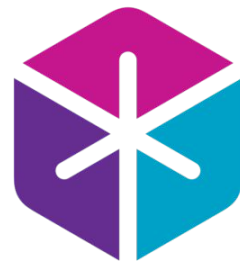# Lecture 03

*Score Matching and Guidance*

**MIT IAP 2026 | Jan 23, 2025**

Peter Holderrieth and Ron Shprints

*Sponsor: Tommi Jaakkola*

# Reminder: Conditional Prob. Path and Cond. Vector Field

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

# Reminder: Marginal Prob. Path and Marginal Vector Field

| | Notation | Key property | Formula |
|---|---|---|---|
| Marginal Probability Path | $p_t$ | Interpolates $p_{\text{init}}$ and $p_{\text{data}}$ | $\int p_t(x|z)p_{\text{data}}(z)\mathrm{d}z$ |
| Marginal Vector Field | $u_t^{\text{target}}(x)$ | ODE follows marginal path | $\int u_t^{\text{target}}(x|z)\dfrac{p_t(x|z)p_{\text{data}}(z)}{p_t(x)}\mathrm{d}z$ |

**Algorithm 3** Flow Matching Training Procedure (General)

---

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

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 $\theta$ via gradient descent on $\mathcal{L}(\theta)$
7: **end for**

---

*We can learn the marginal vector field by approximating the cond. VF for many different data points z.*

# Reminder: Sampling Algorithm for Flow Model

**Algorithm 1** Sampling from a Flow Model with Euler method

**Require:** Neural network vector field $u_t^\theta$, 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** $i = 1, \ldots, n - 1$ **do**
5:    $X_{t+h} = X_t + h u_t^\theta(X_t)$
6:    Update $t \leftarrow t + h$
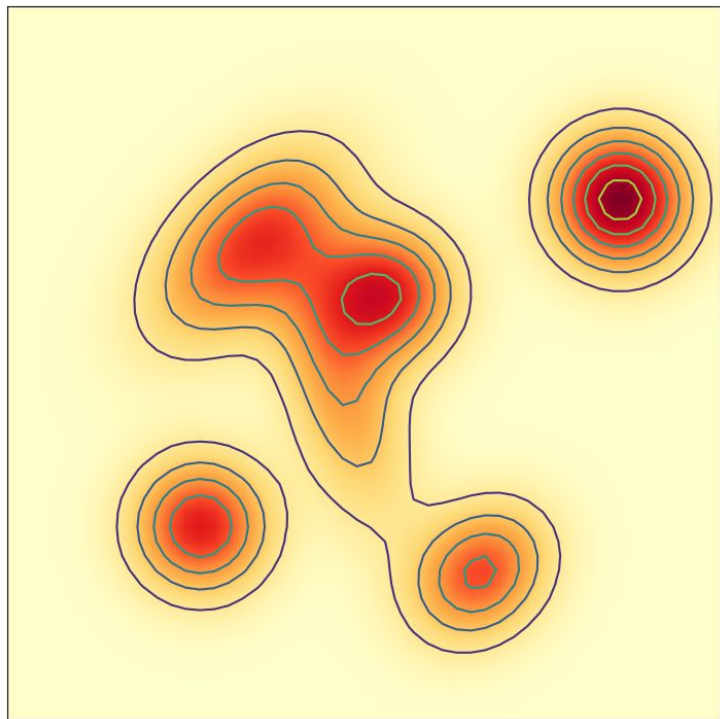7: **end for**
8: **return** $X_1$   *Return final point*

# Section 4:

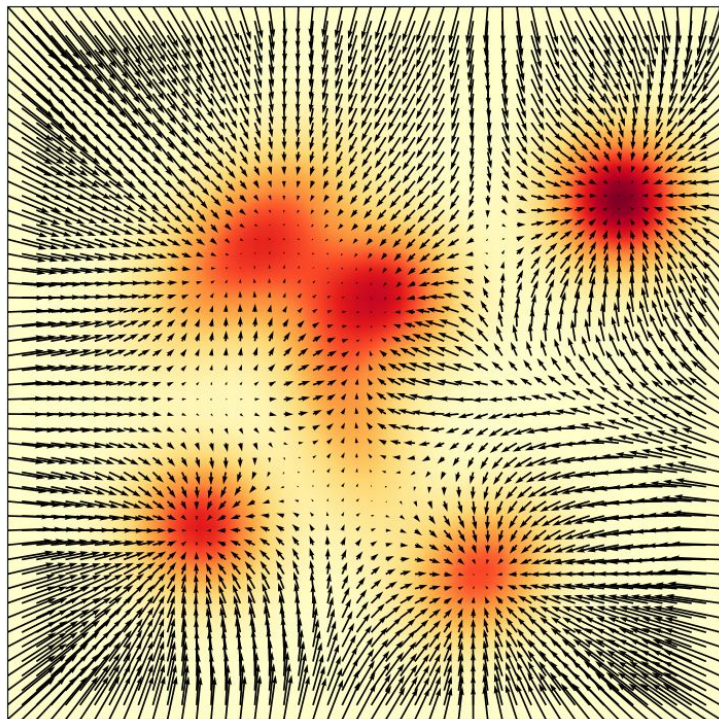# **Training algorithms - 2: Score Matching**

*Goal*: New perspective on flow and diffusion models. SDE/Stochastic Sampling.

# Score Functions = Gradients of the log-likelihood

Log-likelihood: $\log q(x)$

Score function: $\nabla \log q(x)$

# Example - Score of Gaussian Probability Path

$$\nabla \log p_t(x|z) = - \frac{1}{\beta_t^2} x + \frac{\alpha_t}{\beta_t^2} z$$

*Proof:*

$$p_t(x|z) = \mathcal{N}(x; \alpha_t z, \beta_t^2 I_d) = \frac{1}{(2\pi)^{d/2} \beta_t^d} \exp\left(-\frac{1}{2\beta_t^2} \|x - \alpha_t z\|^2\right)$$

$$\log p_t(x|z) = \log \mathcal{N}(x; \alpha_t z, \beta_t^2 I_d) = -\frac{d}{2} \log(2\pi) - d \log \beta_t - \frac{1}{2\beta_t^2} \|x - \alpha_t z\|^2$$

$$\nabla \log p_t(x|z) = \nabla \log \mathcal{N}(x; \alpha_t z, \beta_t^2 I_d) = -\frac{x - \alpha_t z}{\beta_t^2}$$

# Conditional Prob. Path, Vector Field, and Score

|  | Notation | Key property | Gaussian example |
|---|---|---|---|
| Conditional Probability Path | $p_t(\cdot\|z)$ | Interpolates $p_{\text{init}}$ and a data point z | $\mathcal{N}(\alpha_t z, \beta_t^2 I_d)$ |
| Conditional Vector Field | $u_t^{\text{target}}(x\|z)$ | ODE follows conditional path | $\left(\dot{\alpha}_t - \dfrac{\dot{\beta}_t}{\beta_t}\alpha_t\right)z + \dfrac{\dot{\beta}_t}{\beta_t}x$ |
| **Conditional Score Function** | $\nabla \log p_t(x\|z)$ | Gradient of log-likelihood | $\dfrac{\alpha_t}{\beta_t^2}z - \dfrac{1}{\beta_t^2}x$ |

# Marginal Prob. Path, Vector Field, and Score

| | Notation | Key property | Formula |
|---|---|---|---|
| Marginal Probability Path | $p_t$ | Interpolates $p_{\text{init}}$ and $p_{\text{data}}$ | $\int p_t(x|z)p_{\text{data}}(z)\mathrm{d}z$ |
| Marginal Vector Field | $u_t^{\text{target}}(x)$ | ODE follows marginal path | $\int u_t^{\text{target}}(x|z)\dfrac{p_t(x|z)p_{\text{data}}(z)}{p_t(x)}\mathrm{d}z$ |
| **Marginal Score Function** | $\nabla \log p_t(x)$ | Can be used to convert ODE target to SDE | $\int \nabla \log p_t(x|z)\dfrac{p_t(x|z)p_{\text{data}}(z)}{p_t(x)}\mathrm{d}z$ |

# Observation: Both Conditional VF and Cond Score are linear functions! Just with different coefficients!

| | | | |
|---|---|---|---|
| Conditional Vector Field | $u_t^{\mathbf{target}}(x\mid z)$ | ODE follows conditional path | $\left(\dot{\alpha}_t - \dfrac{\dot{\beta}_t}{\beta_t}\alpha_t\right)z + \dfrac{\dot{\beta}_t}{\beta_t}x$ |
| **Conditional Score Function** | $\nabla \log p_t(x\mid z)$ | Gradient of log-likelihood | $\dfrac{\alpha_t}{\beta_t^2}z - \dfrac{1}{\beta_t^2}x$ |

# Reparameterization: Velocity Field –> Score Function

$$a_t = \left( \beta_t^2 \frac{\dot{\alpha}_t}{\alpha_t} - \dot{\beta}_t \beta_t \right), \quad b_t = \frac{\dot{\alpha}_t}{\alpha_t}$$

$$u_t^{\text{target}}(x|z) = a_t \nabla \log p_t(x|z) + b_t x$$

$$u_t^{\text{target}}(x) = a_t \nabla \log p_t(x) + b_t x$$

*Proof: Algebra. Insert formulas. See lecture notes.*

***Early Diffusion Models learnt the score function instead and then just transformed it into the vector field! This is equivalent!***

**Algorithm 6** Score Matching Training Procedure (General)

---

**Require:** A dataset of samples $z \sim p_{\text{data}}$, score network $s_t^\theta$
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

$$\mathcal{L}(\theta) = \|s_t^\theta(x) - \nabla \log p_t(x|z)\|^2$$

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

---

# Denoising Score Matching for Gaussian Prob. Path

$$\nabla \log p_t(x|z) = -\frac{x - \alpha_t z}{\beta_t^2}$$

$$\epsilon \sim \mathcal{N}(0, I_d) \quad \Rightarrow \quad x = \alpha_t z + \beta_t \epsilon \sim \mathcal{N}(\alpha_t z, \beta_t^2 I_d)$$

$$\mathcal{L}_{\mathrm{dsm}}(\theta) = \mathbb{E}_{t \sim \mathrm{Unif}, z \sim p_{\mathrm{data}}, x \sim p_t(\cdot|z)} \left[ \left\| s_t^\theta(x) + \frac{x - \alpha_t z}{\beta_t^2} \right\|^2 \right]$$

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

*Note what the network does: It needs to predict the noise that was used to corrupt the data point! (**DENOISING** diffusion models)*

**Algorithm 5** Score Matching Training Procedure for Gaussian probability path

---

**Require:** A dataset of samples $z \sim p_{\text{data}}$, score network $s_t^\theta$ or noise predictor $\epsilon_t^\theta$
**Require:** Schedulers $\alpha_t, \beta_t$ with $\alpha_0 = \beta_1 = 0, \alpha_1 = \beta_0 = 1$

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 $x_t = \alpha_t z + \beta_t \epsilon$
6:     Compute loss

$$\mathcal{L}(\theta) = \left\| s_t^\theta(x_t) + \frac{\epsilon}{\beta_t} \right\|^2$$

*Numerically unstable for low beta!*

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

# Fokker-Planck equation

*Randomly initialized SDE*

**Given:** $\quad X_0 \sim p_{\mathrm{init}}, \quad \mathrm{d}X_t = u_t(X_t)\mathrm{d}t + \sigma_t \mathrm{d}W_t$

---

**Follow probability path:**

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

*Marginals are p_t*

equivalent

**Fokker-Planck equation holds**

$$\frac{\mathrm{d}}{\mathrm{d}t}p_t(x) = -\mathrm{div}(p_t u_t)(x) + \frac{\sigma_t^2}{2}\Delta p_t(x)$$

*Continuity equ.*      *Heat equ.*

# Fokker-Planck Equation

$$\frac{\mathrm{d}}{\mathrm{d}t} p_t(x) = -\mathrm{div}(p_t u_t)(x)$$

*Change of probability mass at x*

$$+ \frac{\sigma_t^2}{2} \Delta p_t(x)$$

*Heat dispersion*



$u_t(x)$

$\sigma_t$

*Dispersion*

$x$

# Stochastic Sampling of diffusion models

Choose noise level $\sigma_t$. By "SDE extension trick", we can sample from:

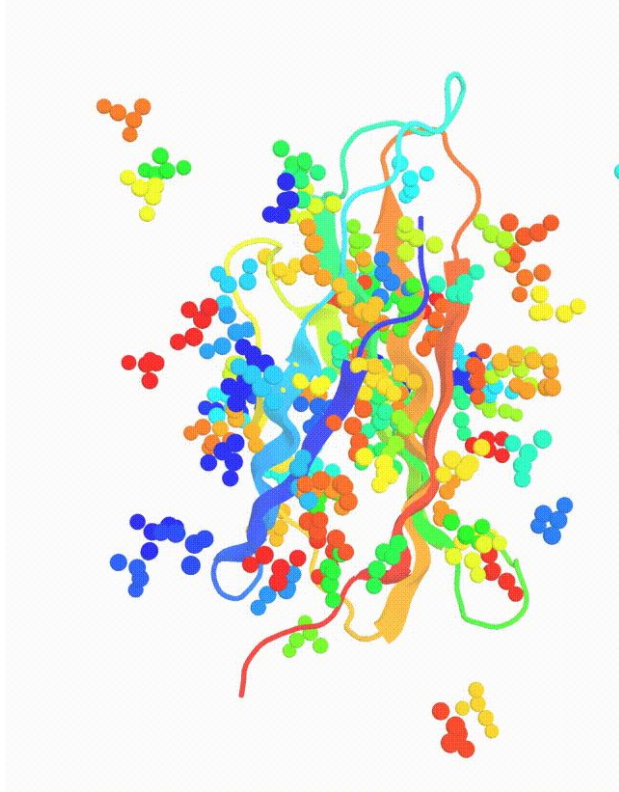$$dX_t = \left[ u_t^{\text{target}}(X_t) + \frac{\sigma_t^2}{2} \nabla \log p_t(X_t) \right] dt + \sigma_t dW_t$$

For Gaussian probability paths, we can express this solely in terms of the score:

$$dX_t = \left[ \left( a_t + \frac{\sigma_t^2}{2} \right) \nabla \log p_t(X_t) + b_t X_t \right] dt + \sigma_t dW_t$$
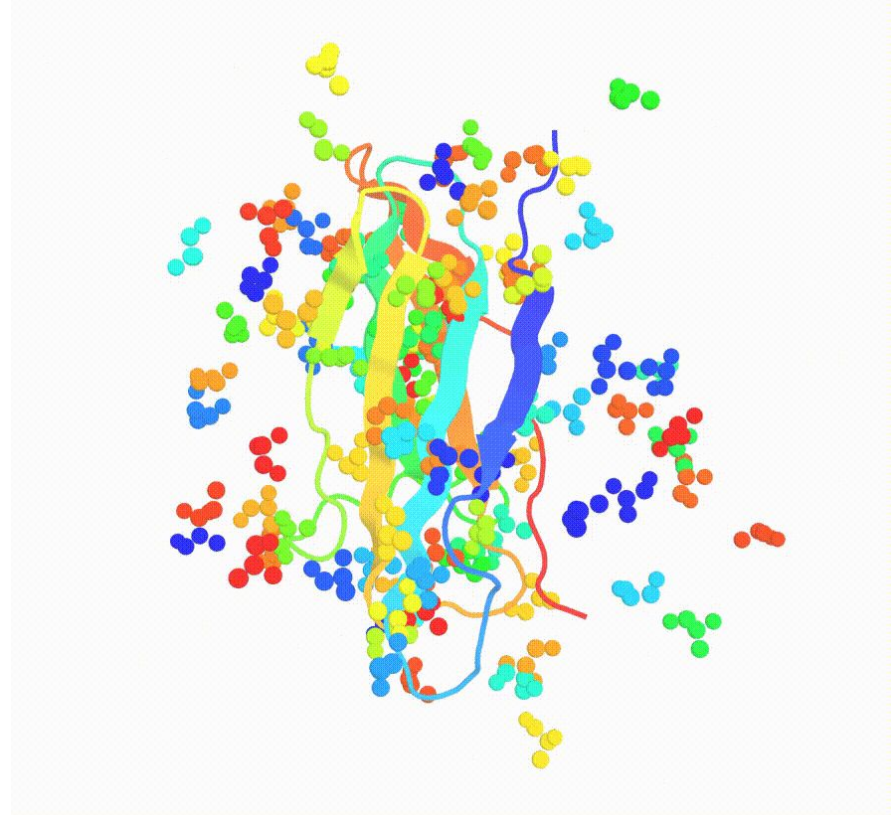
Plugin score network:

$$dX_t = \left[ \left( a_t + \frac{\sigma_t^2}{2} \right) s_t^\theta(X_t) + b_t X_t \right] dt + \sigma_t dW_t$$
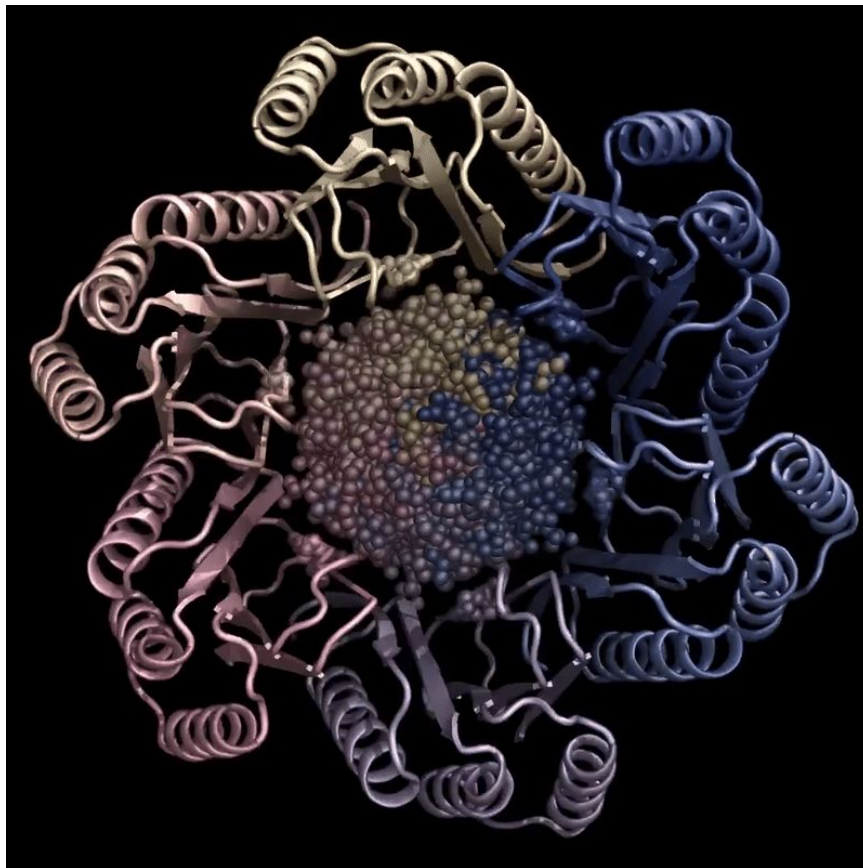
Deterministic Sampling          Stochastic Sampling

# Stochastic (SDE) Sampling with Diffusion Models

Conversion of
of noise into
protein
structure via
SDE
sampling



*Slide credit:*
*Jason Yim*

# Why would we want stochastic/SDE dynamics?

**In theory:** All diffusion coefficients lead to the same result (sample from data distribution).
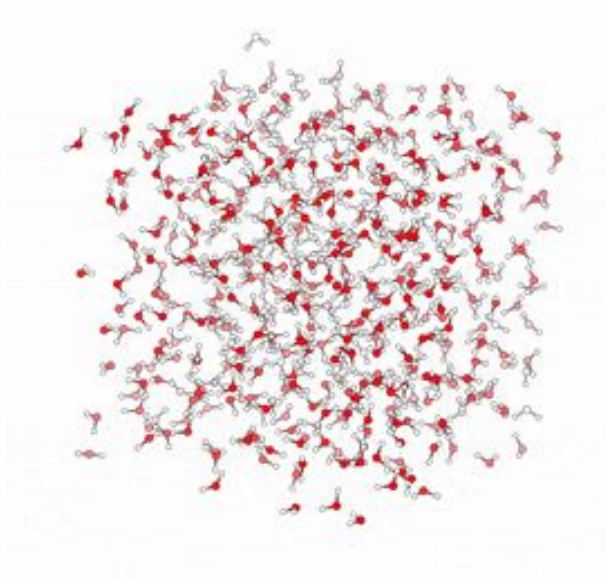
**In practice:**

- **Training error:** Neural network has not perfectly learnt the marginal vector field/score.
- **Simulation error:** We need to simulate SDE/ODE leading to discretization error.

**Downstream applications**: Fine-tuning, inference-time optimization, etc. might require stochastic evolution

**Good news: ODE sampling often leads to the best results. Therefore, SDE sampling is an option, not a must!**

# Aside: Langevin dynamics - Basis of Molecular Dynamics simulation



Molecular dynamics simulate Langevin dynamics. This equals the SDE extension trick for marginal vector field = zero and a constant probability path.

$$\mathrm{d}X_t = \frac{\sigma_t^2}{2} \nabla \log p_t(X_t)\mathrm{d}t + \sigma_t \mathrm{d}W_t$$

$$p_t(x) = p_{\text{Boltzmann}}(x) = \frac{1}{Z} \exp(-U(x))$$

We have shown:

$$X_0 \sim p_{\text{Boltzmann}} \quad \Rightarrow \quad X_t \sim p_{\text{Boltzmann}}$$

*Equilibrium distribution*

# Key takeaway:

- **Conversion formula:** Learning the marginal vector field and learning the score function is equivalent for Gaussian probability paths.

- **Denoising score matching**: Simple way of learning marginal score functions by approximating conditional score functions.

- **Sampling with score models:** Add desired amount of noise + apply correction to vector field

# Section 6:

## **Classifier-free guidance**

*Goal*: Understand how to enforce coherence to prompts

A swamp ogre with a pearl earring by Johannes Vermeer

A car made out of vegetables.

heat death of the universe, line art

**Unguided**: "Generate an image."

**Guided:** "Generate an image of a cat baking a cake."

## Vanilla Guided Sampling

---

**Algorithm 7** Guided Sampling Procedure

---

**Require:** A trained guided vector field $u_t^\theta(x|y)$.
  1: Select a prompt $y \in \mathcal{Y}$, such as "a cat baking a cake".
  2: Initialize $X_0 \sim p_{\text{init}}$.
  3: Simulate $\mathrm{d}X_t = u_t^\theta(X_t|y)\mathrm{d}t$ from $t = 0$ to $t = 1$.
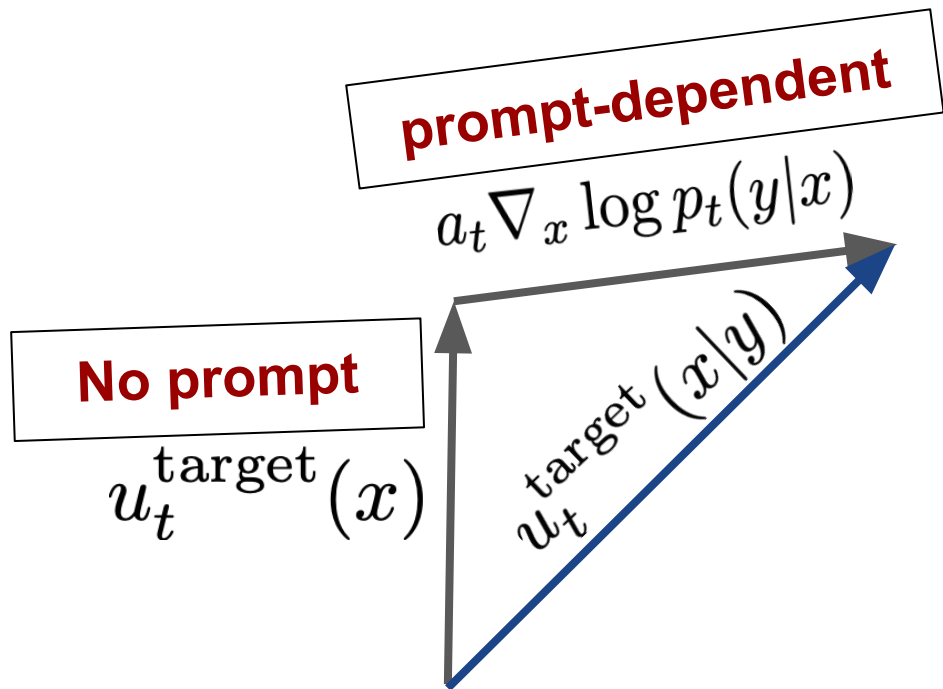
---

# Vanilla Guidance leads to suboptimal results

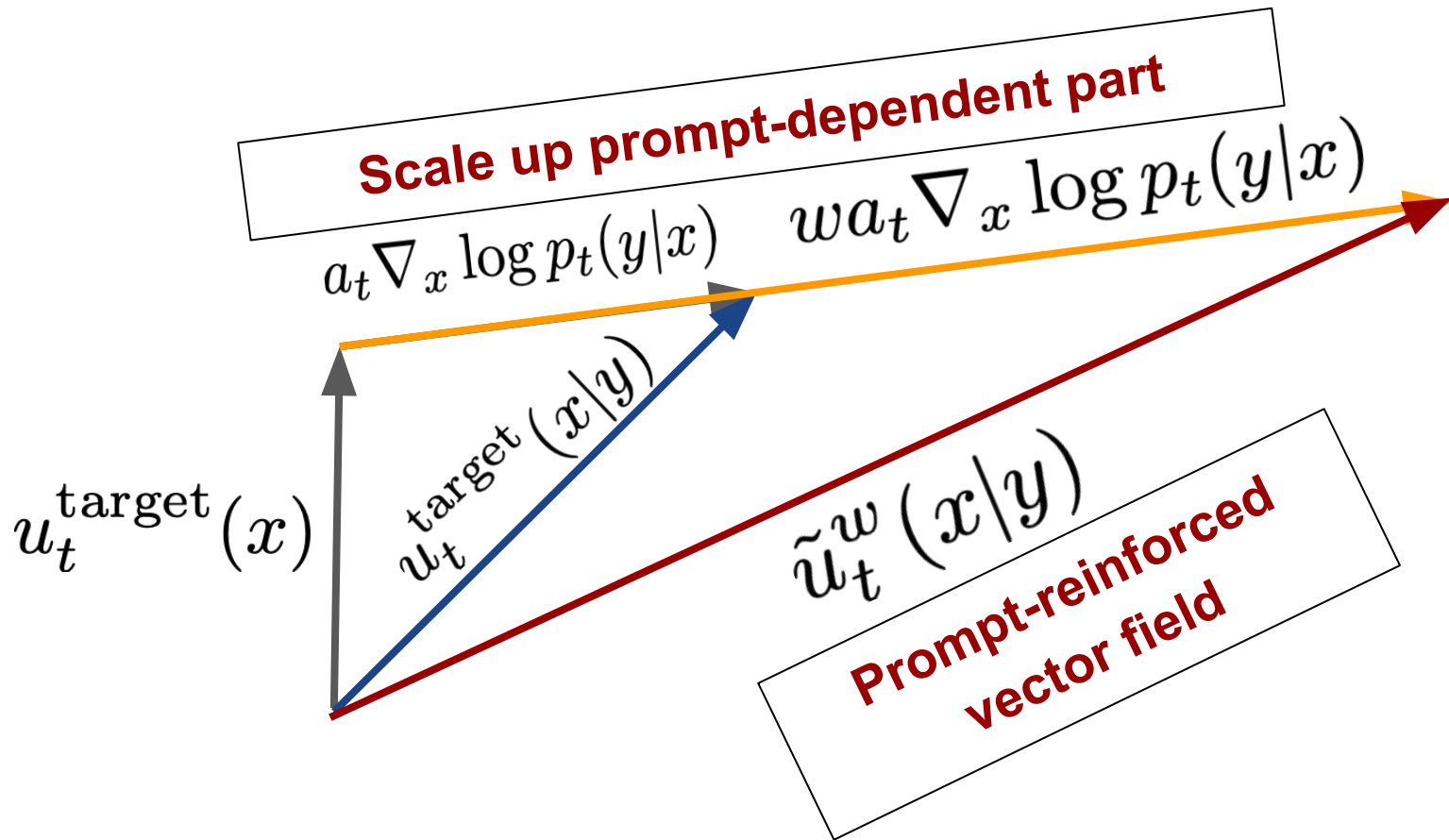*Prompt: "**Corgi dog**"*

***These images do not fit well to the prompt and they have errors!***

# Intuition: Classifier guidance



prompt-dependent

$$a_t \nabla_x \log p_t(y|x)$$

No prompt

$$u_t^{\text{target}}(x)$$

$$u_t^{\text{target}}(x|y)$$

# Intuition: Classifier guidance



Scale up prompt-dependent part

$$a_t \nabla_x \log p_t(y|x) \qquad w a_t \nabla_x \log p_t(y|x)$$

$$u_t^{\text{target}}(x)$$

$$u_t^{\text{target}}(x|y)$$

$$\tilde{u}_t^w(x|y)$$

Prompt-reinforced
vector field

# Classifier-free guidance



$$u_t^{\text{target}}(x|y) - u_t^{\text{target}}(x)$$

$$u_t^{\text{target}}(x)$$

$$u_t^{\text{target}}(x|y)$$

# Classifier-free guidance

# Classifier-free guidance training: Account for empty token $\varnothing$

**Algorithm 5** Classifier-free guidance training

**Require:** Paired dataset $(z, y) \sim p_{\text{data}}$, neural network $u_t^\theta$

1: **for** each mini-batch of data **do**
2:     Sample a data example $(z, y)$ 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 $x = \alpha_t z + \beta_t \epsilon$
6:     With probability $p$ drop label: $y \leftarrow \varnothing$    *Drop label with a certain probability!*
7:     Compute loss

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

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

**Sampling with Classifier-Free Guidance simply is the same as before but we use the weighted vector field:**

$$u_t^{\theta,w}(x) = (1-w)u_t^\theta(x|\varnothing) + wu_t^\theta(x|y)$$

---

**Algorithm 8** Classifier-Free Guidance Sampling Procedure

---

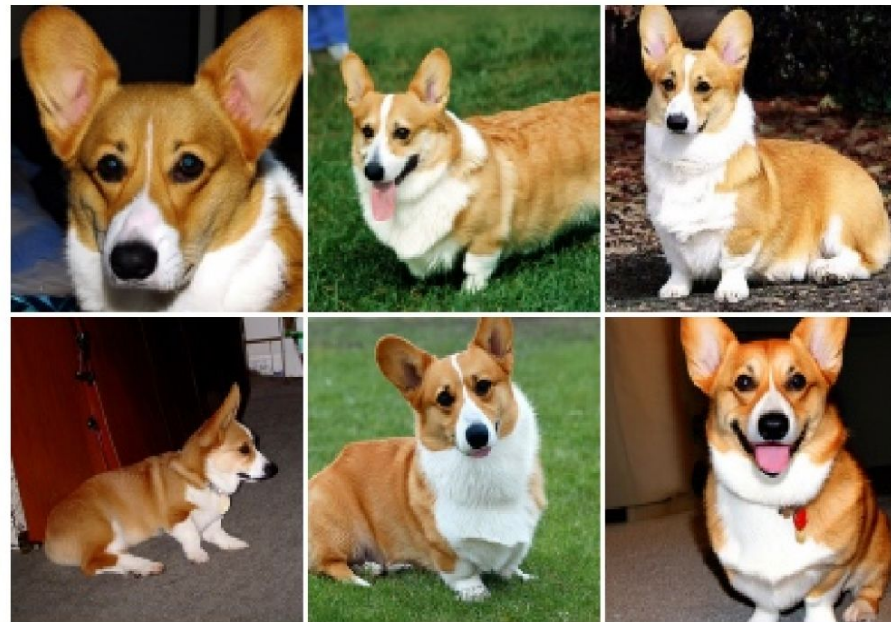**Require:** A trained guided vector field $u_t^\theta(x|y)$.
 1: Select a prompt $y \in \mathcal{Y}$, or take $y = \varnothing$ for unguided sampling.
 2: Select a **guidance scale** $w > 1$.
 3: Initialize $X_0 \sim p_{\text{init}}$.
 4: Simulate $dX_t = \left[(1-w)u_t^\theta(X_t|\varnothing) + wu_t^\theta(X_t|y)\right] dt$ from $t=0$ to $t=1$.

---

# Example: Classifier-Free Guidance

w=1.0

w=4.0

# Example: Classifier-Free Guidance

# Virtually all Images or Videos that you see use CFG!

- **CFG is key:** Without classifier-free guidance (CFG), almost nothing would work.



**Example - Stable Diffusion 3: Classifier-free guidance scale w ~= 4.0**

# CFG does not model the data distribution anymore!

- **CFG is a heuristic:** We do not model the data distribution anymore. In fact, we go beyond it! It is primarily justified by its good empirical results!
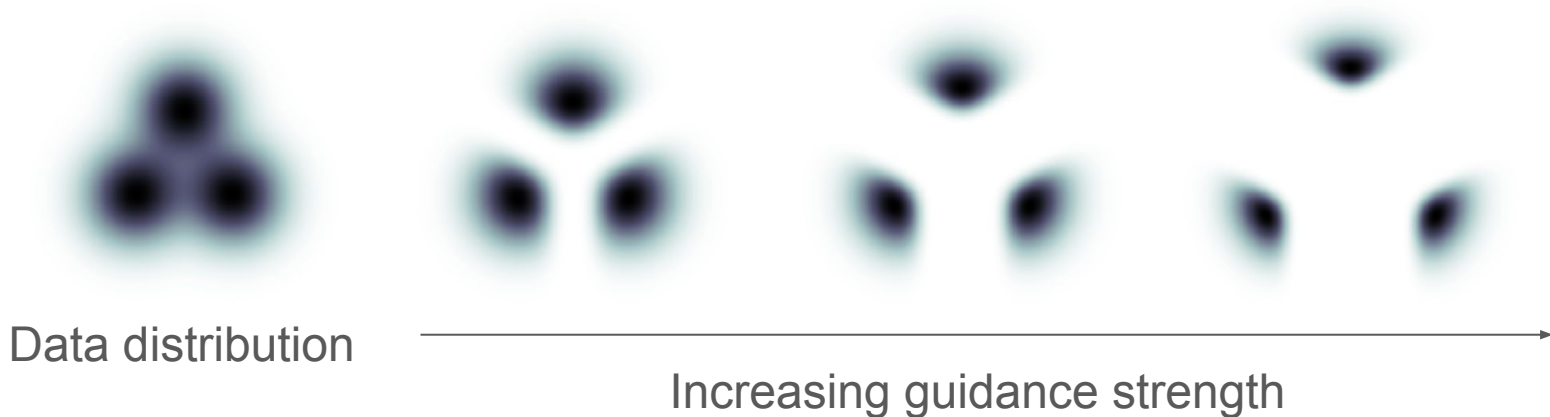
Data distribution

Increasing guidance strength

**Image source**: Classifier-free diffusion guidance [5].

# Section 5:

## **A guide to the diffusion literature**

*Goal*: Understand different interpretations for diffusion models

# Time conventions:

**"Flow time convention":**

- Data: t = 1
- Noise: t = 0

**"Diffusion time convention":**

- Data t = 0
- Noise: t → infinity

**"Discrete time":**

- Use discrete time steps instead of continuous time steps
- No ODE or SDE but Markov chain
- DDIM ~= Probability flow ODE

*Flow matching, rectified flows, stochastic interpolants*

*Score-based diffusion models with SDEs*

*DDPM*
*DDIM*

# Noising Procedure

**Probability path (here):**

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

*Flow matching, rectified flows*

**Interpolant function:**

$$I_t(\epsilon, z) = \alpha_t \epsilon + \beta_t z$$

*Stochastic interpolants*

**"Forward" diffusion process:**

$$dX_t = a_t(X_t)dt + \sigma_t dW_t$$

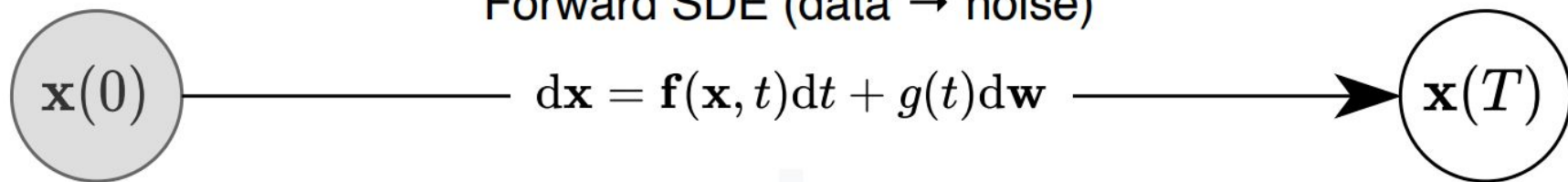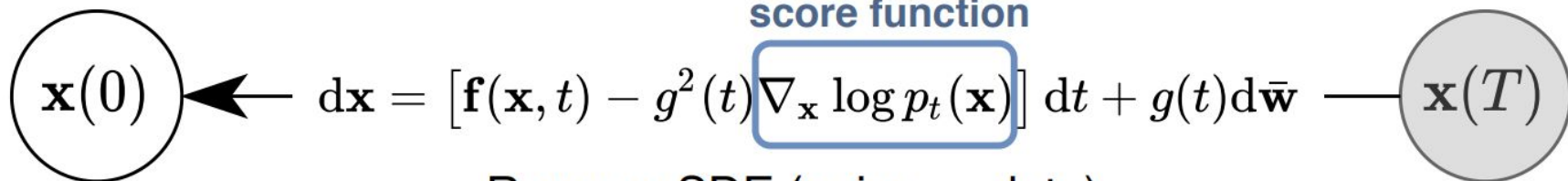*Denoising diffusion models*

Note: For Gaussian probability paths, the above procedures are equivalent.

# Constructing noising procedures via forward noising processes



Forward SDE (data → noise)

$$\mathbf{x}(0) \quad\longrightarrow\quad d\mathbf{x} = \mathbf{f}(\mathbf{x}, t)dt + g(t)d\mathbf{w} \quad\longrightarrow\quad \mathbf{x}(T)$$

score function

$$\mathbf{x}(0) \quad\longleftarrow\quad d\mathbf{x} = \left[\mathbf{f}(\mathbf{x}, t) - g^2(t)\boxed{\nabla_{\mathbf{x}} \log p_t(\mathbf{x})}\right] dt + g(t)d\bar{\mathbf{w}} \quad\longleftarrow\quad \mathbf{x}(T)$$

Reverse SDE (noise → data)

*The time-reversed SDE is a specific solution to the SDE extension trick that we discussed for a specific noise level. Empirically, this is often not the best solution in practice.*

# Here: Flow Matching

- Arguably most simple flow and diffusion algorithms
- Allows you to restrict yourself to flows
- Allows you go from arbitrary $p_{\text{init}}$ to arbitrary $p_{\text{data}}$

**Note: The method presented here allows to convert arbitrary distributions into arbitrary distributions!**

# Bridging arbitrary distributions - Example

Videos without audio → videos with audio

Low resolution images → high resolution images

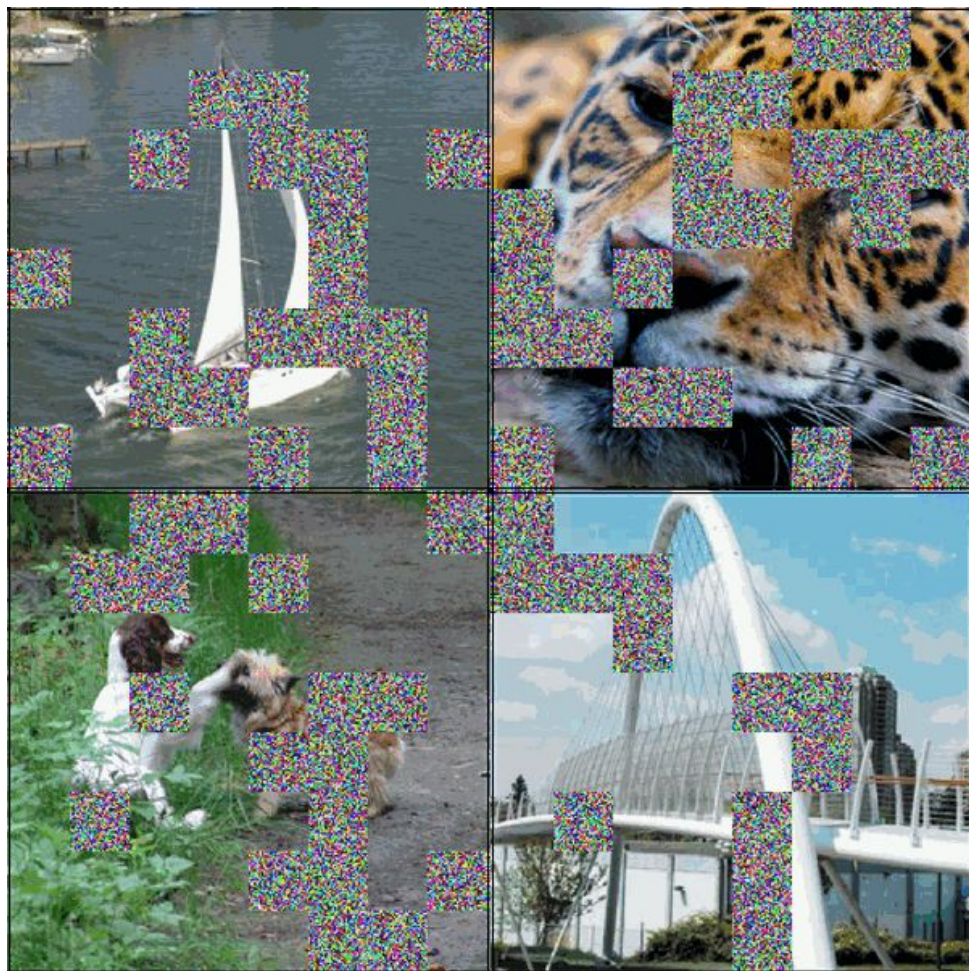Unperturbed cells → perturbed cells

etc.



*Figure credit: Michael Albergo*

# Class Overview

- **Lecture 1 -** Generation as Sampling. Flow and Diffusion Models

- **Lecture 2 - Flow Matching:** Training algorithm.

- **Lecture 3 - Score Matching, Guidance:** How to condition on a prompt.

- **Lecture 4 - Build Image Generators:** Network architectures + Latent spaces

- **Lecture 5 - Advanced Topics:** Discrete diffusion models + distilled models

Next class:

**Monday, 11am-12:30pm**

**Neural network architectures + latent spaces!**

*E25-111 (same room)*

**Office hours: Today, 3pm-4:30pm in 36-156**

# References

1. **Scaling Rectified Flow Transformers for High-Resolution Image Synthesis, https://arxiv.org/abs/2403.03206**
2. **An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale, https://arxiv.org/abs/2010.11929**
3. **Scalable Diffusion Models with Transformers, https://arxiv.org/abs/2212.09748**
4. **High Resolution Image Synthesis with Latent Diffusion Models, https://arxiv.org/abs/2112.10752**
5. **Classifier-Free Diffusion Guidance, https://arxiv.org/abs/2207.12598**