

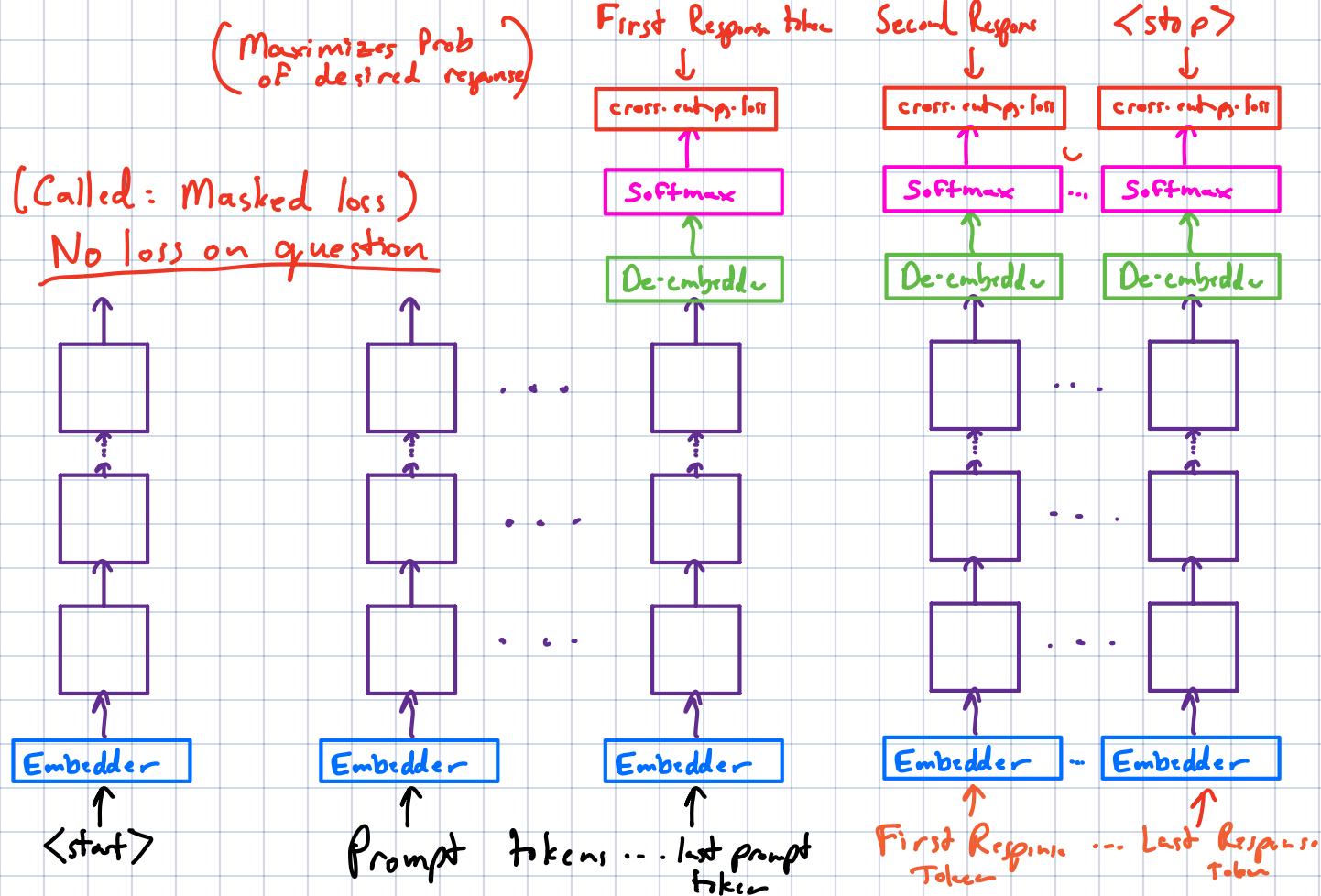
Today: RL Post-training (finish)
Diffusion Models

Readings in Textbook: Ch 17 (VAE)
Ch 19.5 (Some RL...)
Ch 18 (Diffusion)

Announcement: 1) Survey: 40%. 😐
We'll award 15% extra credit
2) Final Project Report: Due Dec 5th Sun AM
Poster Session Tue Dec 9th
Final Peer Review Thu Dec 11th
Final Final Report Sun Dec 14th
3) Final Exam: Wed Dec 17th 8-11AM

Post-training:

Recall basic SFT for instruction following.



But how can we make a model better at solving problems?

Two parts to the answer:

A) Be willing to spend more compute while answering.
Fast time compute

B) Train it to be better.

RLVR : Reinforcement Learning with Verifiable Rewards

Key new paper: Khatri, et al "The Art of Scaling Reinforcement Learning Compute for LLMs", Oct 15, 2025.

Treatment here influenced by this, along with DeepSeek paper.

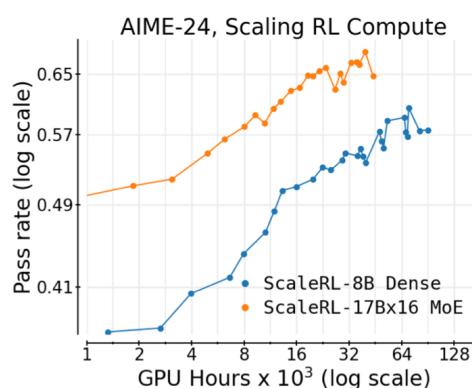
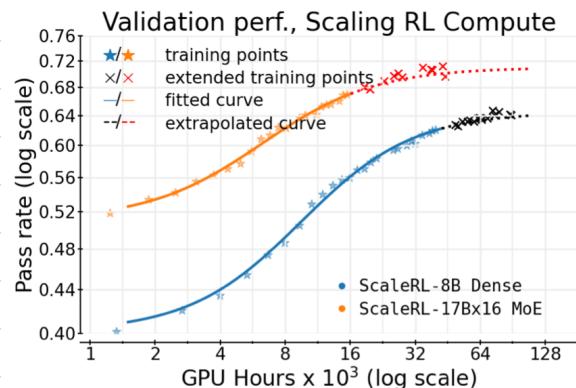


Fig 1 in ScaleRL...

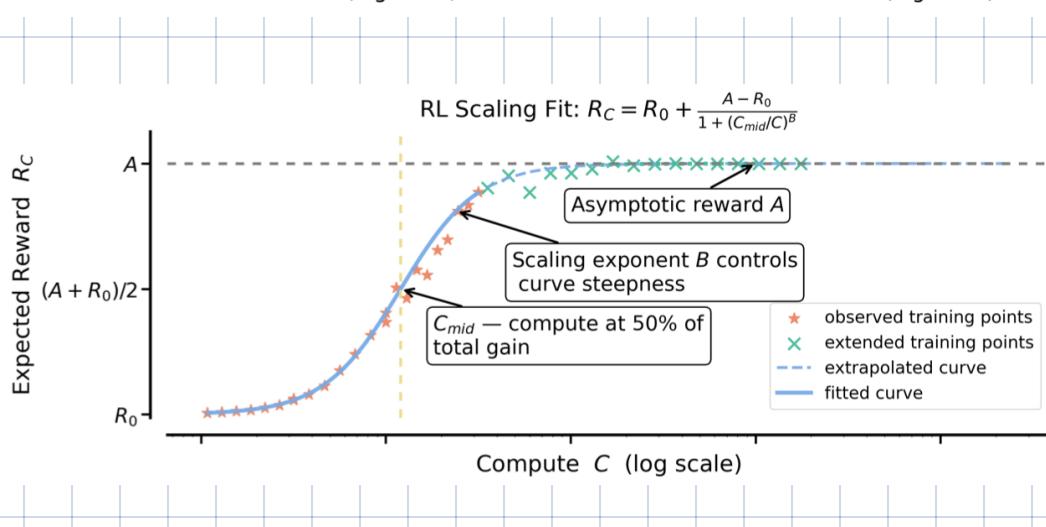


Fig 3 in ScaleRL

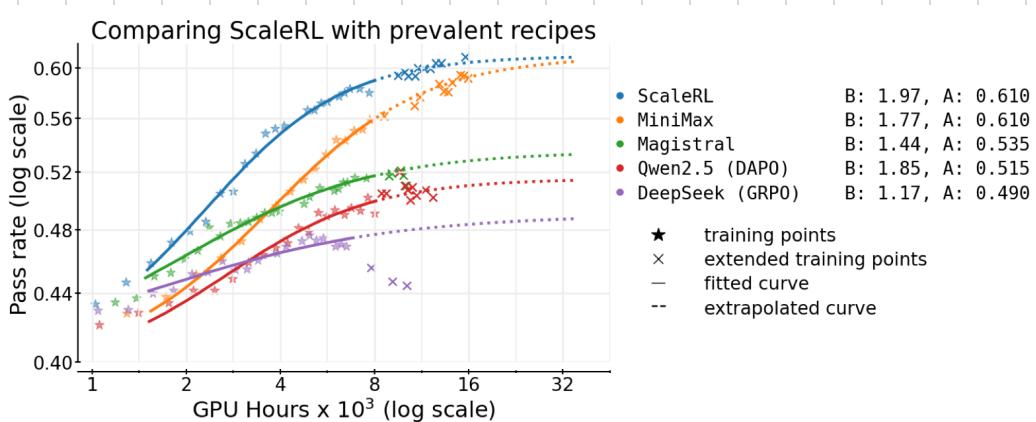


Fig 2 in ScaleRL...

"Bias of tokens perspective"

Unpacking what to do....

$$\mathcal{J}_{\text{ScaleRL}}(\theta) = \mathbb{E}_{\substack{x \sim D, \\ \{y_i\}_{i=1}^G \sim \pi_{\text{gen}}^{\theta_{\text{old}}}(\cdot | x)}} \left[\frac{1}{\sum_{g=1}^G |y_g|} \sum_{i=1}^G \sum_{t=1}^{|y_i|} \text{sg}(\min(\rho_{i,t}, \epsilon)) \hat{A}_i^{\text{norm}} \log \pi_{\text{train}}^{\theta}(y_{i,t}) \right],$$

$$\rho_{i,t} = \frac{\pi_{\text{train}}^{\theta}(y_{i,t})}{\pi_{\text{gen}}^{\theta_{\text{old}}}(y_{i,t})}, \quad \hat{A}_i^{\text{norm}} = \hat{A}_i / \hat{A}_{\text{std}}, \quad 0 < \text{mean}(\{r_j\}_{j=1}^G) < 1, \quad \text{pass_rate}(x) < 0.9,$$

"Input-Sample"
Distribution Mismatch.

Generations Stop Gradient

HW 3 log P_θ

Demean across 6 generations round

Don't bother if we get 90% mld.

VR gives this.

Recall HW3 Trick (Reparameterization Gradient Estimator)

$$\nabla_{\theta} \mathbb{E}_{y \sim P_{\theta}} [f(y)] = \mathbb{E}_{y \sim P_{\theta}} [f(y) \nabla_{\theta} \log P_{\theta}]$$

$\text{sg}(\dots) \leftarrow$ To stop the gradient from inadvertently going into $\rho_{i,t}$.

$$\min(\rho_{i,t}, \epsilon) = \begin{cases} \epsilon & \text{if } \rho_{i,t} \geq \epsilon \rightarrow \text{Caps the overweights of this sample in case it is strongly preferred by } \pi_{\text{true}} \text{ as opposed to } \pi_{\text{gen}}. \\ \rho_{i,t} & \text{if } \rho_{i,t} < \epsilon \end{cases}$$

In general $\pi_{\text{gen}}^{\theta_{\text{old}}}$ closely tracks $\pi_{\text{train}}^{\theta}$ since it gets updated with a lag.

KL Regularization with a reference policy π_{ref}
and reward function r

$$\underset{\pi}{\operatorname{argmax}} \underset{x, y \sim \pi(\cdot | x)}{E[r(x, y)]} - \beta D_{\text{KL}}[\pi || \pi_{\text{ref}}]$$

has solution

$$\frac{1}{Z(x)} \pi_{\text{ref}}(y | x) e^{\frac{1}{\beta} r(x, y)}$$

Moves mass towards higher reward

e.g. GPRO has KL Regularization.

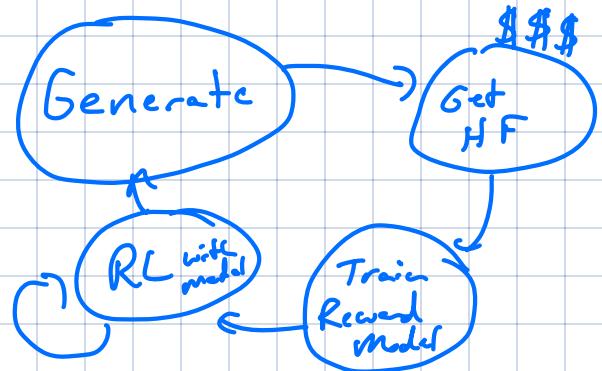
RLHF: RL from human feedback.

Key constraints/considerations:

A) Humans are expensive \Rightarrow Use human responses to train AI - feedback / reward models

B) Humans aren't good at real-valued feedback \Rightarrow Ask for binary comparison

Pretrain \rightarrow SFT on Good Examples \Rightarrow



DPO: Direct Preference Optimization

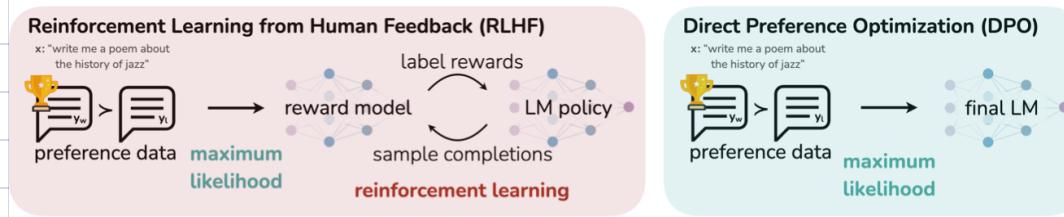


Fig 1 from the paper
"Direct Preference Optimisation
Your LM is secretly a
reward model" by
Rafailov, et al. 2023

Can we work with only the preference outputs?

Idea 0: Just maximise $\log \pi_B(\vec{y}_w) - \log \pi_\theta(\vec{y}_e)$ over dataset

Table 3 from paper

Total Meltdown 😞

Idea 1: Just SFT on winners... Not horrible.

Idea 2: (DPO) Use Idea 1 to get $\pi_{\text{ref}} \leftarrow A$ reference policy

Then

$$\mathcal{L}_{\text{DPO}}(\pi_\theta; \pi_{\text{ref}}) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[\log \sigma \left(\beta \log \frac{\pi_\theta(y_w | x)}{\pi_{\text{ref}}(y_w | x)} - \beta \log \frac{\pi_\theta(y_l | x)}{\pi_{\text{ref}}(y_l | x)} \right) \right].$$

(start at $\pi_0 = \pi_{\text{ref}}$)

prompt

Sigmoid $\sigma(x) = \frac{1}{1 + e^{-x}}$

Just SGD on preference pairs

Question: If no gradients enter π_{ref} and $\log \frac{a}{b} = \log a - \log b$, what's the point of the π_{ref} ?

We're maximizing $\log \sigma \left(\beta \log \frac{\pi_\theta(y_w | x)}{\pi_0(y_w | x)} + \underbrace{\beta \log \frac{\pi_{\text{ref}}(y_l | x)}{\pi_{\text{ref}}(y_w | x)}}_{\text{Bias term for sigmoid nonlinearity}} \right)$

$$\nabla_\theta \mathcal{L}_{\text{DPO}}(\pi_\theta; \pi_{\text{ref}}) =$$

$$-\beta \mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[\underbrace{\sigma(\hat{r}_\theta(x, y_l) - \hat{r}_\theta(x, y_w))}_{\text{higher weight when reward estimate is wrong}} \left[\underbrace{\nabla_\theta \log \pi(y_w | x)}_{\text{increase likelihood of } y_w} - \underbrace{\nabla_\theta \log \pi(y_l | x)}_{\text{decrease likelihood of } y_l} \right] \right],$$

where $\hat{r}_\theta(x, y) = \beta \log \frac{\pi_\theta(y | x)}{\pi_{\text{ref}}(y | x)}$ *i.e. What is good? what π_θ likes to generate!*

Idea 3 : (TR-DPO) Keep updating π_{ref} every τ timesteps to be π_θ
e.g. $\tau = 512$

Diffusion Models for Generation... "Test-time compute"

Autoregressive Generation builds generations one chunk/token at a time.
Learns a conditional distribution. Leverages inherent sequentiality.

What if we don't have inherent sequentiality? e.g. images

Recall from Lecture 18:

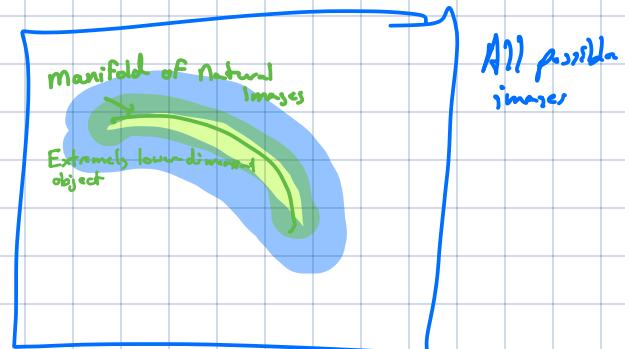
Ideas that don't work

A) Use a classifier



Try: Random Uniform Image
Followed by Gradient Ascent
on the Cat Score

Result: Noise-like image that classifier confidently classifies as a cat.



Bring insights from VAE story:

Need to see noisy things during training

But how noisy?

Two extremes clear:

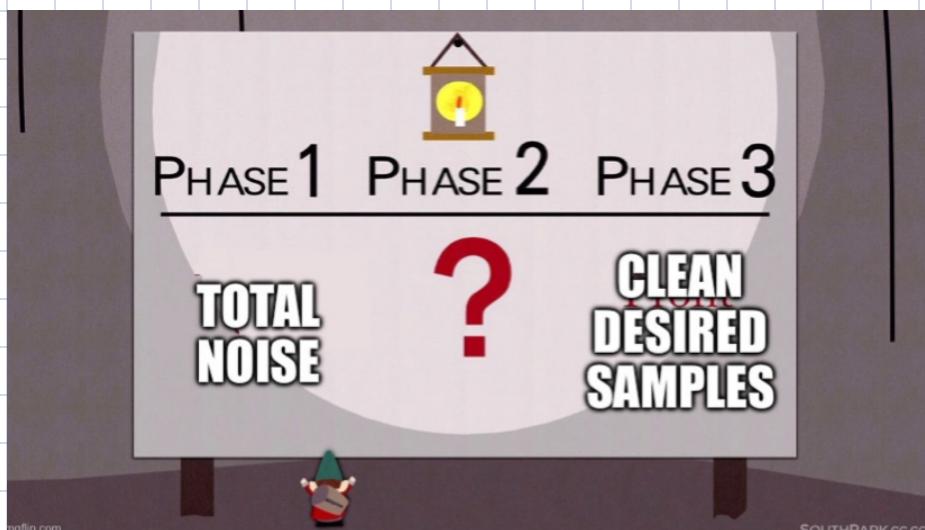


Figure adapted from
South Park to parody/meme

Why not everything in between too?

Can sample from new natural images
by adding noise.

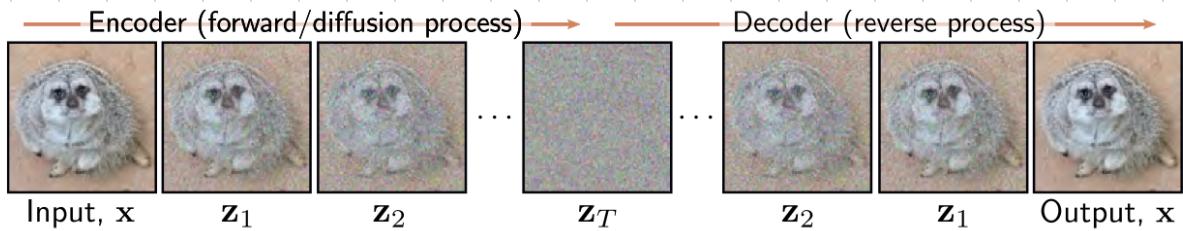
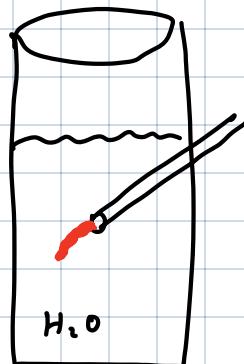
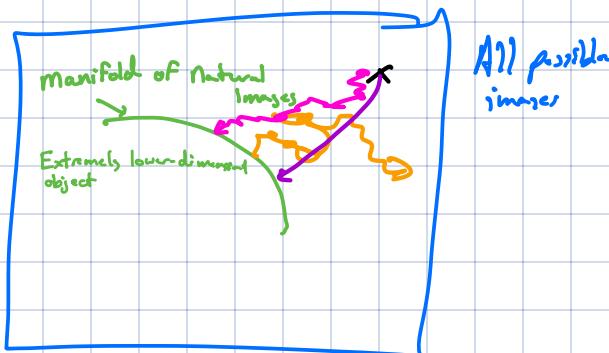


Fig 18.1
in Prince

$$\begin{aligned}\mathbf{z}_1 &= \sqrt{1 - \beta_1} \cdot \mathbf{x} + \sqrt{\beta_1} \cdot \boldsymbol{\epsilon}_1 \\ \mathbf{z}_t &= \sqrt{1 - \beta_t} \cdot \mathbf{z}_{t-1} + \sqrt{\beta_t} \cdot \boldsymbol{\epsilon}_t\end{aligned}\quad \forall t \in 2, \dots, T,$$

Traditional Perspective from
Sohl-Dickstein, et.al.'s
"Deep Unsupervised learning using
Non-equilibrium Thermodynamics"
in ICMC 2015.

For clarity, we'll mostly follow arxiv 2406.08929 from June 2024:
Nakkiran, et.al. "Step-by-step Diffusion: An elementary tutorial"



Idea: Take sample paths of forward diffusion starting at $\{\mathbf{x}_{\text{train}, i}\}$

$$\mathbf{x}_0 \rightarrow \mathbf{x}_1 \rightarrow \mathbf{x}_2 \rightarrow \mathbf{x}_3 \rightarrow \dots \dots \dots \dots \dots \dots \mathbf{x}_{T-1} \rightarrow \mathbf{x}_T$$

Train a network to predict backwards in time: $\mathbf{x}_{t+1} \rightarrow \hat{\mathbf{x}}_t$

Somehow use this learned network to simulate a path backwards from \mathbf{x}_T to \mathbf{x}_0 .

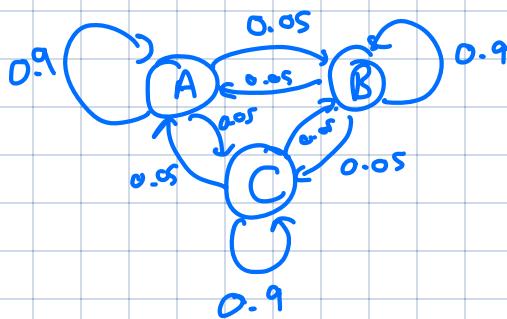
Key Design Question: Do I need randomness on the backward path

Two approaches: stochastic sampling (DDPM)

deterministic sampling (DDIM - next time)

Stochastic Case → Simulate Revers. Process

For Intuition. Consider discrete-time discrete-state



$$\text{Start } [1, 0, 0] \xrightarrow{\text{(A) } p(A) \text{ or } p(C)} \left[\frac{1}{3}, \frac{1}{3}, \frac{1}{3} \right]$$

Is $p(X_{t+1} | X_t)$ like a Gaussian?

(Assume $X_t \rightarrow X_{t+1}$ is conditionally Gaussian. e.g. $X_{t+1} = X_t + N_t \stackrel{\uparrow}{\sim} N(0, \sigma^2 I)$)

As long as density for x_0 sufficiently smooth, if σ^2 small enough,

$$p(x_{t+1} | x_t = \bar{x}) \approx N(\hat{\mu}_{t+1}(\bar{x}), \sigma^2 I)$$

↑
Need to learn this from data.

Let's take a continuous-time perspective. $T \rightarrow 1$. Δt small step.

$$x_0 \rightarrow x_{\Delta t} \rightarrow x_{2\Delta t} \rightarrow x_{3\Delta t} \rightarrow \dots \dots \dots \dots \dots \quad x_{1-\Delta t} \rightarrow x_1$$

Noise added in one step is $N(0, \sigma^2 \Delta t)$

Want to get $\hat{x}_{t-\Delta t}$ given \hat{x}_t going backwards

$$\text{To do this in DDPM, } \hat{x}_{t-\Delta t} = \hat{\mu}_t(\hat{x}_t) + \underbrace{N(0, \sigma^2 \Delta t)}_{\text{Add noise going backward}}$$

How to train $\vec{\mu}_t(\vec{x}_t)$ ← Deep Network

1) Pick \vec{x}_0 from training set

2) Pick time t in $[0, T]$

3) Construct $\vec{x}_{t-\Delta t} = \vec{x}_0 + N(0, \sigma^2(t-\Delta t))$

$$\vec{x}_t = \vec{x}_{t-\Delta t} + N(0, \sigma^2 \Delta t)$$

4) SGD step for $\vec{\mu}_t(\vec{x}_t)$ to predict $\vec{x}_{t-\Delta t}$ under squared-loss