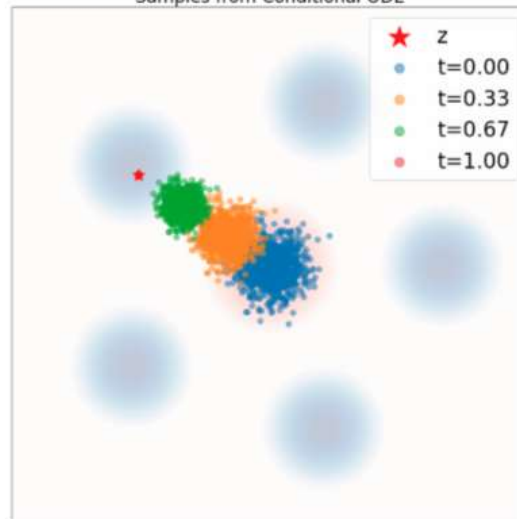
Figure 7 A velocity field u_t (in blue) *generates* a probability path p_t (PDFs shown as contours) if the flow defined by u_t (square grid) reshapes p (left) to p_t at all times $t \in [0, 1)$.

$$p_t(\cdot|z)$$

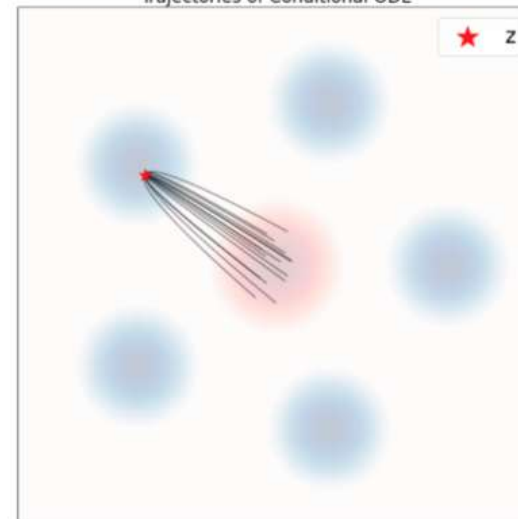
Ground-Truth Conditional Probability Path



Samples from Conditional ODE

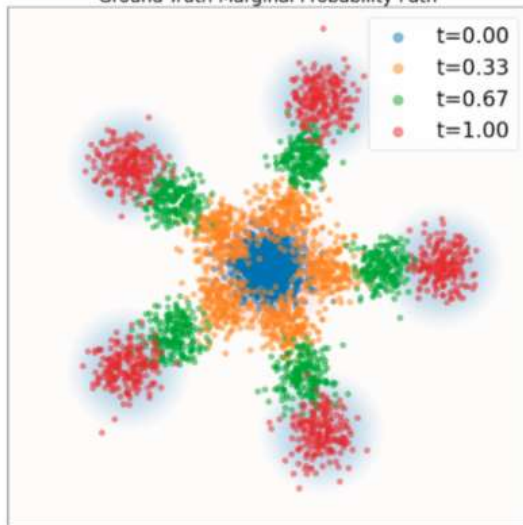


Trajectories of Conditional ODE

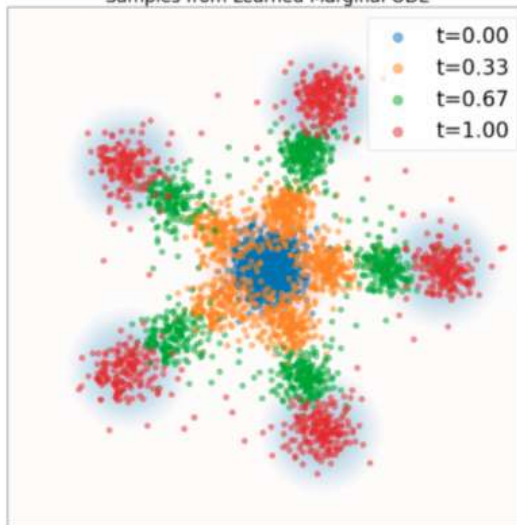


$$p_t$$

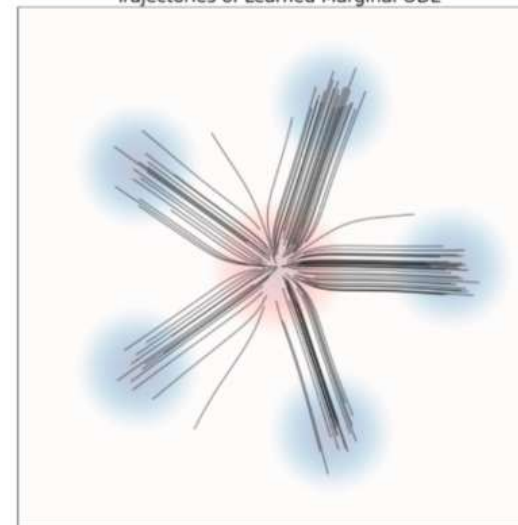
Ground-Truth Marginal Probability Path



Samples from Learned Marginal ODE

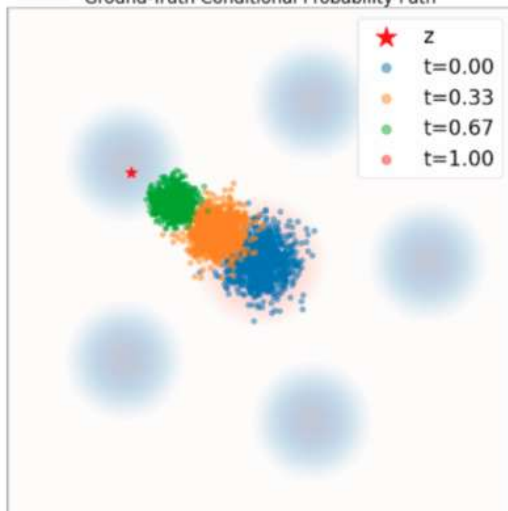


Trajectories of Learned Marginal ODE

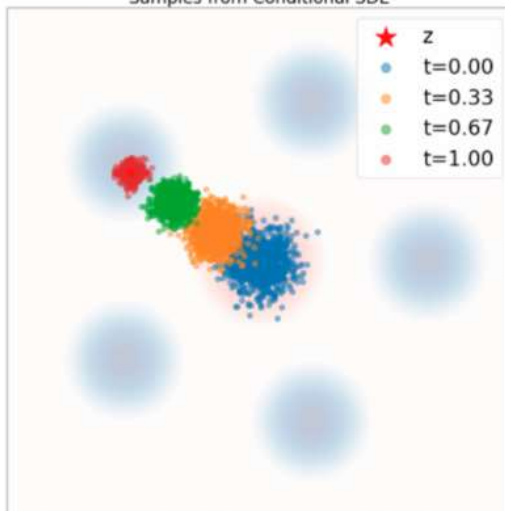


$$p_t(\cdot|z)$$

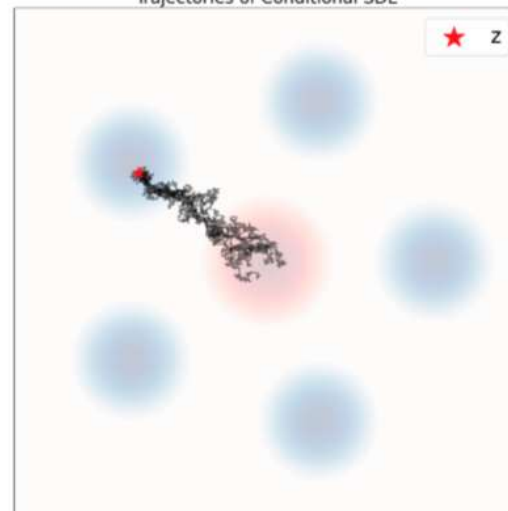
Ground-Truth Conditional Probability Path



Samples from Conditional SDE

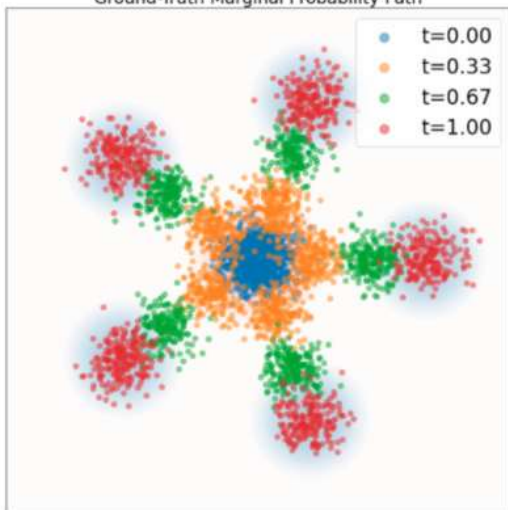


Trajectories of Conditional SDE

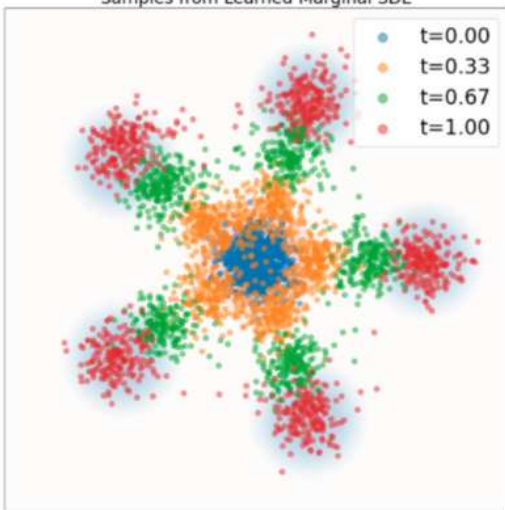


$$p_t$$

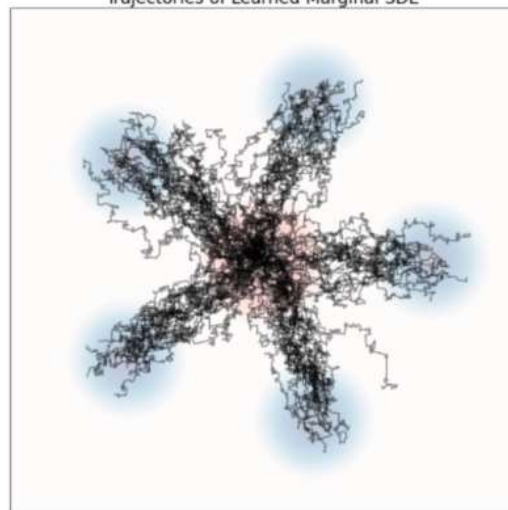
Ground-Truth Marginal Probability Path



Samples from Learned Marginal SDE



Trajectories of Learned Marginal SDE



Continuity Equation

Randomly initialized ODE

Given: $X_0 \sim p_{\text{init}}, \quad \frac{d}{dt}X_t = u_t(X_t)$

Follow probability path:

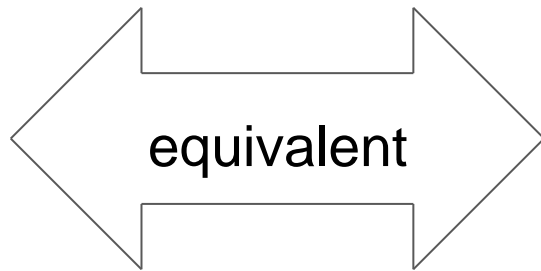
$$X_t \sim p_t \quad (0 \leq t \leq 1)$$

Marginals are p_t

Continuity equation holds

$$\frac{d}{dt}p_t(x) = -\text{div}(p_t u_t)(x)$$

PDE holds

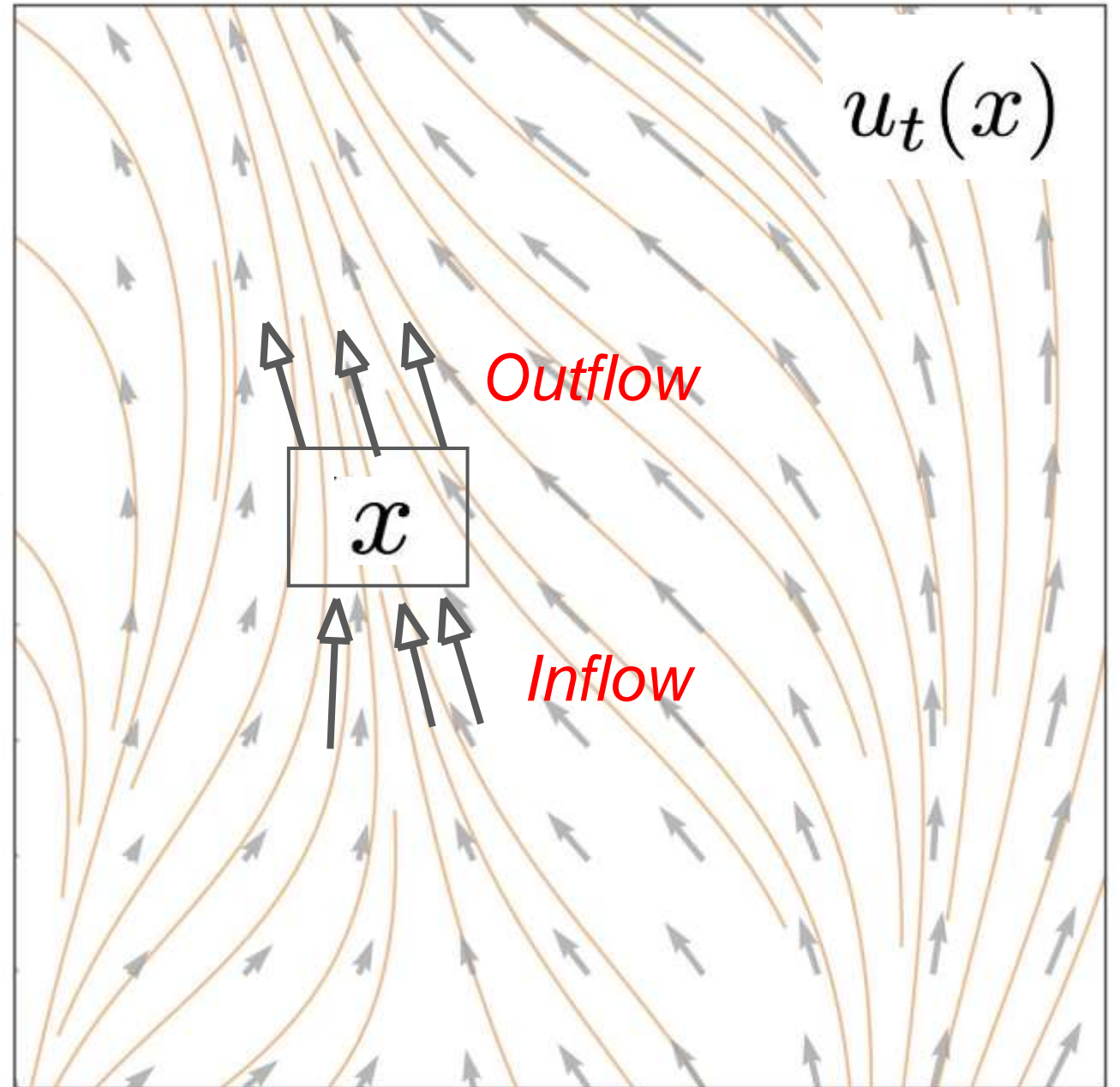


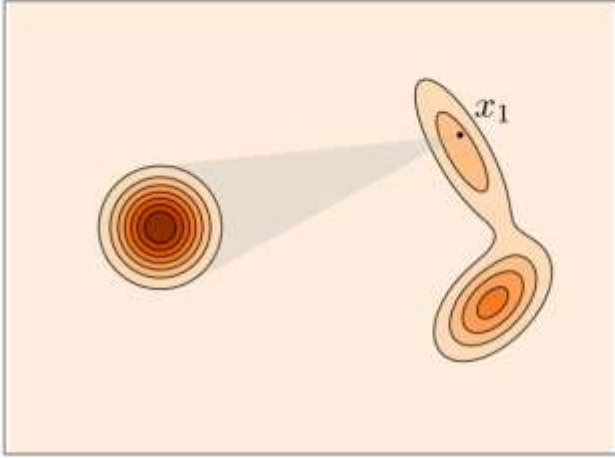
Continuity Equation

$$\frac{d}{dt}p_t(x) = -\operatorname{div}(p_t u_t)(x)$$

*Change of
probability
mass at x*

*Outflow - inflow
of probability
mass from u*

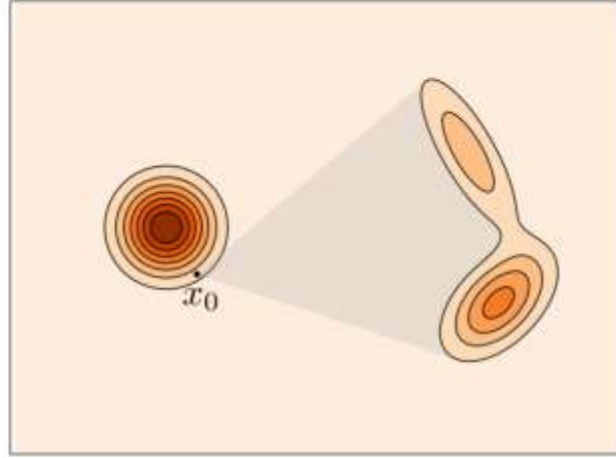


X_1 conditioning

$$\psi_t(X_0|x_1) \sim p_{t|1}(\cdot|x_1)$$

$$p_t(x) = \int p_{t|1}(x|x_1)q(x_1)dx_1$$

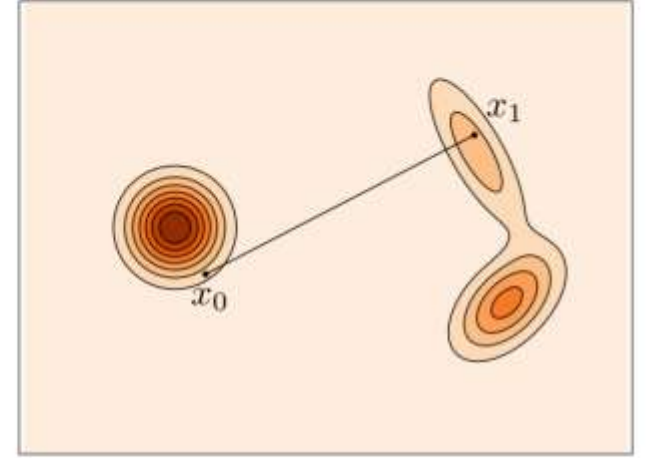
$$u_t(x) = \mathbb{E}[u_t(X_t|X_1)|X_t = x]$$

 X_0 conditioning

$$\psi_t(X_1|x_0) \sim p_{t|0}(\cdot|x_0)$$

$$= \int p_{t|0}(x|x_0)p(x_0)dx_1$$

$$= \mathbb{E}[u_t(X_t|X_0)|X_t = x]$$

 (X_0, X_1) conditioning

$$\psi_t(x_0, x_1) \sim p_{t|0,1}(\cdot|x_0, x_1)$$

$$= \int p_{t|0,1}(x|x_0, x_1)\pi_{0,1}(x_0, x_1)dx_0dx_1$$

$$= \mathbb{E}[u_t(X_t|X_0, X_1)|X_t = x]$$

Algorithm 3 Flow Matching Training Procedure (General)

Require: A dataset of samples $z \sim p_{\text{data}}$, neural network u_t^θ

- 1: **for** each mini-batch of data **do**
- 2: Sample a data example z from the dataset.
- 3: Sample a random time $t \sim \text{Unif}_{[0,1]}$.
- 4: Sample $x \sim p_t(\cdot|z)$
- 5: Compute loss

$$\mathcal{L}(\theta) = \|u_t^\theta(x) - u_t^{\text{target}}(x|z)\|^2$$

- 6: Update the model parameters θ via gradient descent on $\mathcal{L}(\theta)$
 - 7: **end for**
-

Conditional Flow Matching for Gaussian probability path

Prob. path

$$\mathcal{N}(\alpha_t z, \beta_t^2 I_d)$$

Conditional VF

$$u_t^{\text{target}}(x|z) = \left(\dot{\alpha}_t - \frac{\dot{\beta}_t}{\beta_t} \alpha_t \right) z + \frac{\dot{\beta}_t}{\beta_t} x$$

Noise Sampling

$$x \sim p_t(\cdot|z) \iff x = \alpha_t z + \beta_t \epsilon, \quad \epsilon \sim \mathcal{N}(0, I_d)$$

Plugging in Noise Sampling into CFM Loss results in:

$$L_{\text{CFM}}(\theta)$$

$$= \mathbb{E}_{t \sim \text{Unif}, z \sim p_{\text{data}}, x \sim p_t(\cdot|z)} \left[\|u_t^\theta(x) - u_t^{\text{target}}(x|z)\|^2 \right]$$

$$= \mathbb{E}_{t \sim \text{Unif}, z \sim p_{\text{data}}, \epsilon \sim \mathcal{N}(0, I_d)} \left[\|u_t^\theta(\alpha_t z + \beta_t \epsilon) - u_t^{\text{target}}(\alpha_t z + \beta_t \epsilon|z)\|^2 \right]$$

$$= \mathbb{E}_{t \sim \text{Unif}, z \sim p_{\text{data}}, \epsilon \sim \mathcal{N}(0, I_d)} \left[\|u_t^\theta(\alpha_t z + \beta_t \epsilon) - (\dot{\alpha}_t z + \dot{\beta}_t \epsilon)\|^2 \right]$$

noise+data

velocity

Straight Line Schedule

$$L_{\text{CFM}}(\theta)$$

$$= \mathbb{E}_{t \sim \text{Unif}, z \sim p_{\text{data}}, \epsilon \sim \mathcal{N}(0, I_d)} \left[\left\| u_t^\theta (\alpha_t z + \beta_t \epsilon) - (\dot{\alpha}_t z + \dot{\beta}_t \epsilon) \right\|^2 \right]$$

$$= \mathbb{E}_{t \sim \text{Unif}, z \sim p_{\text{data}}, \epsilon \sim \mathcal{N}(0, I_d)} \left[\left\| u_t^\theta (tz + (1-t)\epsilon) - (z - \epsilon) \right\|^2 \right]$$

**Linear
interpolation
of noise and
data**

**Difference
between
noise and
data**

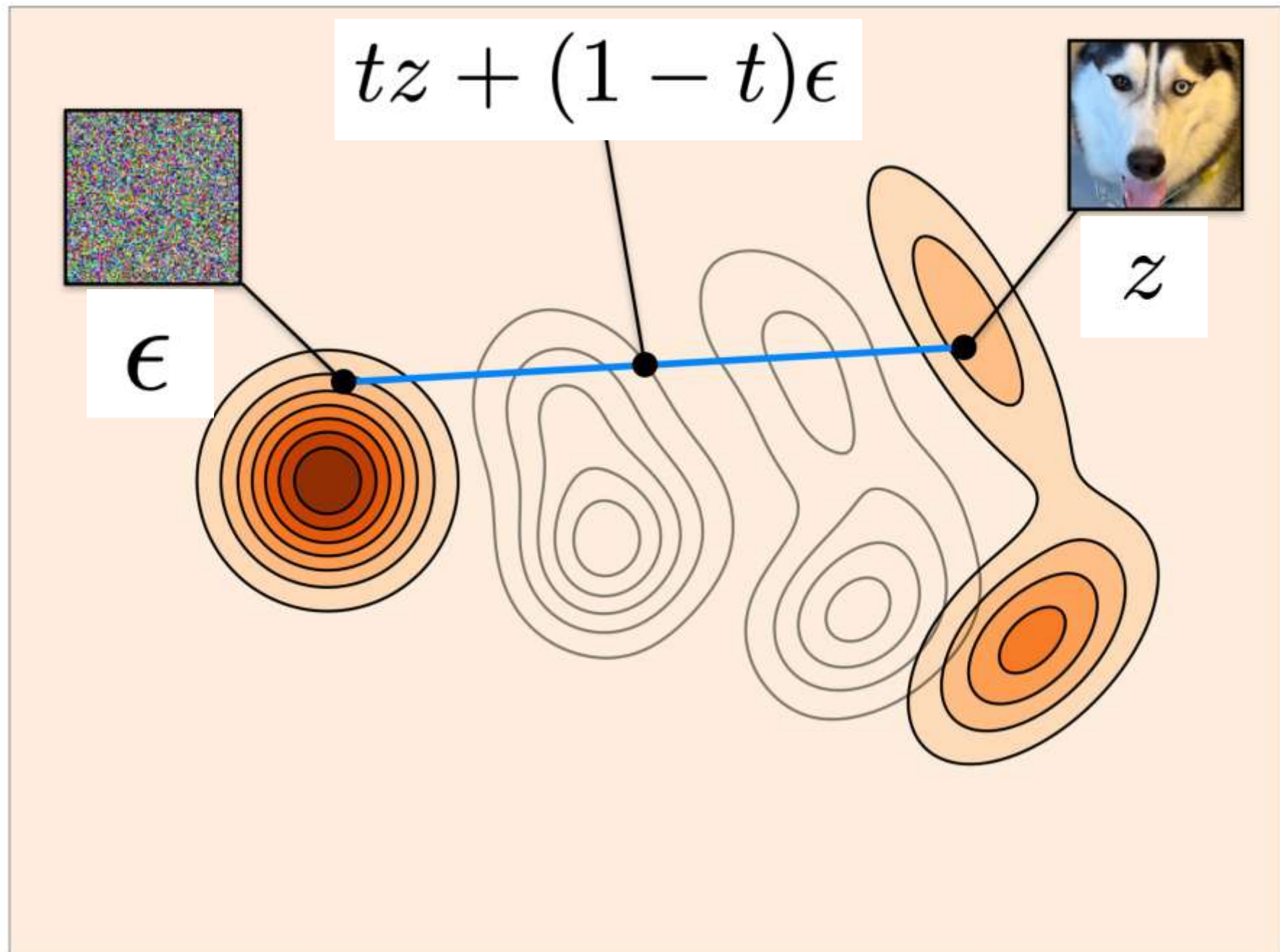


Figure
credit:
Yaron
Lipman

Algorithm 4 Flow Matching Training for CondOT path

Require: A dataset of samples $z \sim p_{\text{data}}$, neural network u_t^θ

- 1: **for** each mini-batch of data **do**
- 2: Sample a data example z from the dataset.
- 3: Sample a random time $t \sim \text{Unif}_{[0,1]}$.
- 4: Sample noise $\epsilon \sim \mathcal{N}(0, I_d)$
- 5: Set $x = tz + (1 - t)\epsilon$
- 6: Compute loss

$$\mathcal{L}(\theta) = \|u_t^\theta(x) - (z - \epsilon)\|^2$$

- 7: Update the model parameters θ via gradient descent on $\mathcal{L}(\theta)$.
 - 8: **end for**
-

Reminder: Sampling Algorithm for Flow Model

Algorithm 1 Sampling from a Flow Model with Euler method

Require: Neural network vector field u_t^θ , number of steps n

- 1: Set $t = 0$
 - 2: Set step size $h = \frac{1}{n}$
 - 3: Draw a sample $X_0 \sim p_{\text{init}}$ *Random initialization!*
 - 4: **for** $i = 1, \dots, n - 1$ **do**
 - 5: $X_{t+h} = X_t + hu_t^\theta(X_t)$
 - 6: Update $t \leftarrow t + h$
 - 7: **end for**
 - 8: **return** X_1 *Return final point*
-

The Flow Matching Matrix

**Conditional
Probability Path**

**Conditional
Vector Field**

**Conditional
Flow Matching Loss**

**Marginal
Probability Path**

**Marginal
Vector Field**

**Marginal
Flow Matching Loss**

**Defines
distributions
from noise to
data**

**Defines training
target that we
want to learn**

**Loss function that
we want to
minimize during
traing**

Conditional Prob. Path, Vector Field, and FM Loss

	Notation	Key property	Gaussian example
Conditional Probability Path	$p_t(\cdot z)$	Interpolates p_{init} and a data point z	$p_t(\cdot z) = \mathcal{N}(\alpha_t z, \beta_t^2 I_d)$
Conditional Vector Field	$u_t^{\text{target}}(x z)$	ODE follows conditional path	$\left(\dot{\alpha}_t - \frac{\dot{\beta}_t}{\beta_t} \alpha_t \right) z + \frac{\dot{\beta}_t}{\beta_t} x$
Conditional FM Loss	$L_{\text{CFM}}(\theta)$	Loss we minimize during training	$\mathbb{E}_{t,z,x} [\ u_t^\theta(x) - u_t^{\text{target}}(x z)\ ^2]$

All these objects are tractable. Just analytical formulas!

Marginal Prob. Path, Vector Field, and FM Loss

	Notation	Key property	Formula
Marginal Probability Path	p_t	Interpolates p_{init} and p_{data}	$\int p_t(x z)p_{\text{data}}(z)dz$
Marginal Vector Field	$u_t^{\text{target}}(x)$	ODE follows marginal path	$\int u_t^{\text{target}}(x z) \frac{p_t(x z)p_{\text{data}}(z)}{p_t(x)} dz$
Marginal FM Loss	$L_{\text{FM}}(\theta)$	Implicitly minimized via cond FM loss	$\mathbb{E}_{t,z,x} [\ u_t^\theta(x) - u_t^{\text{target}}(x)\ ^2]$

None of these objects are tractable. But we can still learn them!

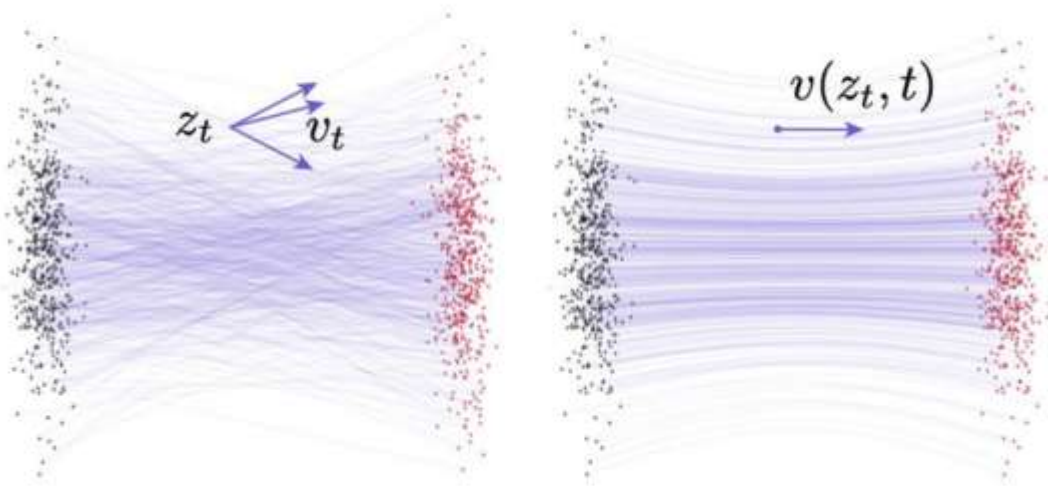


Figure 2: **Velocity fields in Flow Matching** [28]. **Left:** *conditional* flows [28]. A given z_t can arise from different (x, ϵ) pairs, resulting in different conditional velocities v_t . **Right:** *marginal* flows [28], obtained by marginalizing over all possible conditional velocities. The marginal velocity field serves as the underlying ground-truth field for network training. All velocities shown here are essentially *instantaneous* velocities. Illustration follows [12]. (Gray dots: samples from prior; red dots: samples from data.)

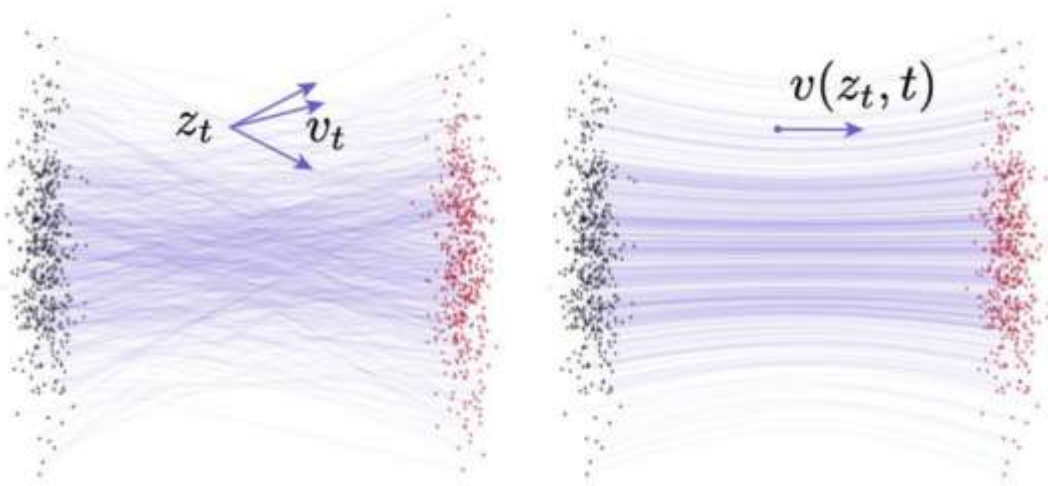
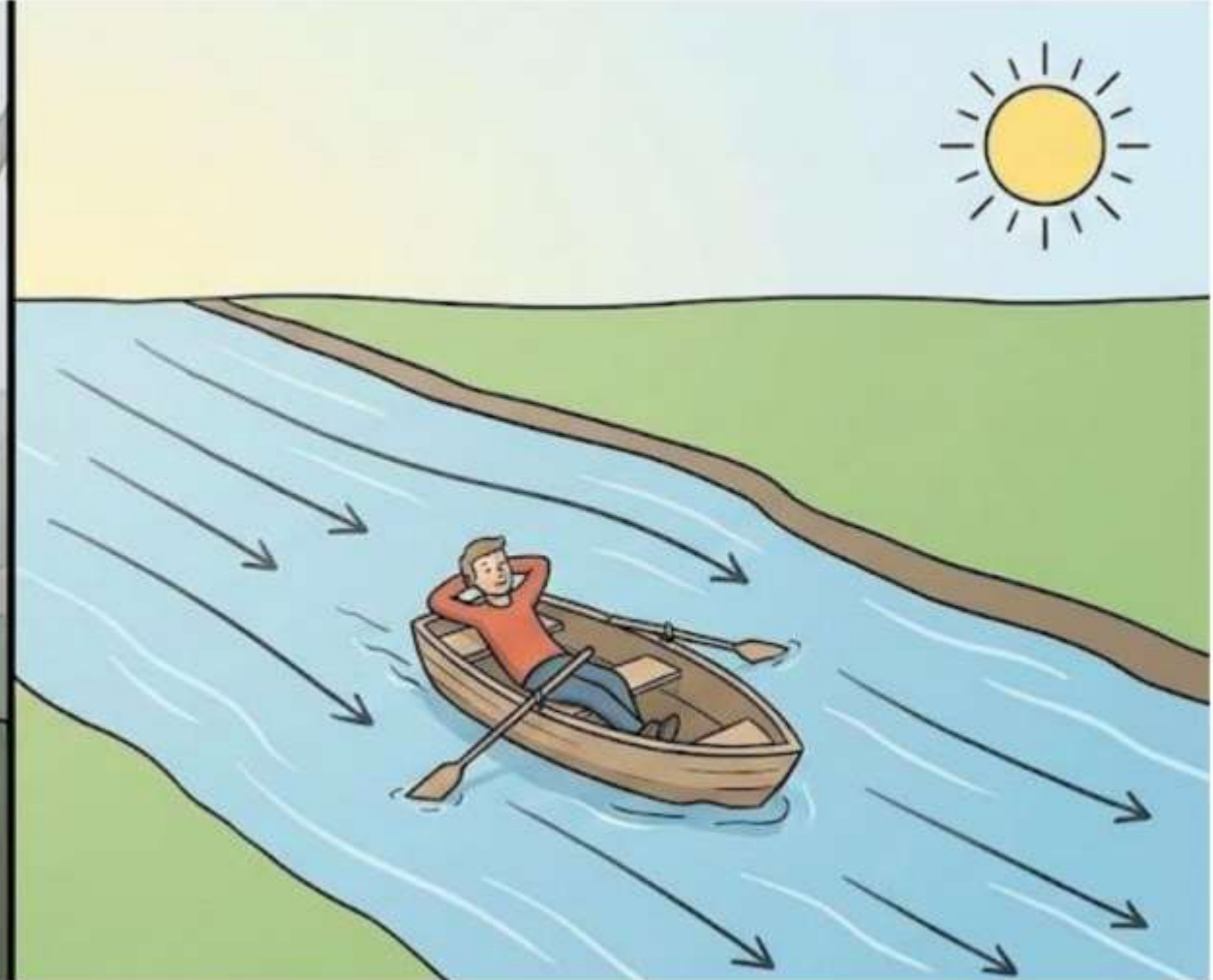


Figure 2: **Velocity fields in Flow Matching** [28]. **Left:** *conditional* flows [28]. A given z_t can arise from different (x, ϵ) pairs, resulting in different conditional velocities v_t . **Right:** *marginal* flows [28], obtained by marginalizing over all possible conditional velocities. The marginal velocity field serves as the underlying ground-truth field for network training. All velocities shown here are essentially *instantaneous* velocities. Illustration follows [12]. (Gray dots: samples from prior; red dots: samples from data.)

Diffusion Models



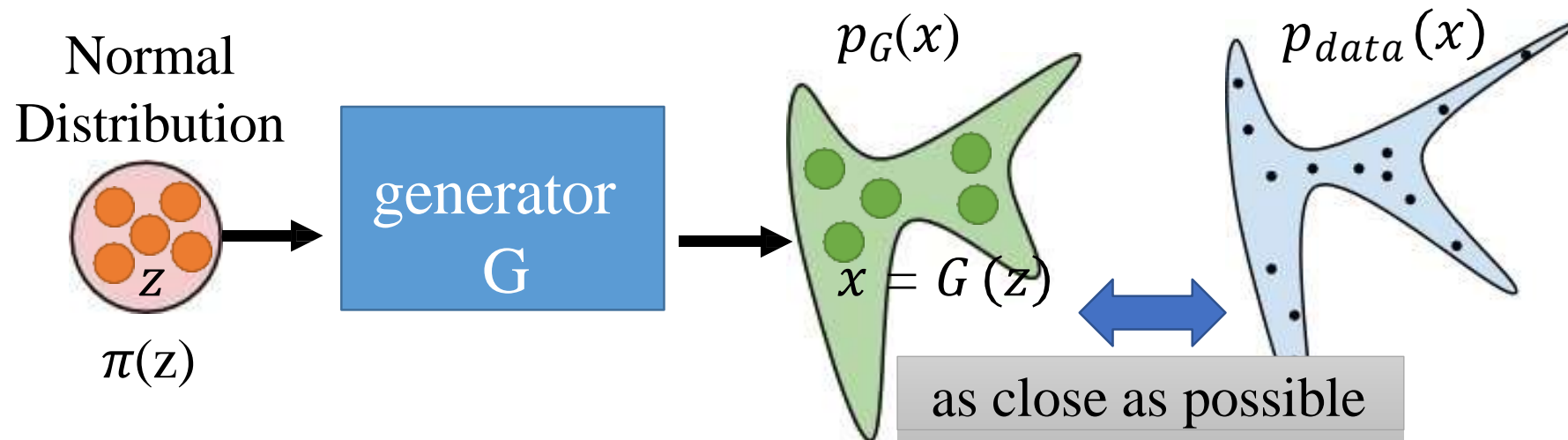
Flow Matching



Flow-Based Models

Generator

- A generator G is a network. The network defines a probability distribution p_G

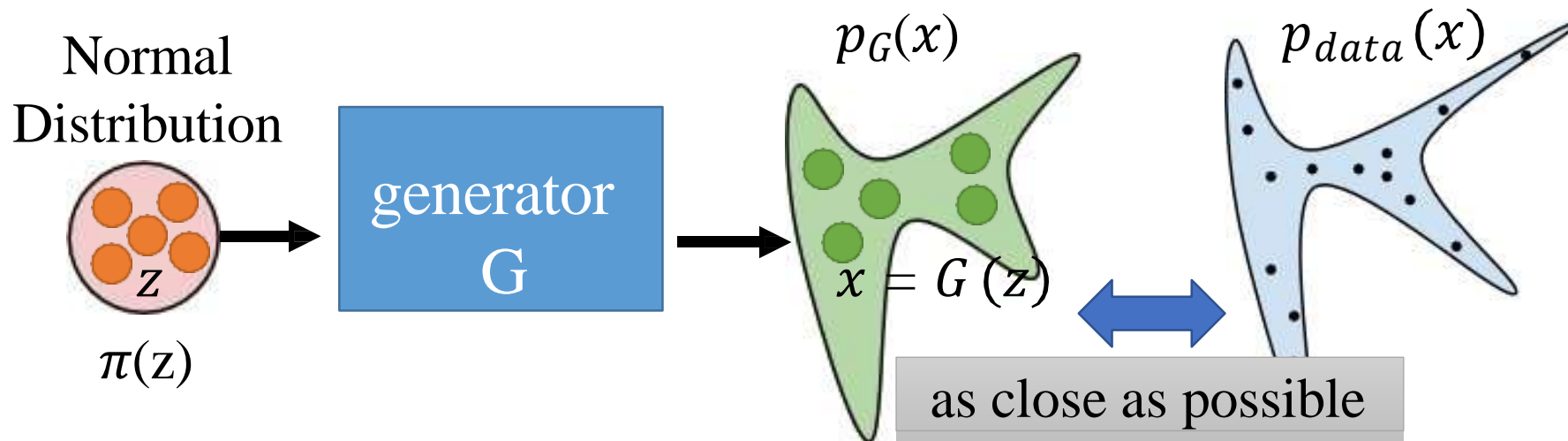


$$G^* = \arg \max_G \sum_{i=1}^m \log P_G(x^i)$$
$$\approx \arg \min_G KL(P_{data} || P_G)$$

$\{x^1, x^2, \dots, x^m\}$ from $P_{data}(x)$

Generator

- A generator G is a network. The network defines a probability distribution p_G



$$G^* = \arg \max_G \sum_{i=1}^m \log P_G(x^i) \quad \{x^1, x^2, \dots, x^m\} \text{ from } P_{data}(x)$$

Flow-based model directly optimizes the objective function.

Math Background

Jacobian, Determinant, Change of Variable Theorem

Jacobian Matrix

$$z = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix} \quad x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$$x = f(z) \quad z = f^{-1}(x)$$

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = f \left(\begin{bmatrix} z_1 \\ z_2 \end{bmatrix} \right)$$

$$\begin{bmatrix} x_2/2 \\ x_1 - x_2/2 \end{bmatrix} = f^{-1} \left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \right)$$

$$J_f = \begin{array}{c} \text{input} \\ \left[\begin{array}{cc} \frac{\partial x_1}{\partial z_1} & \frac{\partial x_1}{\partial z_2} \\ \frac{\partial x_2}{\partial z_1} & \frac{\partial x_2}{\partial z_2} \end{array} \right] \\ \text{output} \end{array}$$

$$J_{f^{-1}} = \begin{bmatrix} \frac{\partial z_1}{\partial x_1} & \frac{\partial z_1}{\partial x_2} \\ \frac{\partial z_2}{\partial x_1} & \frac{\partial z_2}{\partial x_2} \end{bmatrix}$$

$$J_f = \begin{bmatrix} 1 & 1 \\ 2 & 0 \end{bmatrix}$$

$$J_{f^{-1}} = \begin{bmatrix} 0 & 1/2 \\ 1 & -1/2 \end{bmatrix}$$

$$J_f J_{f^{-1}} = I$$

Determinant

The determinant of a **square matrix** is a **scalar** that provides information about the matrix.

- 2 X 2

$$A = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$$

$$\det(A) = ad - bc$$

$$\det(A) = 1/\det(A^{-1})$$

$$\det(J_f) = 1/\det(J_f^{-1})$$

- 3 x 3

$$A = \begin{bmatrix} a_1 & a_2 & a_3 \\ a_4 & a_5 & a_6 \\ a_7 & a_8 & a_9 \end{bmatrix}$$

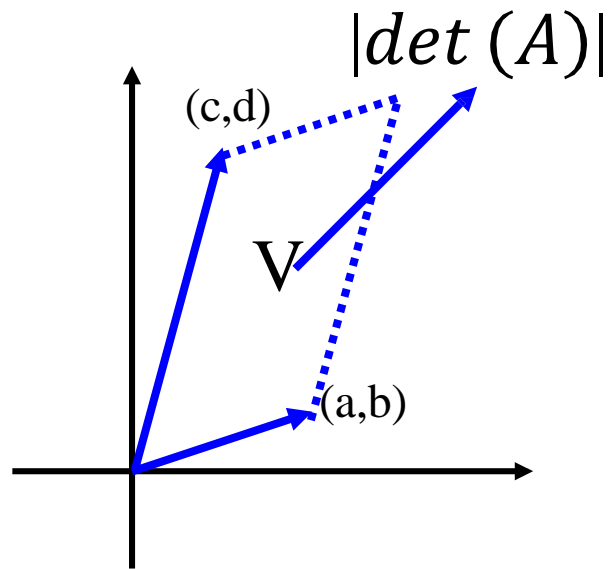
$$\det(A) =$$

$$a_1 a_5 a_9 + a_2 a_6 a_7 + a_3 a_4 a_8 \\ - a_3 a_5 a_7 - a_2 a_4 a_9 - a_1 a_6 a_8$$

Determinant

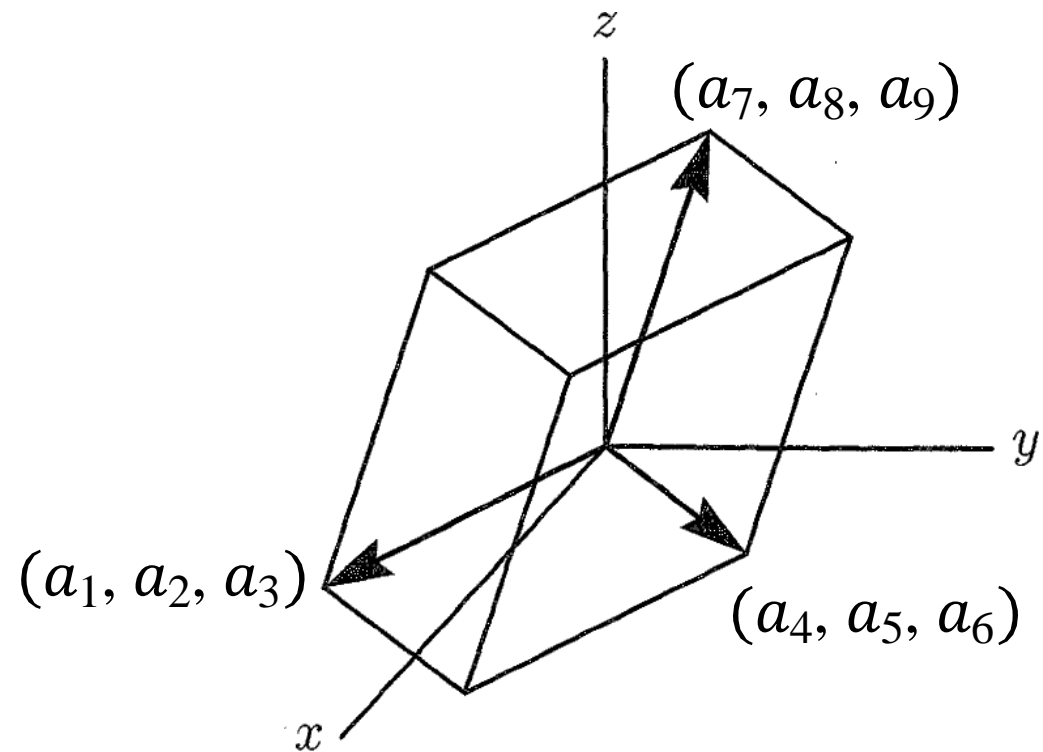
• 2 X 2

$$A = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$$

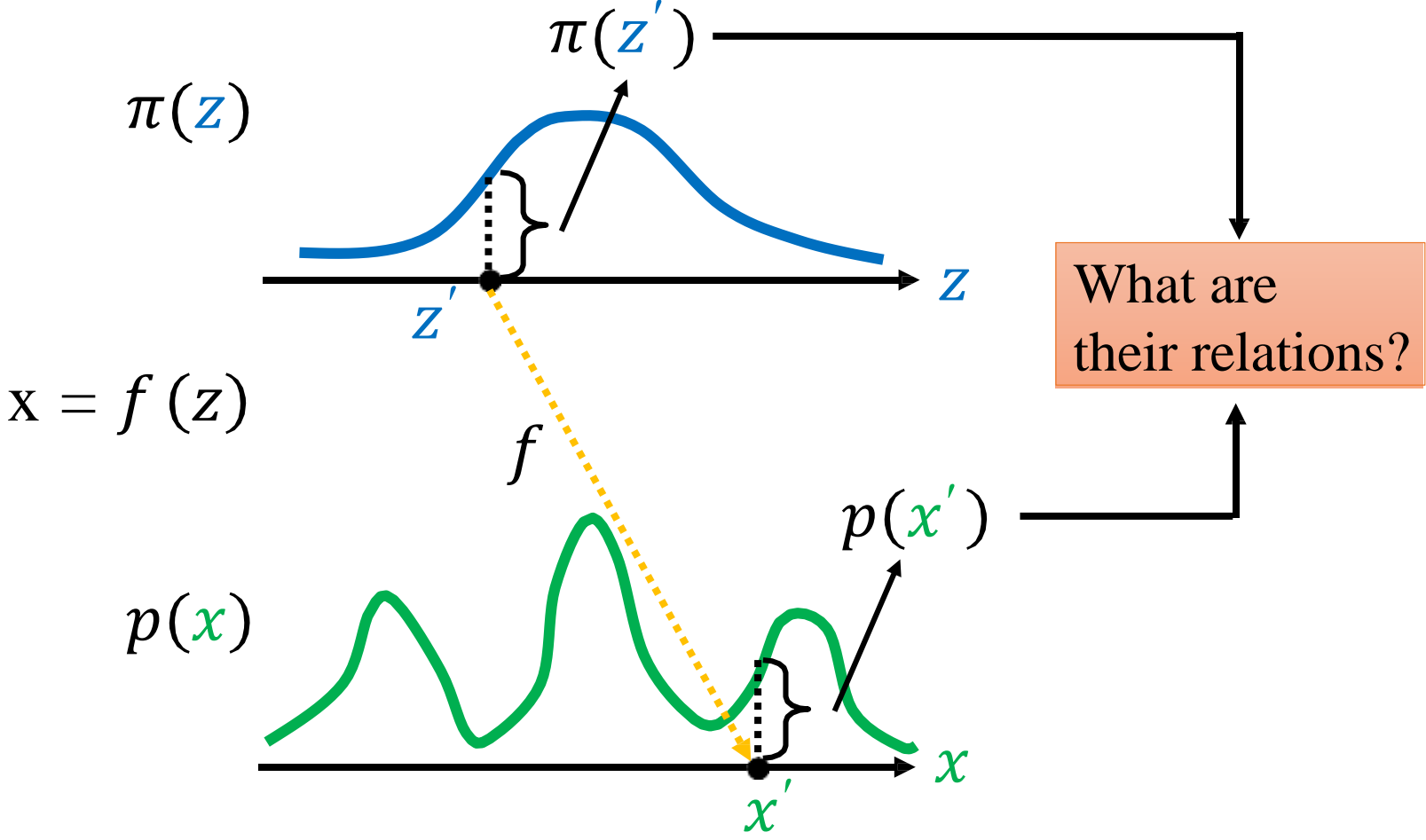


• 3 x 3

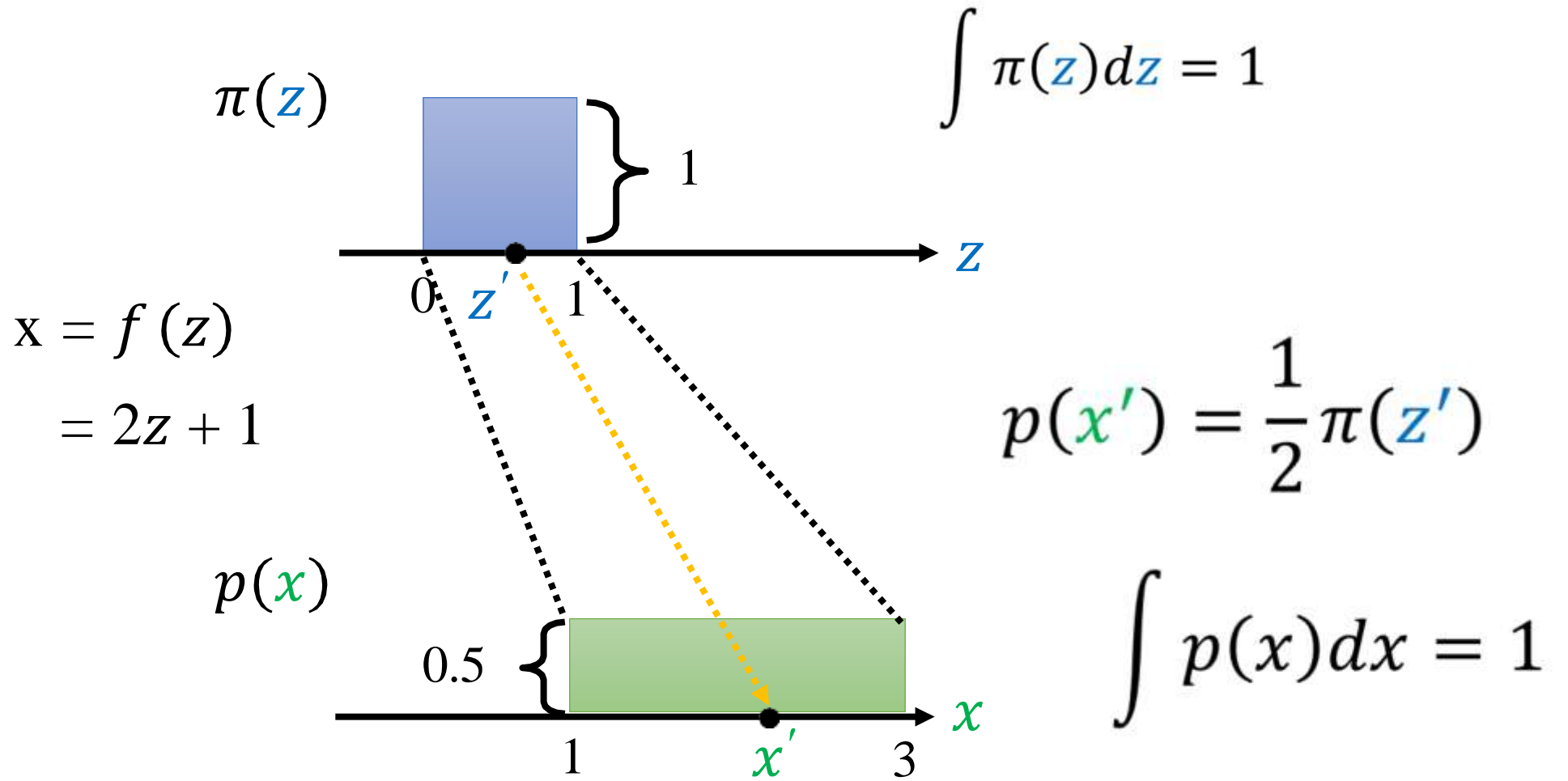
$$A = \begin{bmatrix} a_1 & a_2 & a_3 \\ a_4 & a_5 & a_6 \\ a_7 & a_8 & a_9 \end{bmatrix}$$



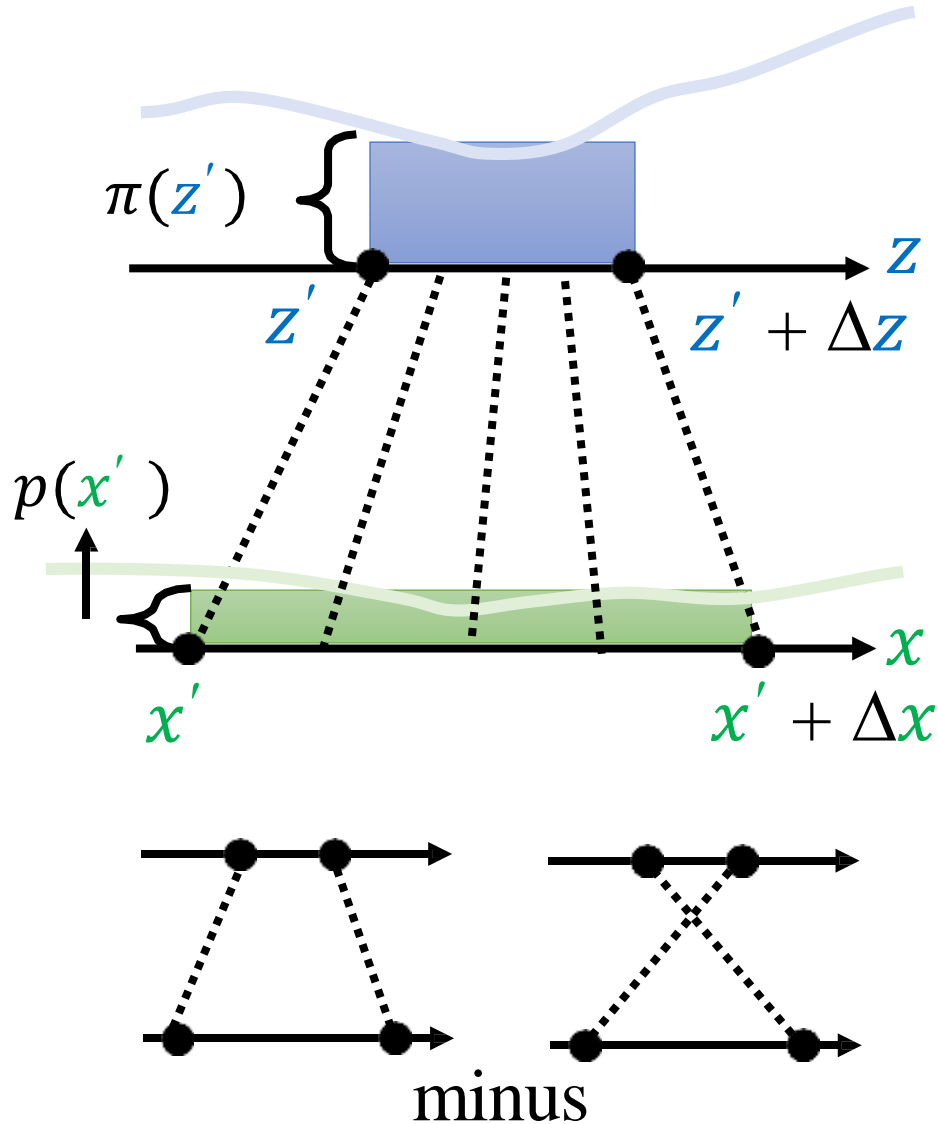
Change of Variable Theorem



Change of Variable Theorem



Change of Variable Theorem

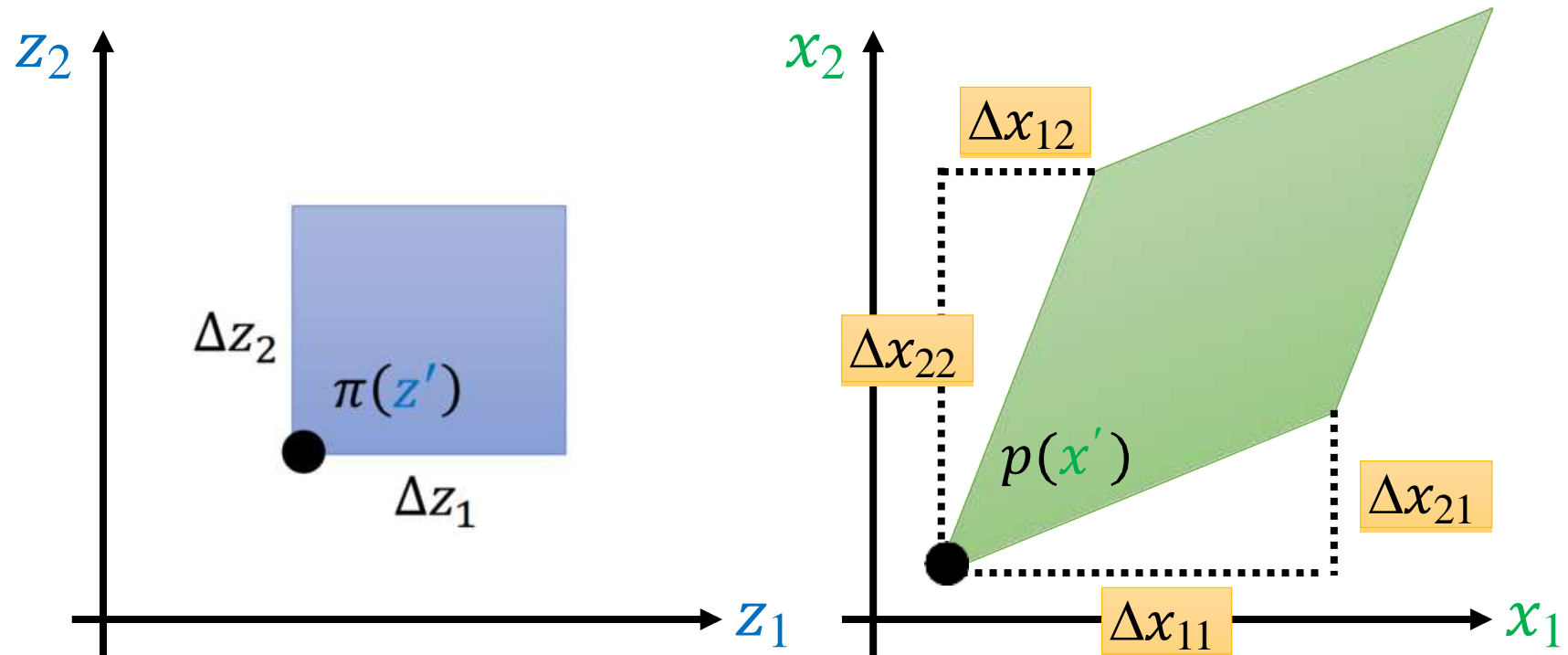


$$p(x') \Delta x = \pi(z') \Delta z$$

$$p(x') = \pi(z') \frac{\Delta z}{\Delta x}$$

$$p(x') = \pi(z') \left| \frac{dz}{dx} \right|$$

Change of Variable Theorem-2D



$$p(x') \left| \det \begin{bmatrix} \Delta x_{11} & \Delta x_{21} \\ \Delta x_{12} & \Delta x_{22} \end{bmatrix} \right| = \pi(z') \Delta z_1 \Delta z_2$$

$$p(\mathbf{x}') \left| \det \begin{bmatrix} \Delta x_{11} & \Delta x_{21} \\ \Delta x_{12} & \Delta x_{22} \end{bmatrix} \right| = \pi(\mathbf{z}') \Delta z_1 \Delta z_2 \quad \mathbf{x} = f(\mathbf{z})$$

$$p(\mathbf{x}') \left| \frac{1}{\Delta z_1 \Delta z_2} \det \begin{bmatrix} \Delta x_{11} & \Delta x_{21} \\ \Delta x_{12} & \Delta x_{22} \end{bmatrix} \right| = \pi(\mathbf{z}')$$

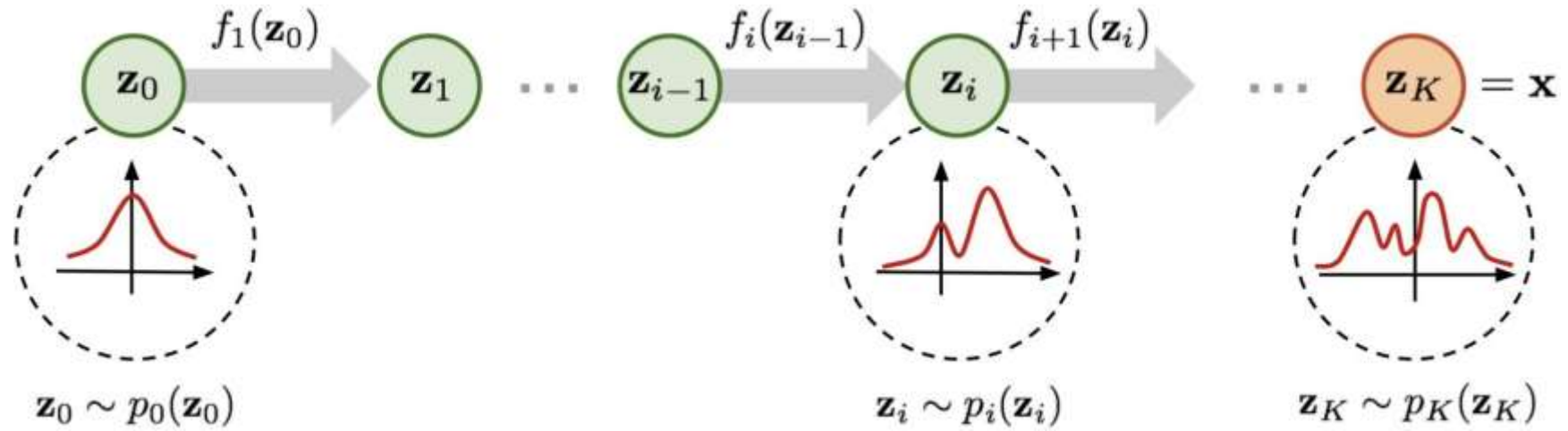
$$p(\mathbf{x}') \left| \det \begin{bmatrix} \Delta x_{11}/\Delta z_1 & \Delta x_{21}/\Delta z_1 \\ \Delta x_{12}/\Delta z_2 & \Delta x_{22}/\Delta z_2 \end{bmatrix} \right| = \pi(\mathbf{z}')$$

$$p(\mathbf{x}') \left| \det \begin{bmatrix} \partial x_1/\partial z_1 & \partial x_2/\partial z_1 \\ \partial x_1/\partial z_2 & \partial x_2/\partial z_2 \end{bmatrix} \right| = \pi(\mathbf{z}')$$

$$p(\mathbf{x}') \left| \det \begin{bmatrix} \partial x_1/\partial z_1 & \partial x_1/\partial z_2 \\ \partial x_2/\partial z_1 & \partial x_2/\partial z_2 \end{bmatrix} \right| = \pi(\mathbf{z}')$$

$$p(\mathbf{x}') | \det(J_f) | = \pi(\mathbf{z}') \quad p(\mathbf{x}') = \pi(\mathbf{z}') \left| \frac{1}{\det(J_f)} \right|$$

$$p(\mathbf{x}') = \pi(\mathbf{z}') | \det(J_{f^{-1}}) |$$

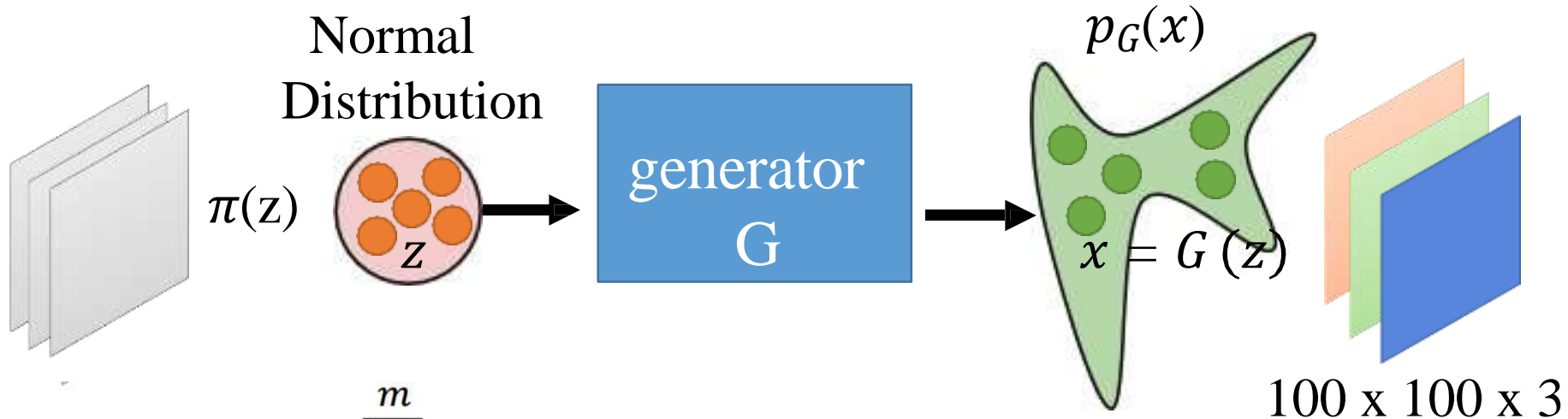


Formal Explanation

Flow-based Model

$$p(x') | \det(J_f) | = \pi(z')$$

$$p(x') = \pi(z') | \det(J_{f^{-1}}) |$$



$$G^* = \arg \max_G \sum_{i=1}^m \log p_G(x^i)$$

$$p_G(x^i) = \pi(z^i) | \det(J_{G^{-1}}) |$$

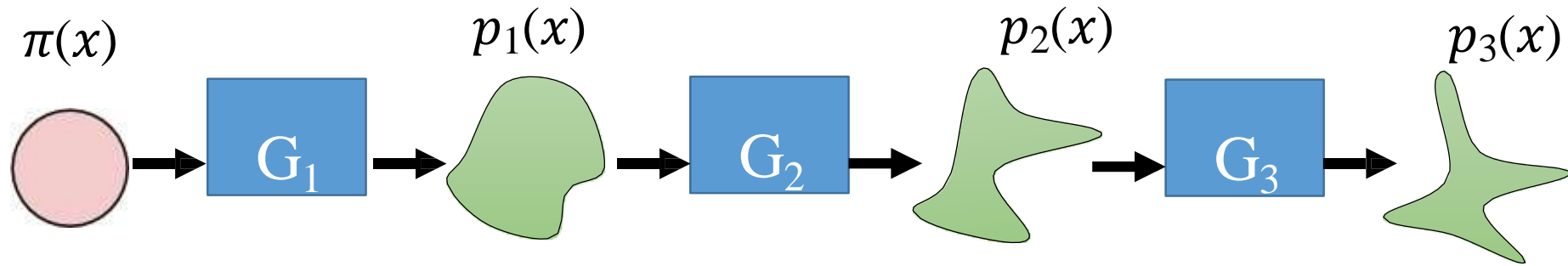
$$z^i = G^{-1}(x^i)$$

You can compute $\det(J_G)$
 You know G^{-1}

G has limitation

$$\log p_G(x^i) = \log \pi(G^{-1}(x^i)) + \log | \det(J_{G^{-1}}) |$$

Multiple Generators



$$p_1(x^i) = \pi(z^i) \left(\left| \det \left(J_{G_1^{-1}} \right) \right| \right) \quad z^i = G_1^{-1} \left(\dots G_K^{-1} (x^i) \right)$$

$$p_2(x^i) = \pi(z^i) \left(\left| \det \left(J_{G_1^{-1}} \right) \right| \right) \left(\left| \det \left(J_{G_2^{-1}} \right) \right| \right)$$

\vdots
 \vdots

$$p_K(x^i) = \pi(z^i) \left(\left| \det \left(J_{G_1^{-1}} \right) \right| \right) \dots \left(\left| \det \left(J_{G_K^{-1}} \right) \right| \right)$$

$$\log p_K(x^i) = \log \pi(z^i) + \sum_{h=1}^K \log \left| \det \left(J_{G_h^{-1}} \right) \right| \quad \text{Maximize}$$

What you actually do?



$$\log p_G(x^i) = \log \pi(G^{-1}(x^i)) + \log |\det(J_{G^{-1}})|$$

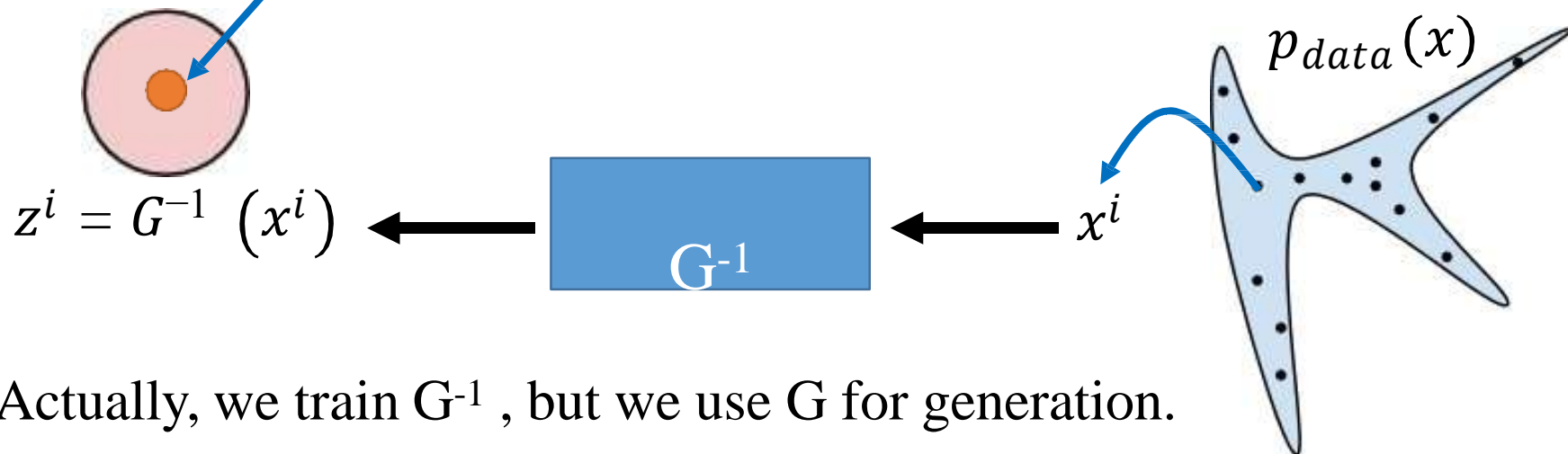
↑ -inf

Make z^i become zero vector

If z^i is always zero:

$J_{G^{-1}}$ would be zero matrix

$$\det(J_{G^{-1}}) = 0$$



Actually, we train G^{-1} , but we use G for generation.

Coupling Layer

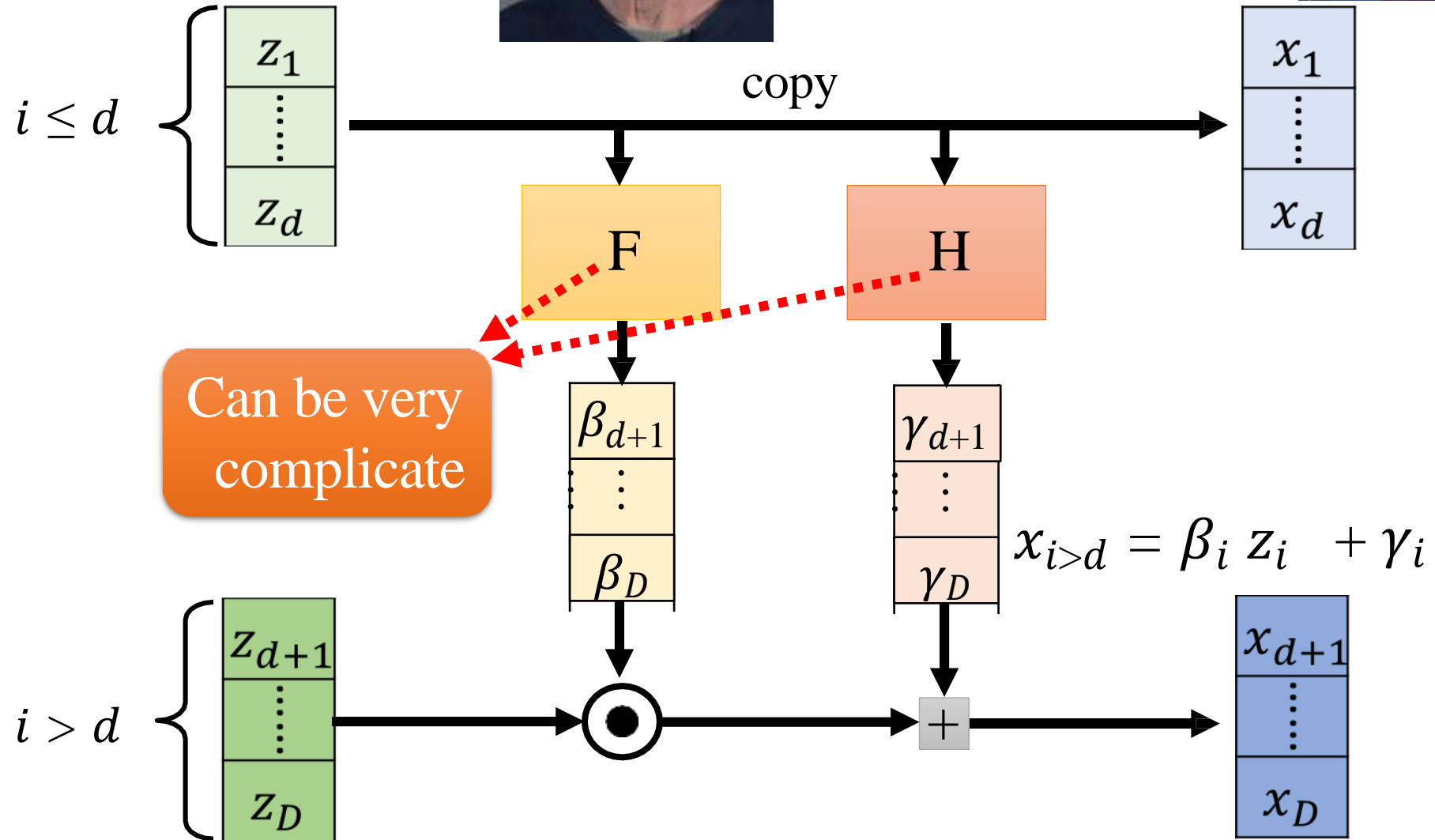
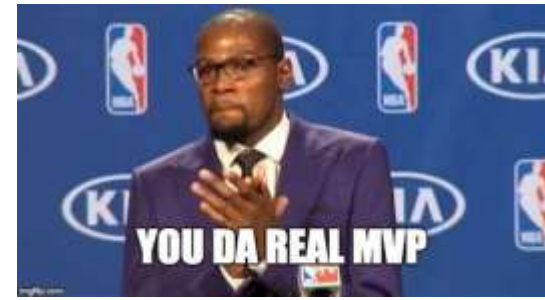


NICE

<https://arxiv.org/abs/1410.8516>

Real NVP

<https://arxiv.org/abs/1605.08803>



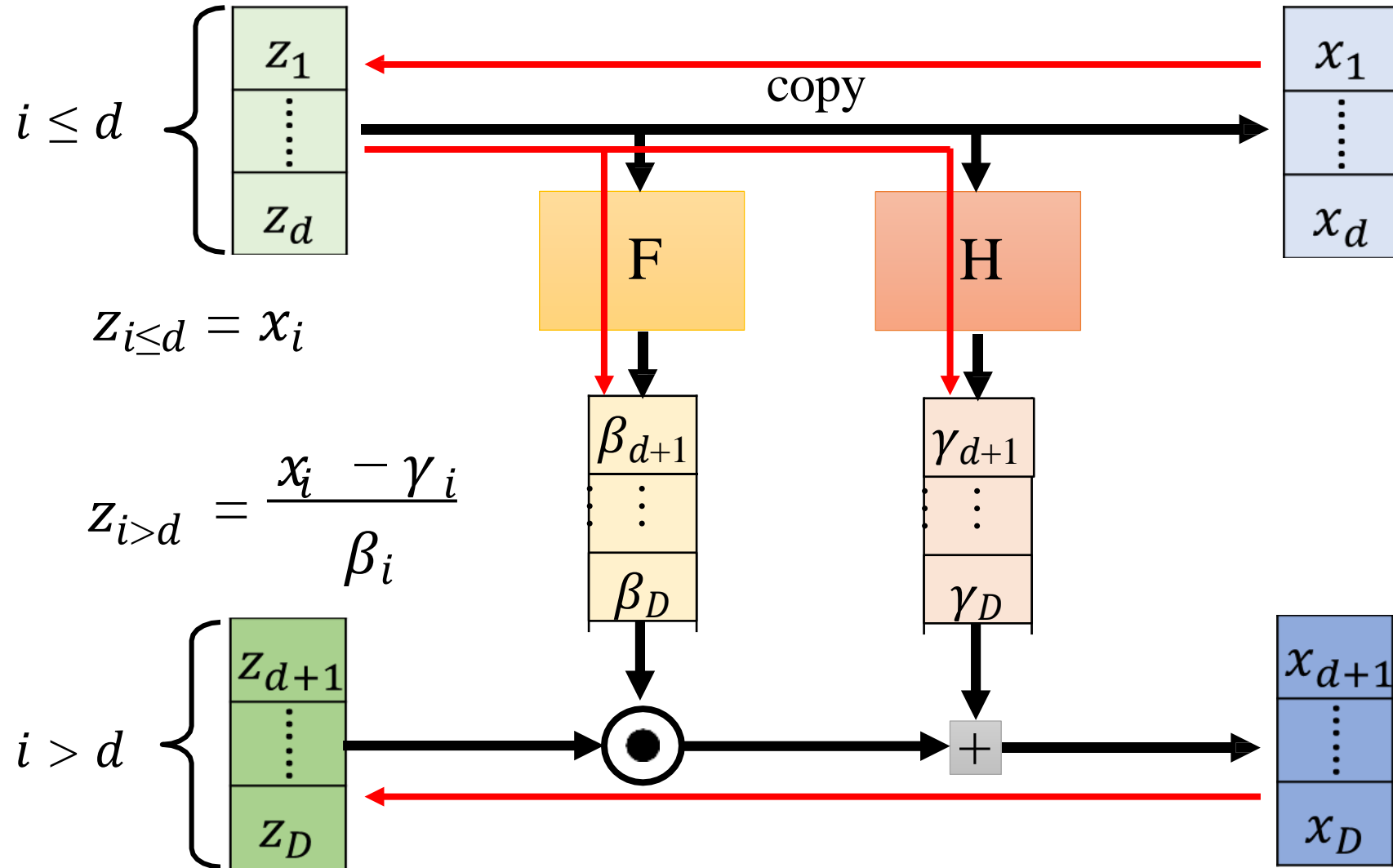
Coupling Layer

NICE

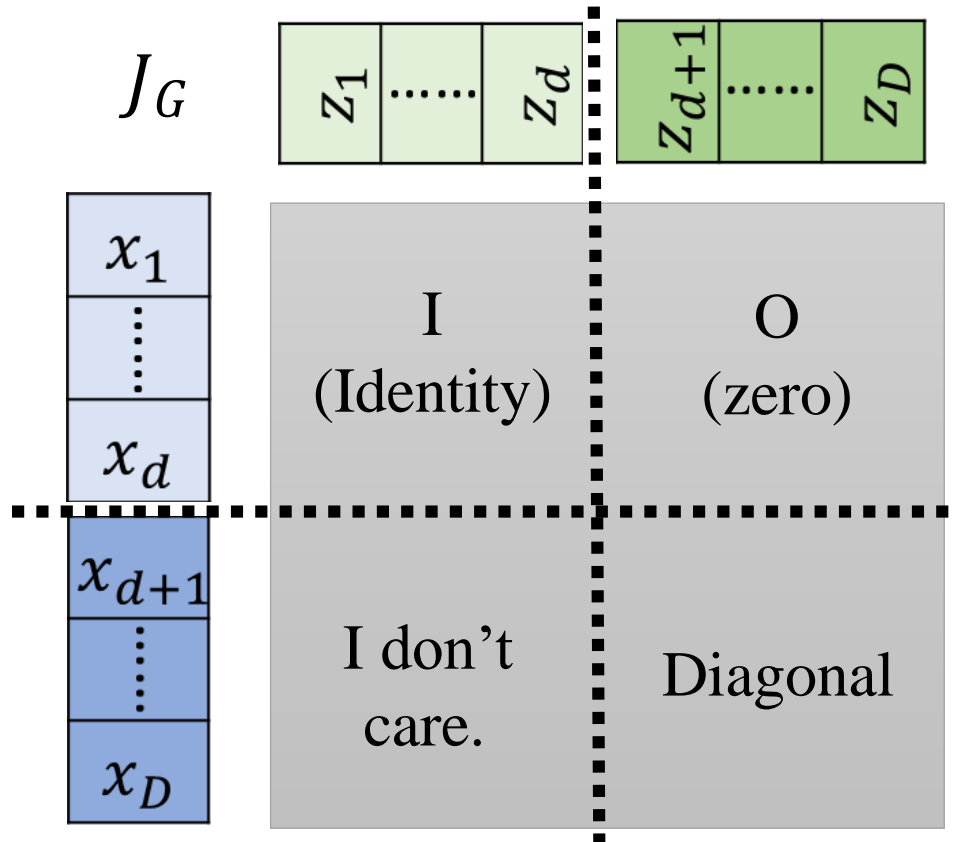
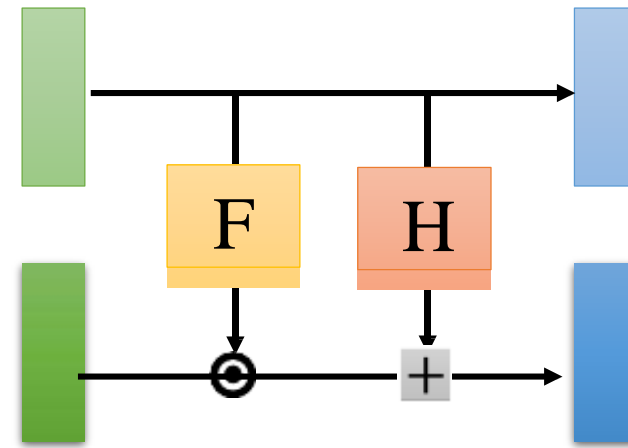
<https://arxiv.org/abs/1410.8516>

Real NVP

<https://arxiv.org/abs/1605.08803>



Coupling Layer



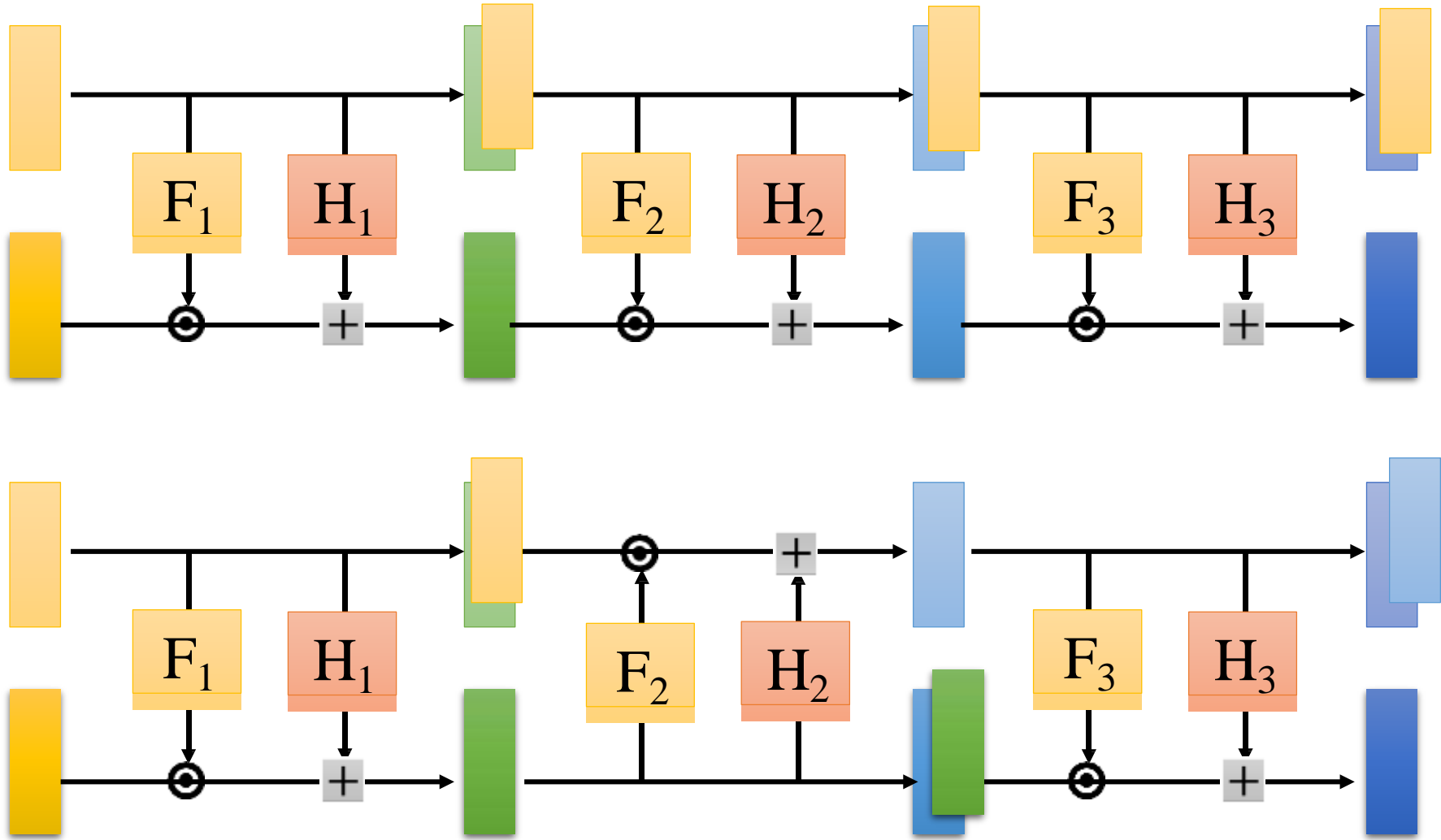
$$\det(J_G)$$

$$= \frac{\partial x_{d+1}}{\partial z_{d+1}} \frac{\partial x_{d+2}}{\partial z_{d+2}} \dots \frac{\partial x_D}{\partial z_D}$$

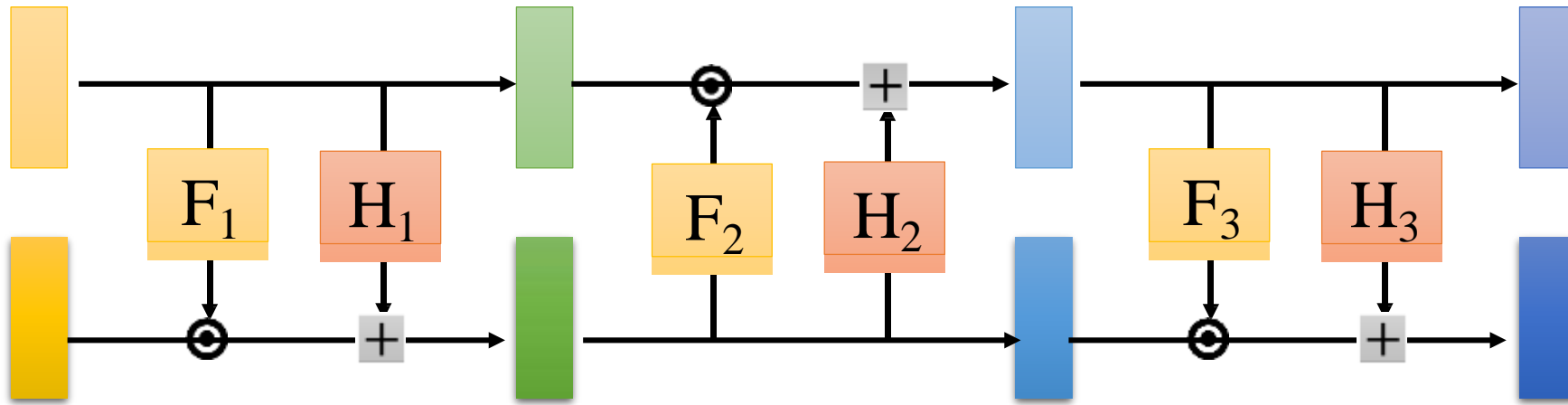
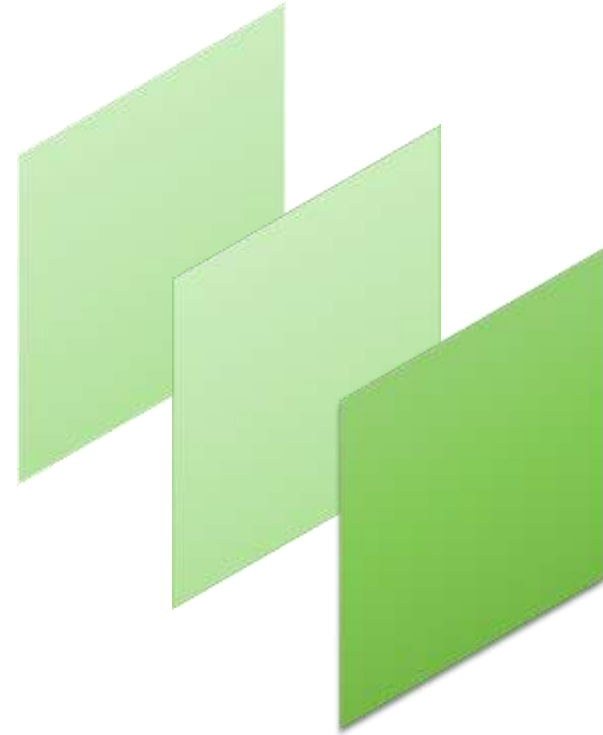
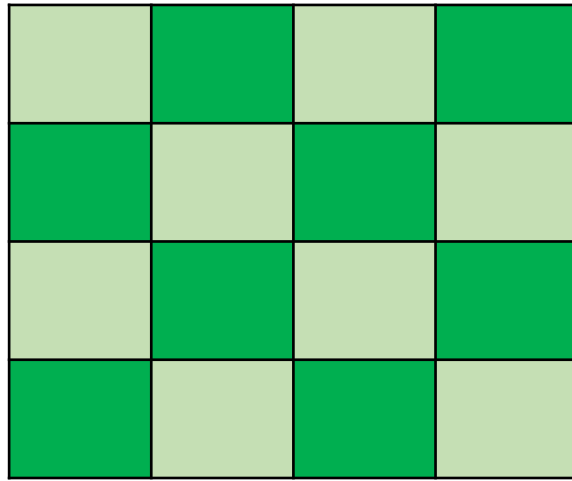
$$= \beta_{d+1} \beta_{d+2} \dots \beta_D$$

$$x_{i>d} = \beta_i z_i + \gamma_i$$

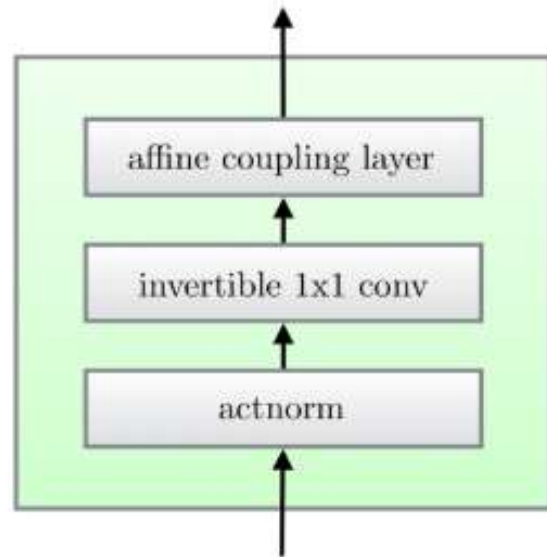
Coupling Layer - Stacking



Coupling Layer

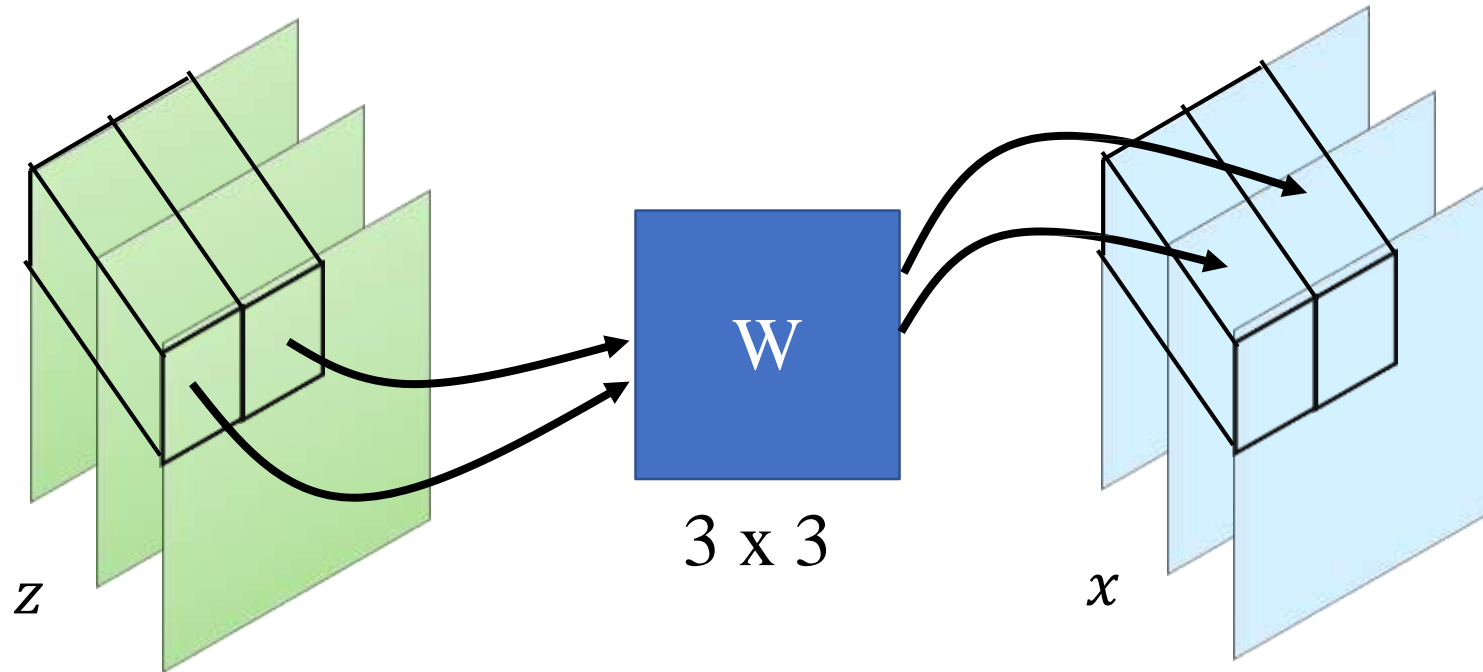


1x1 Convolution



Description	Function	Reverse Function	Log-determinant
Actnorm. See Section 3.1.	$\forall i, j : \mathbf{y}_{i,j} = \mathbf{s} \odot \mathbf{x}_{i,j} + \mathbf{b}$	$\forall i, j : \mathbf{x}_{i,j} = (\mathbf{y}_{i,j} - \mathbf{b})/\mathbf{s}$	$h \cdot w \cdot \text{sum}(\log \mathbf{s})$
Invertible 1×1 convolution. $\mathbf{W} : [c \times c]$. See Section 3.2.	$\forall i, j : \mathbf{y}_{i,j} = \mathbf{W}\mathbf{x}_{i,j}$	$\forall i, j : \mathbf{x}_{i,j} = \mathbf{W}^{-1}\mathbf{y}_{i,j}$	$h \cdot w \cdot \log \det(\mathbf{W}) $ or $h \cdot w \cdot \text{sum}(\log \mathbf{s})$ (see eq. (10))
Affine coupling layer. See Section 3.3 and (Dinh et al., 2014)	$\mathbf{x}_a, \mathbf{x}_b = \text{split}(\mathbf{x})$ $(\log \mathbf{s}, \mathbf{t}) = \text{NN}(\mathbf{x}_b)$ $\mathbf{s} = \exp(\log \mathbf{s})$ $\mathbf{y}_a = \mathbf{s} \odot \mathbf{x}_a + \mathbf{t}$ $\mathbf{y}_b = \mathbf{x}_b$ $\mathbf{y} = \text{concat}(\mathbf{y}_a, \mathbf{y}_b)$	$\mathbf{y}_a, \mathbf{y}_b = \text{split}(\mathbf{y})$ $(\log \mathbf{s}, \mathbf{t}) = \text{NN}(\mathbf{y}_b)$ $\mathbf{s} = \exp(\log \mathbf{s})$ $\mathbf{x}_a = (\mathbf{y}_a - \mathbf{t})/\mathbf{s}$ $\mathbf{x}_b = \mathbf{y}_b$ $\mathbf{x} = \text{concat}(\mathbf{x}_a, \mathbf{x}_b)$	$\text{sum}(\log(\mathbf{s}))$

1x1 Convolution



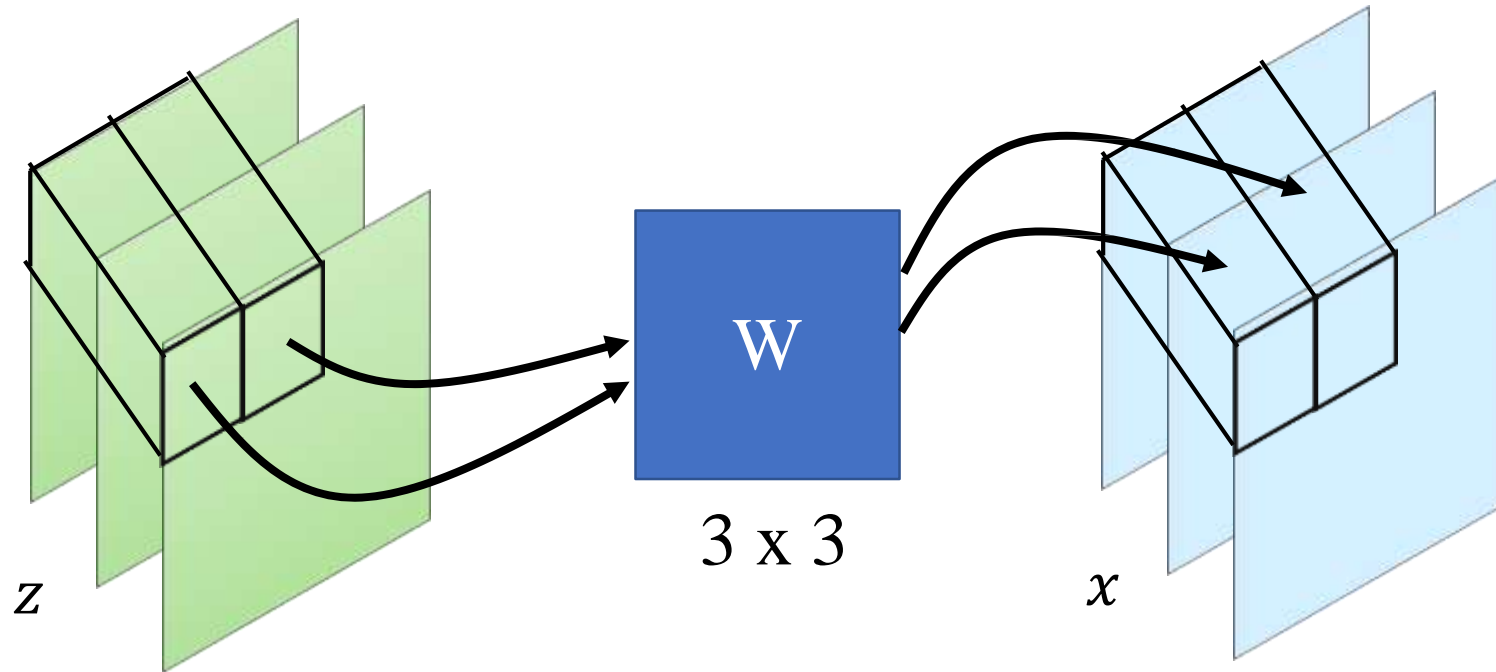
W can shuffle the **channels**.

If W is invertible, it is easy to compute W^{-1} .

$$\begin{array}{|c|} \hline 3 \\ \hline 1 \\ \hline 2 \\ \hline \end{array} = \begin{array}{|c|c|c|} \hline 0 & 0 & 1 \\ \hline 1 & 0 & 0 \\ \hline 0 & 1 & 0 \\ \hline \end{array} \begin{array}{|c|} \hline 1 \\ \hline 2 \\ \hline 3 \\ \hline \end{array}$$

1x1 Convolution

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ w_{31} & w_{32} & w_{33} \end{bmatrix} \begin{bmatrix} z_1 \\ z_2 \\ z_3 \end{bmatrix}$$



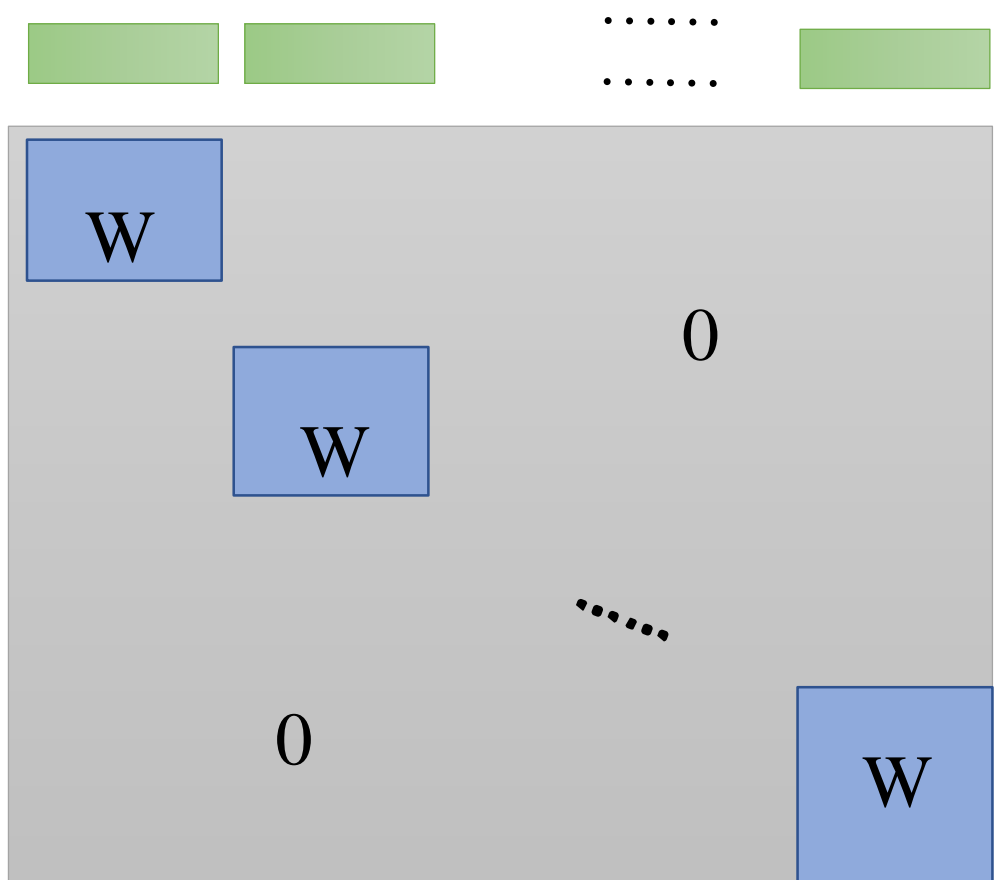
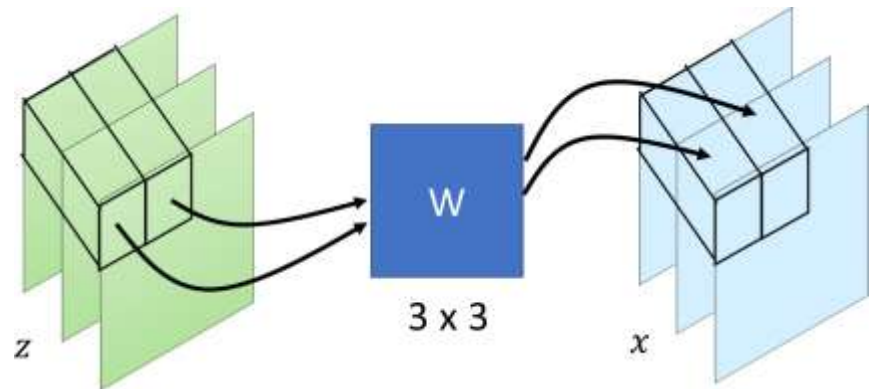
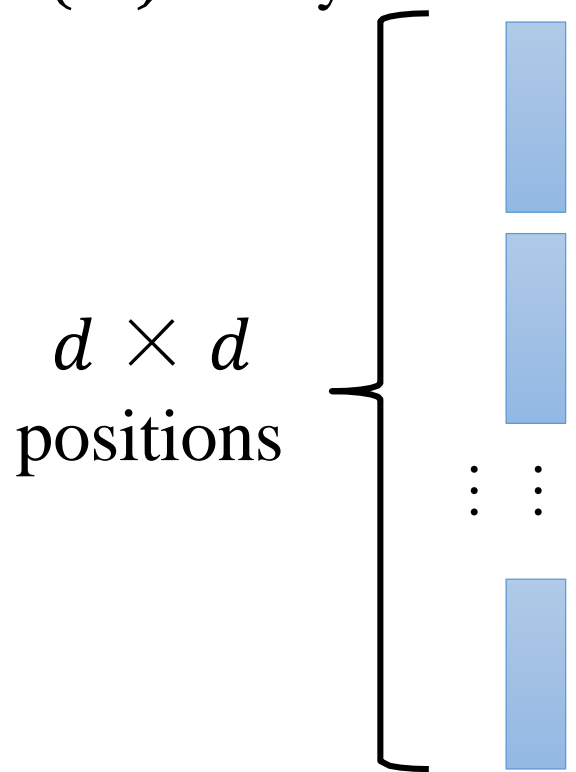
$$x = f(z) = Wz$$

$$J_f = \begin{bmatrix} \partial x_1 / \partial z_1 & \partial x_1 / \partial z_2 & \partial x_1 / \partial z_3 \\ \partial x_2 / \partial z_1 & \partial x_2 / \partial z_2 & \partial x_2 / \partial z_3 \\ \partial x_3 / \partial z_1 & \partial x_3 / \partial z_2 & \partial x_3 / \partial z_3 \end{bmatrix} = \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ w_{31} & w_{32} & w_{33} \end{bmatrix} = W$$

1x1 Convolution

$$(\det(W))^{d \times d}$$

If W is 3x3, computing $\det(W)$ is easy.



Demo of OpenAI

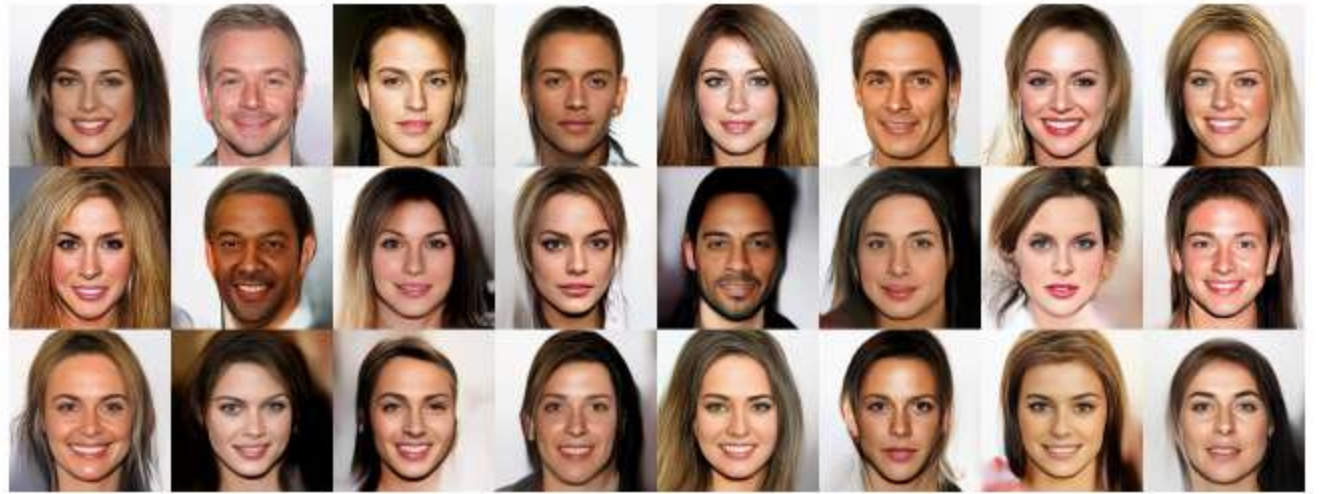


Figure 4: Random samples from the model, with temperature 0.7.

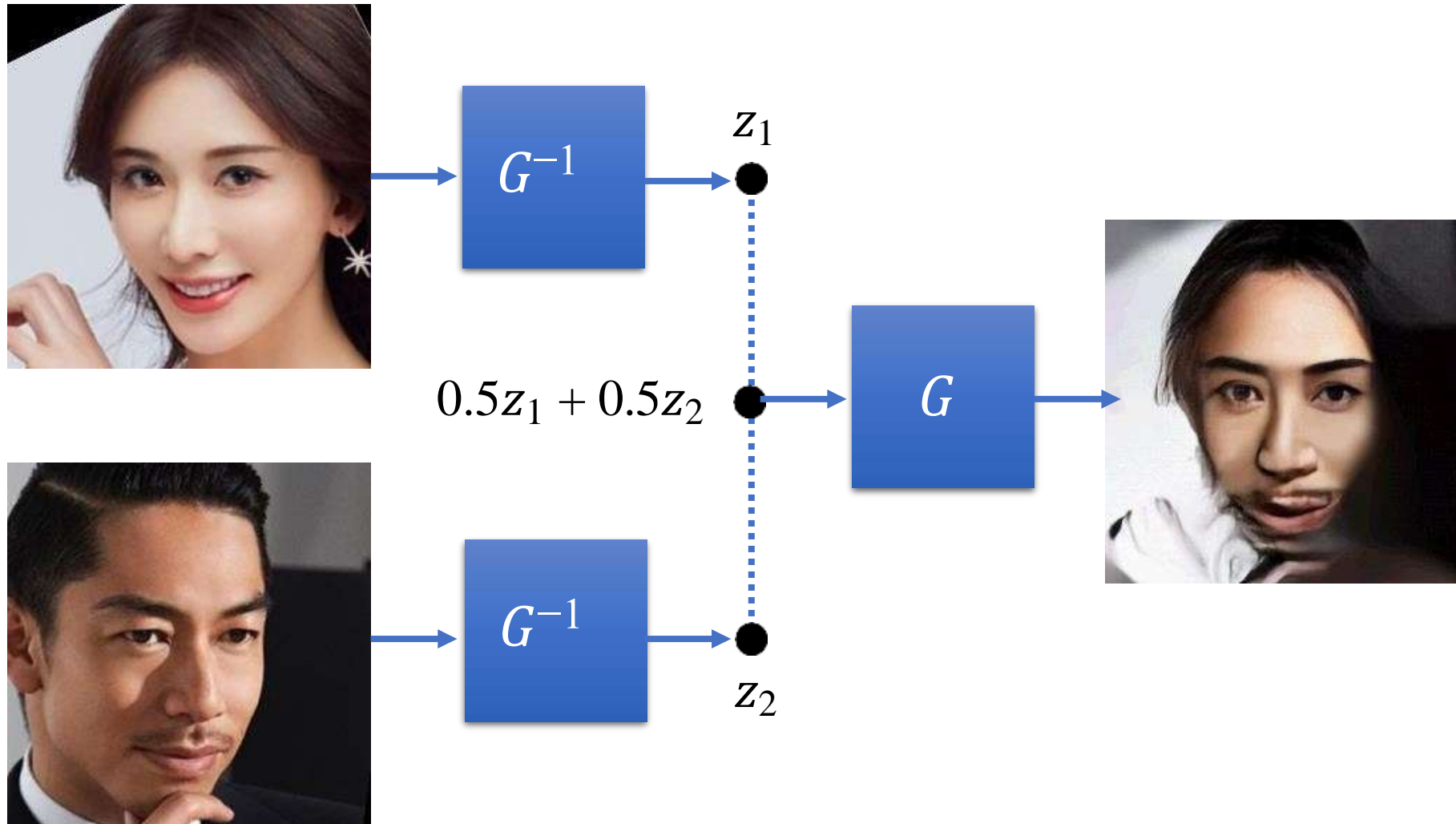


Figure 5: Linear interpolation in latent space between real images.

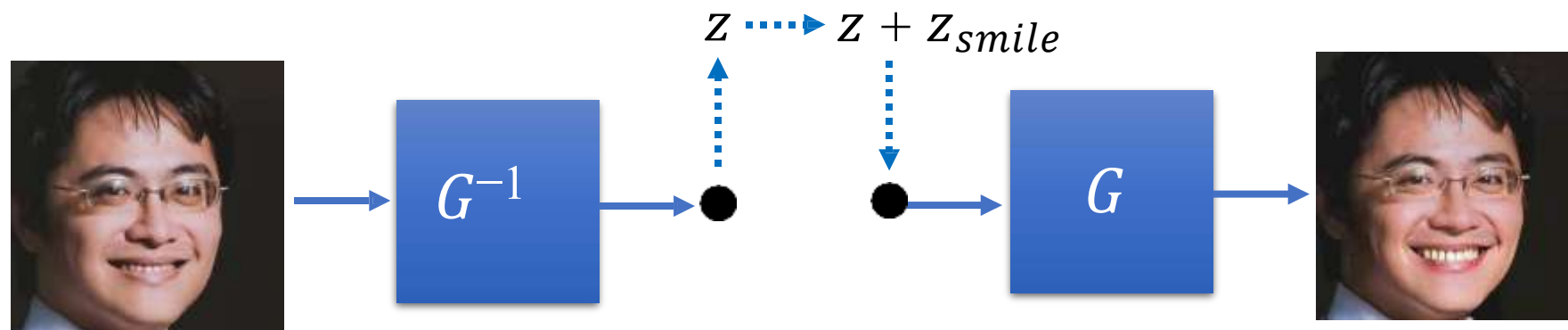
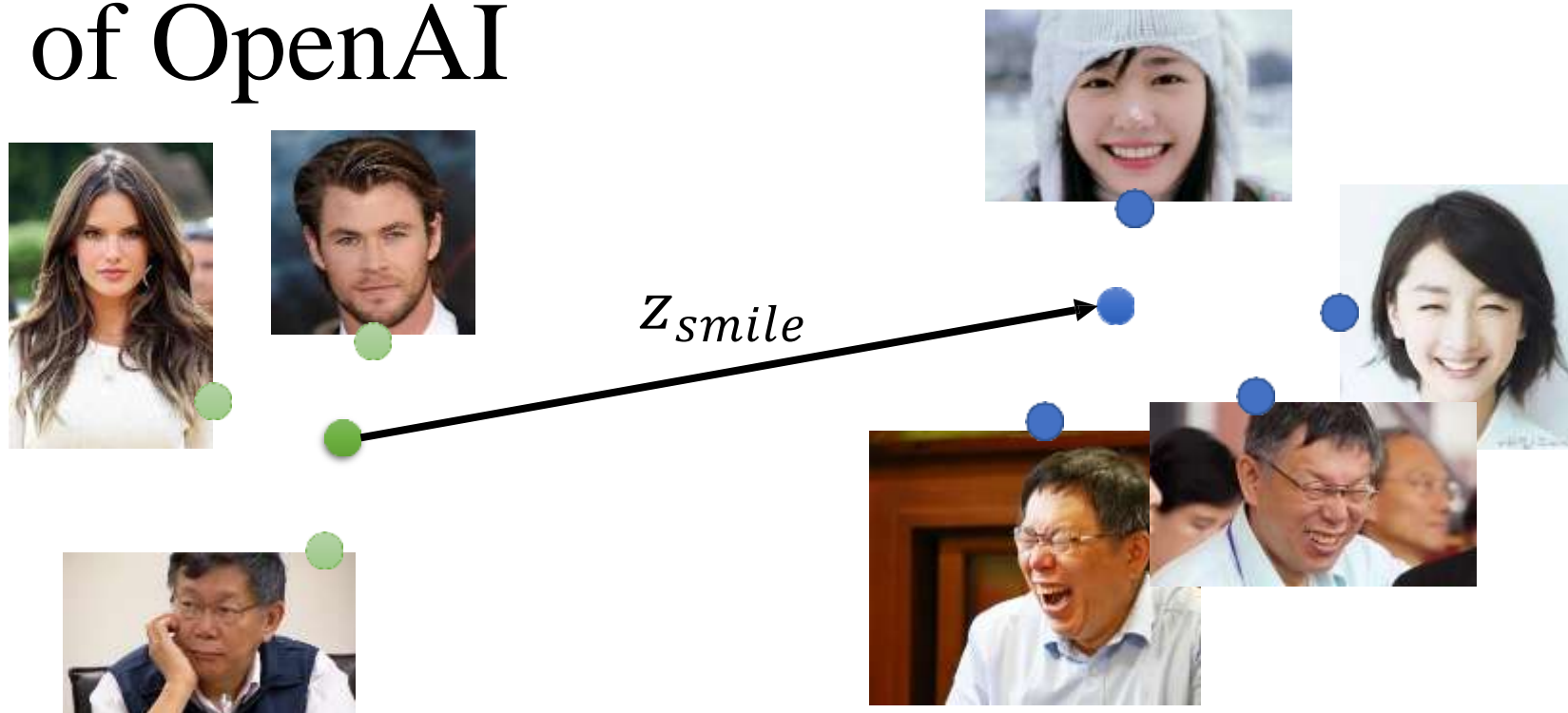
Source of image:

<https://hd.stheadline.com/life/ent/realtime/1517562/>

Demo of OpenAI



Demo of OpenAI



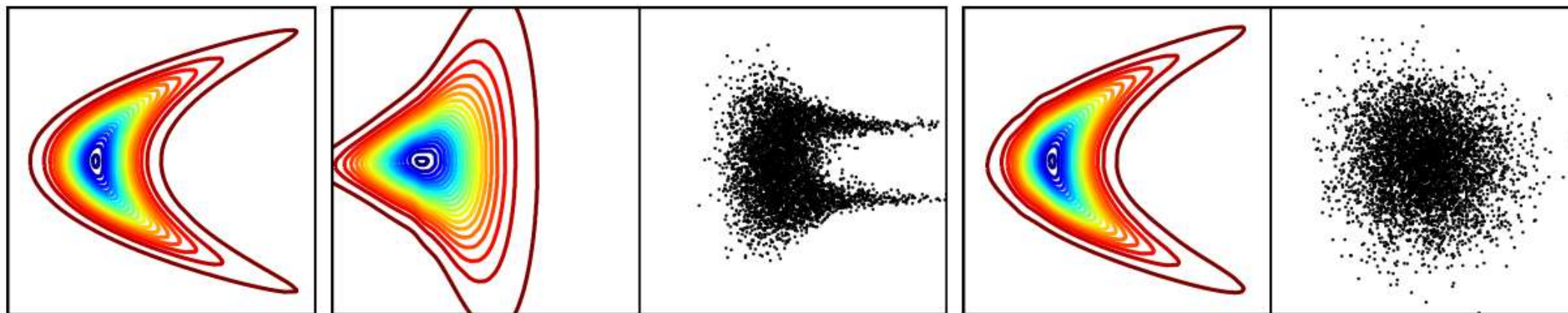
Masked Autoregressive Flow for Density Estimation

- Data generation, producing a new \mathbf{x} :

$$x_i \sim p(x_i | \mathbf{x}_{1:i-1}) = z_i \odot \sigma_i(\mathbf{x}_{1:i-1}) + \mu_i(\mathbf{x}_{1:i-1}), \text{ where } \mathbf{z} \sim \pi(\mathbf{z})$$

- Density estimation, given a known \mathbf{x} :

$$p(\mathbf{x}) = \prod_{i=1}^D p(x_i | \mathbf{x}_{1:i-1})$$

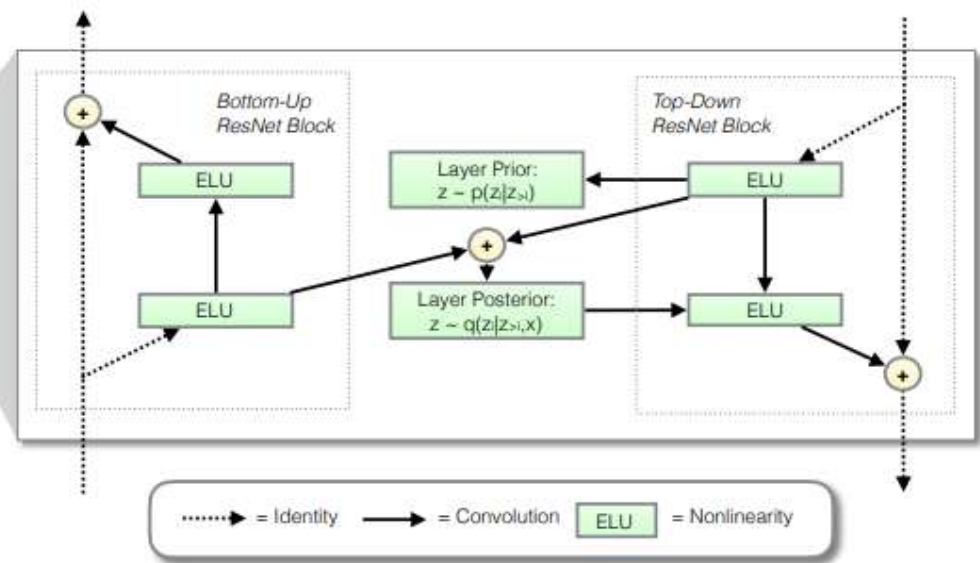
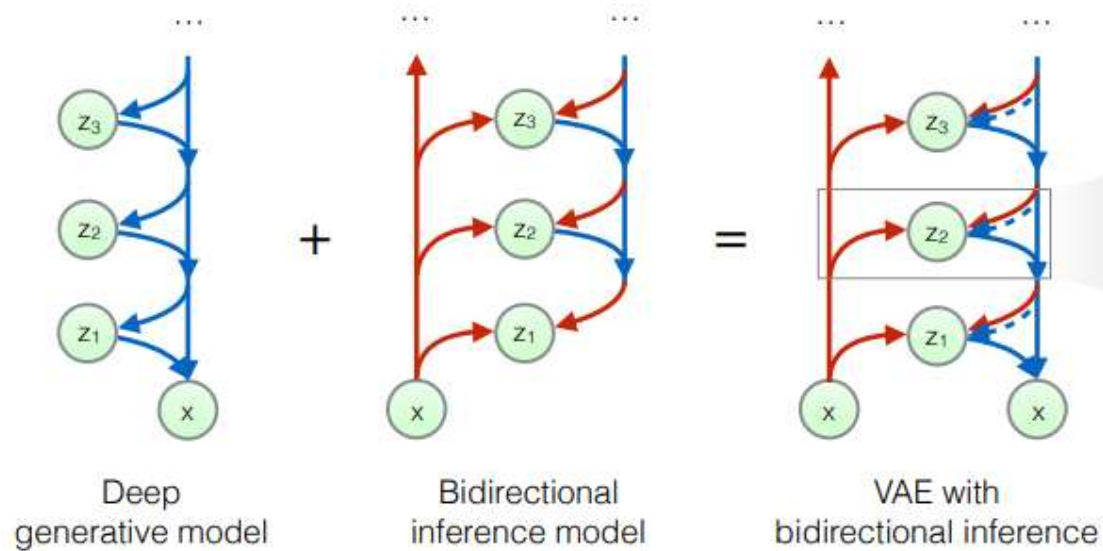
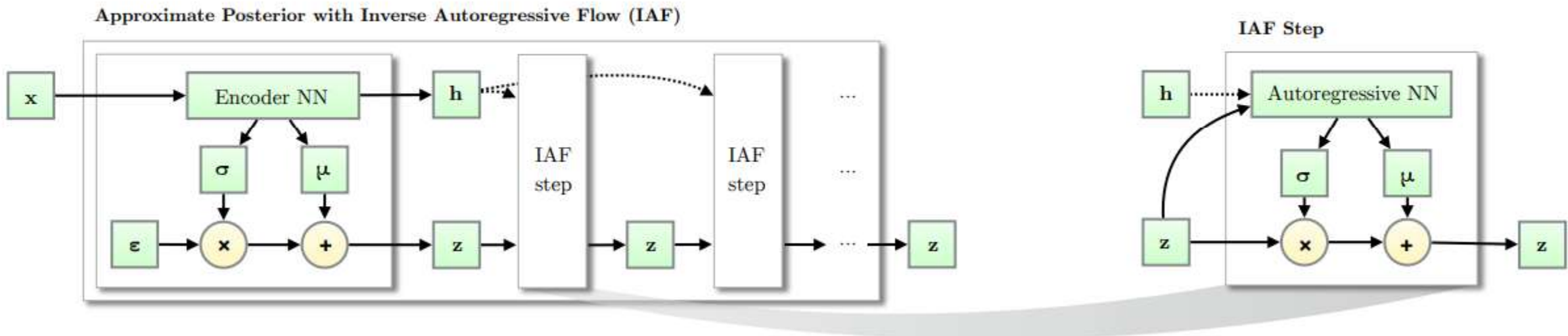


(a) Target density

(b) MADE with Gaussian conditionals

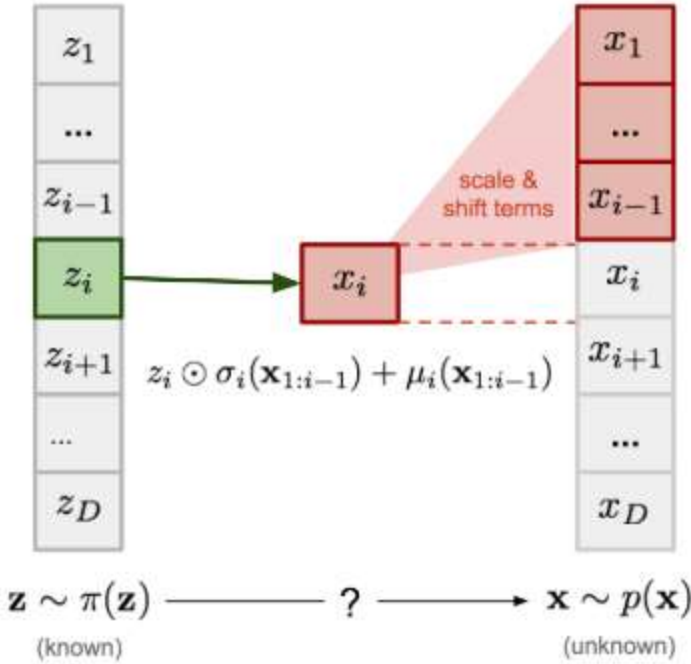
(c) MAF with 5 layers

Improved Variational Inference with Inverse Autoregressive Flow

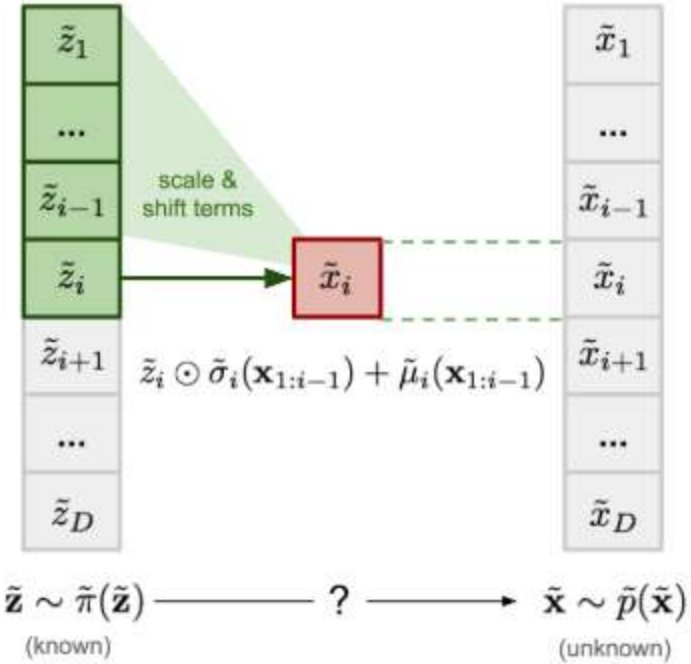


MAF VS IAF

$$z_i = \frac{x_i - \mu_i(\mathbf{x}_{1:i-1})}{\sigma_i(\mathbf{x}_{1:i-1})} = -\frac{\mu_i(\mathbf{x}_{1:i-1})}{\sigma_i(\mathbf{x}_{1:i-1})} + x_i \odot \frac{1}{\sigma_i(\mathbf{x}_{1:i-1})}$$

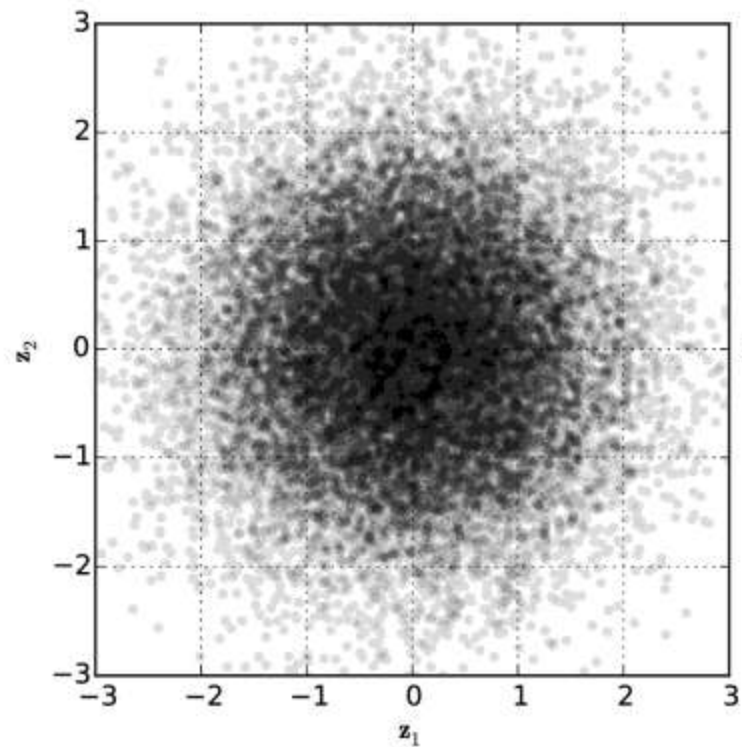


Masked Autoregressive Flow (MAF)

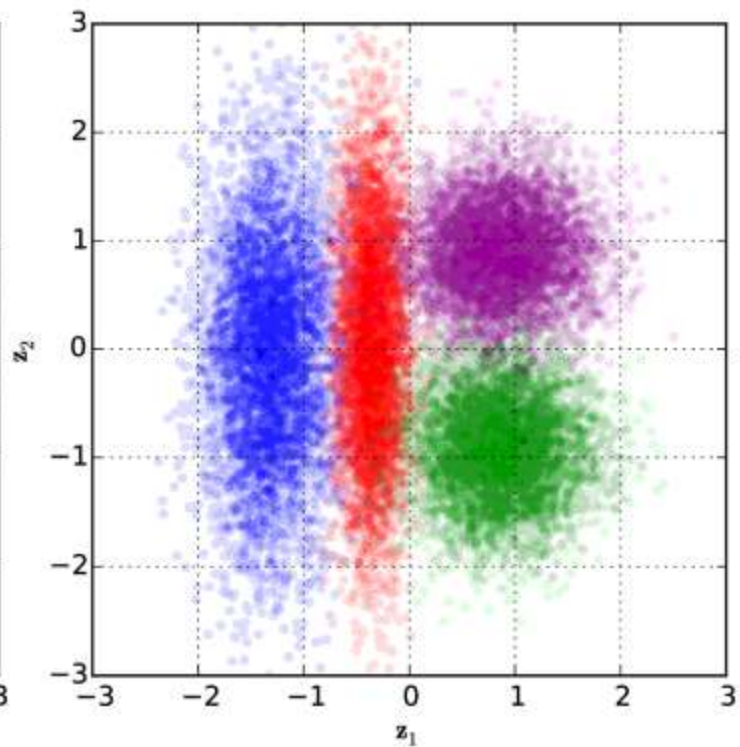


Inverse Autoregressive Flow (IAF)

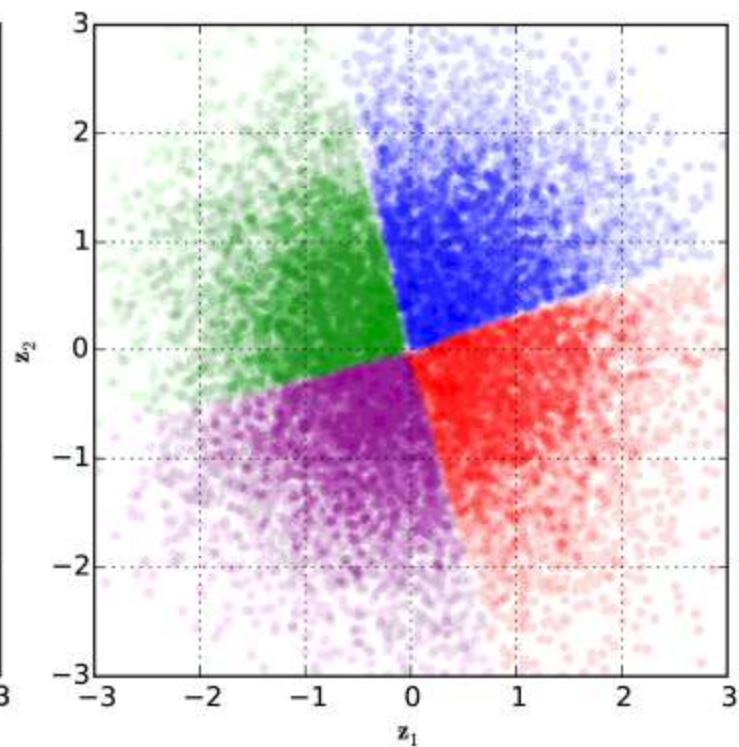
Improved Variational Inference with Inverse Autoregressive Flow



(a) Prior distribution

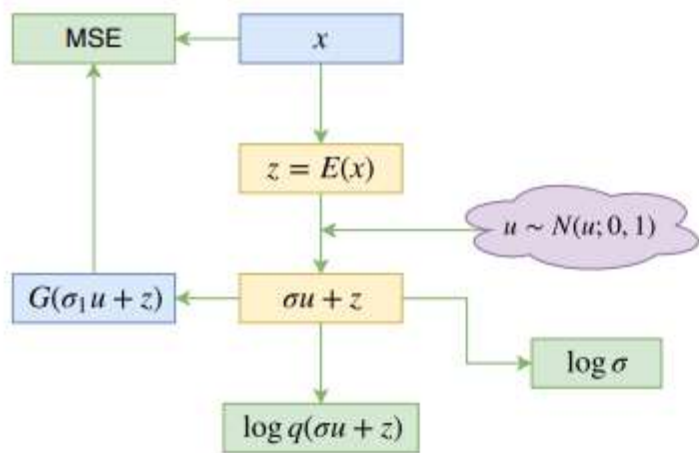


(b) Posteriors in standard VAE

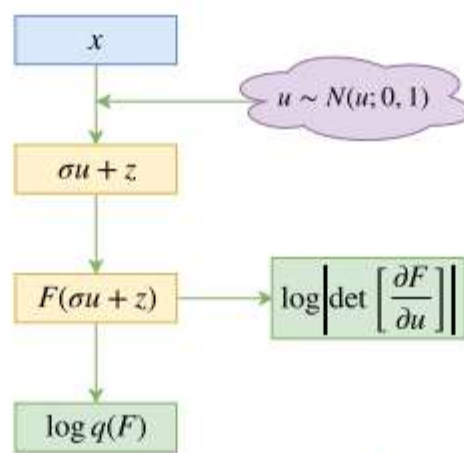


(c) Posteriors in VAE with IAF

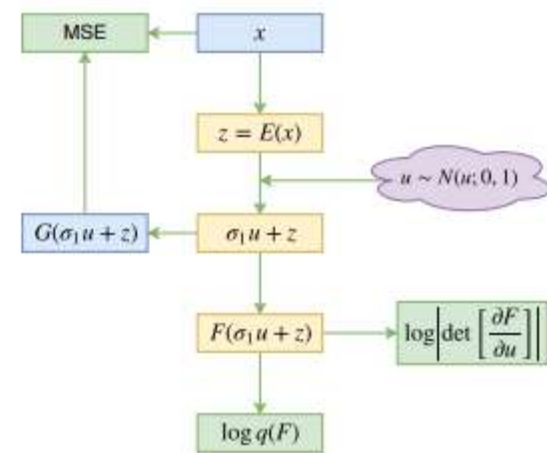
f-VAE



(a) Ordinary VAEs



(b) flow-based models



(c) f-VAEs (ours)

f-VAE



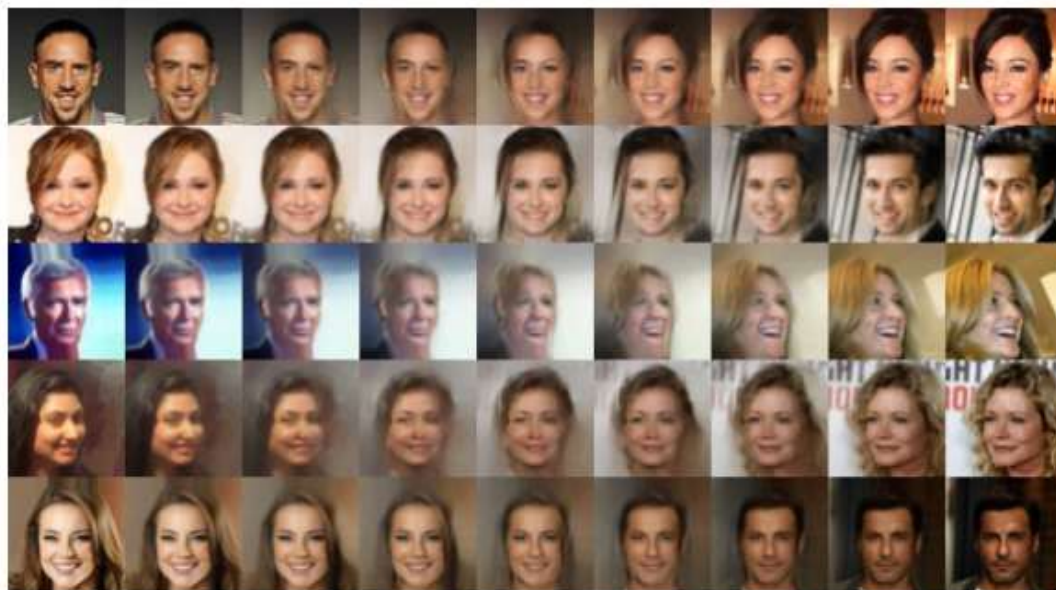
(a) Samples from VAEs



(b) Samples from Glow



(c) Samples from f-VAEs



Variational Inference with Normalizing Flows

