

# Today: Diffusion Models (Last Lecture!)

Readings in Textbook: Ch 18 (Diffusion)

arxiv 2406.08929 from June 2024  
(First sections only)

## Announce:

Final Project Report: **Extended Snd Dec**  
Poster Session Tue Dec 9<sup>th</sup>  
Final Peer Review Thu Dec 11<sup>th</sup>  
Final Final Report Sun Dec 16<sup>th</sup>  
Final Exam: **Wed Dec 17<sup>th</sup> 8-11AM**  
No Notes. Just Pencil.

Also: EECS Colloquium Next Week  
Ece Kamar MSR  
Dec 10<sup>th</sup> 4-5 PM

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

Recall from Lecture 18:

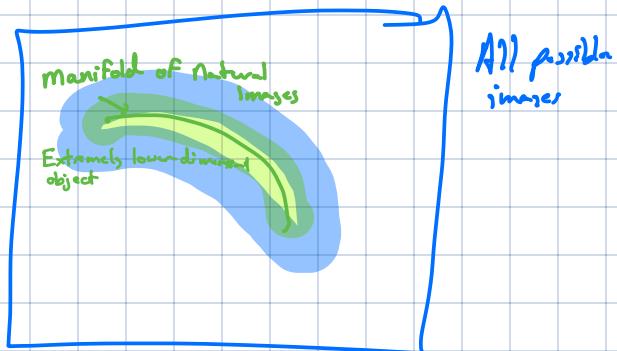
Ideas that don't work

A) Use a classifier



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

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



Bring insights from VAE story:

Need to see noisy things during training



Figure adapted from  
South Park to parody/mem

Why not everything in between too?

Can sample from new natural images  
by adding noise.

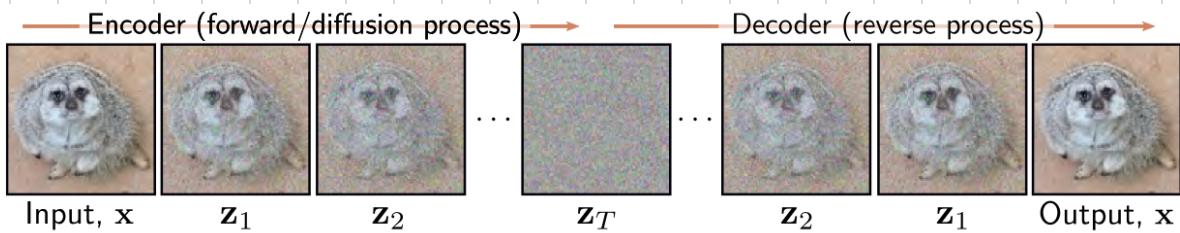
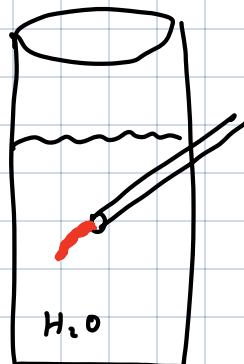
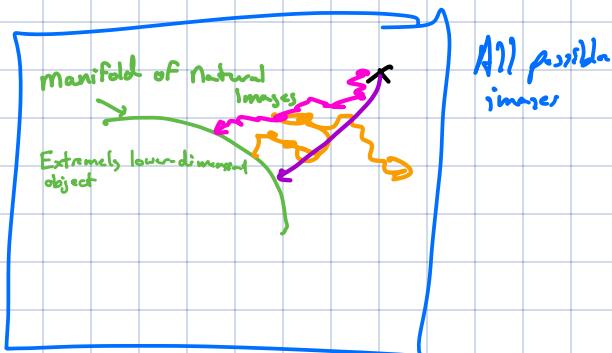


Fig 18.1  
in Prince

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

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  $\{\mathbf{x}_{\text{train},i}\}$

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

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  $\mathbf{x}_0$ .

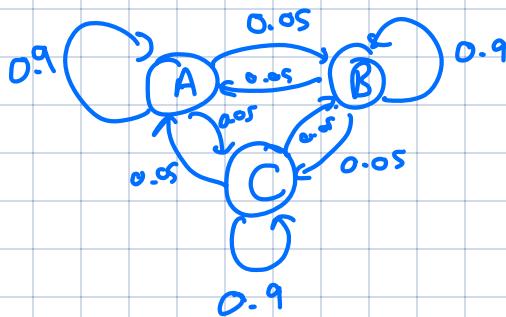
Key Design Question: Do I need randomness on the backward path

Two approaches: stochastic sampling (DDPM)

deterministic sampling (DDIM - next time)

## Stochastic Case $\rightarrow$ Simulate Revers. Process

For Intuition. Consider discrete-time discrete-state



$$\text{Start } [k_1, 0, k_2] \xrightarrow{\text{(A) } p(A) \text{ or } p(C)} \left[ \frac{k_1}{3}, \frac{1}{3}, \frac{k_2}{3} \right]$$

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

Is  $p(x_{t+1} | x_t)$  like a Gaussian? Yes!

As long as density for  $x_0$  sufficiently smooth, if  $\sigma^2$  small enough,

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

↑  
Need to learn this from data.

Before we argue why (HW has you do it for a special case precisely)

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

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

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

Big      Tiny  
→ very small

To do generation...

First draw  $\hat{x}_1 \sim N(0, \sigma^2)$

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

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

Now, refer to claim...

**Claim 1 (Informal).** Let  $p_{t-\Delta t}(x)$  be an arbitrary, sufficiently-smooth density over  $\mathbb{R}^d$ . Consider the joint distribution of  $(x_{t-\Delta t}, x_t)$ , where  $x_{t-\Delta t} \sim p_{t-\Delta t}$  and  $x_t \sim x_{t-\Delta t} + \mathcal{N}(0, \sigma_q^2 \Delta t)$ . Then, for sufficiently small  $\Delta t$ , the following holds. For all conditionings  $z \in \mathbb{R}^d$ , there exists  $\mu_z$  such that:

From  
Paper..

$$p(x_{t-\Delta t} | x_t = z) \approx \mathcal{N}(x_{t-\Delta t}; \underbrace{\mu_z}_{\text{mean}}, \underbrace{\sigma_q^2 \Delta t}_{\text{Variance}}). \quad (16)$$

Idea: Bayes Rule  $\rightarrow$  logs  $\rightarrow$  drop terms that  $\rightarrow 0$  as  $\Delta t \rightarrow 0$

$$p(x_{t-\Delta t} | x_t) = p(x_t | x_{t-\Delta t}) p_{t-\Delta t}(x_{t-\Delta t}) / p_t(x_t)$$

$$\log p(x_{t-\Delta t} | x_t) = \log p(x_t | x_{t-\Delta t}) + \log p_{t-\Delta t}(x_{t-\Delta t}) - \log p_t(x_t)$$

But  $p_{t-\Delta t}(x_{t+\Delta t}) = p_t(x_{t+\Delta t}) + \Delta t \cdot \frac{\partial}{\partial t} p_t(\dots)$

↑ Droppins temporally because it is just normalizing factor

So

$$= \log p(x_t | x_{t-\Delta t}) + \log p_t(x_{t-\Delta t}) + \mathcal{O}(\Delta t) \quad + \text{Normaliza}$$

$$= -\frac{1}{2\sigma_q^2 \Delta t} \|x_{t-\Delta t} - x_t\|_2^2 + \log p_t(x_{t-\Delta t}) \quad \leftarrow \text{dropped}$$

Gaussian Density of forward diffusion..

Now, Taylor expand  $\log p_t(\vec{x}_{t+\Delta t})$  around  $\vec{x}_t$

$$= -\frac{1}{2\sigma_q^2 \Delta t} \|x_{t-\Delta t} - x_t\|_2^2 + \log p_t(x_t) + \langle \nabla_x \log p_t(x_t), (x_{t-\Delta t} - x_t) \rangle + \mathcal{O}(\Delta t) \quad \text{as } \Delta t \rightarrow 0$$

Normalization

$$= -\frac{1}{2\sigma_q^2 \Delta t} \left( \|x_{t-\Delta t} - x_t\|_2^2 - 2\sigma_q^2 \Delta t \langle \nabla_x \log p_t(x_t), (x_{t-\Delta t} - x_t) \rangle \right)$$

$$= -\frac{1}{2\sigma_q^2 \Delta t} \|x_{t-\Delta t} - x_t - \underbrace{\sigma_q^2 \Delta t \nabla_x \log p_t(x_t)}_{\text{"score"}}\|_2^2 + C$$

Complete system

$$= -\frac{1}{2\sigma_q^2 \Delta t} \|x_{t-\Delta t} - \mu\|_2^2$$

mean..

) Variance

How to train  $\vec{\mu}_+(\vec{x}_t) \leftarrow \text{Deep Network}$

- 1) Pick  $\vec{x}_0$  from training set
- 2) Pick time  $t$  in  $[0, 1]$
- 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  $\vec{\mu}_+(\vec{x}_t)$  to predict  $\vec{x}_{t-\Delta t}$  under squared-loss

---

**Pseudocode 1:** DDPM train loss

**Input:** Neural network  $f_\theta$ ; Sample-access to target distribution  $p$ .  
**Data:** Terminal variance  $\sigma_q$ ; step-size  $\Delta t$ .  
**Output:** Stochastic loss  $L_\theta$

```

1  $x_0 \leftarrow \text{Sample}(p)$ 
2  $t \leftarrow \text{Unif}[0, 1]$ 
3  $x_t \leftarrow x_0 + N(0, \sigma_q^2 t)$ 
4  $x_{t+\Delta t} \leftarrow x_t + N(0, \sigma_q^2 \Delta t)$ 
5  $L \leftarrow \|f_\theta(x_{t+\Delta t}, t + \Delta t) - x_t\|_2^2$ 
6 return  $L_\theta$   $\vec{\mu}_+(\vec{x}_{t+\Delta t})$ 
```

---



---

**Pseudocode 2:** DDPM sampling (Code for Algorithm 1)

**Input:** Trained model  $f_\theta$ .  
**Data:** Terminal variance  $\sigma_q$ ; step-size  $\Delta t$ .  
**Output:**  $x_0$

```

1  $x_1 \leftarrow N(0, \sigma_q^2)$ 
2 for  $t = 1, (1 - \Delta t), (1 - 2\Delta t), \dots, \Delta t$  do
3    $\eta \leftarrow N(0, \sigma_q^2 \Delta t)$ 
4    $x_{t-\Delta t} \leftarrow f_\theta(x_t, t) + \eta$ 
5 end
6 return  $x_0$ 
```

---

Figure from  
paper

Note:  $\vec{\mu}_+(\vec{x}_t)$  ideally approximates  $E[\vec{x}_{t-\Delta t} | \vec{x}_t]$

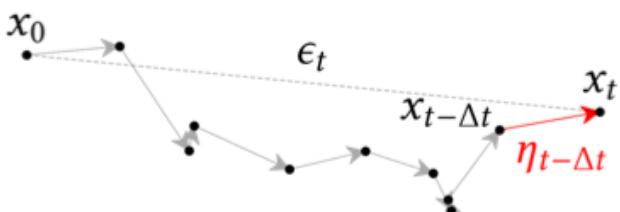
"

$$E[\vec{x}_{t-\Delta t} - \vec{x}_t | \vec{x}_t] + \vec{x}_t$$

negative of  
added noise  
in step from  $t-\Delta t$  to  $t$

So  $\vec{\mu}_+(\vec{x}_t) - \vec{x}_t$  is the expected noise noise.

By iid symmetry of added noise...



$$E[\vec{x}_{t-\Delta t} - \vec{x}_t | \vec{x}_t] = \frac{\Delta t}{t} E[\vec{x}_0 - \vec{x}_t | \vec{x}_t]$$

Can think of the network  $\vec{\mu}_+(\vec{x}_t)$  as learning to predict original image.

i.e.  $E[\vec{x}_{t-\Delta t} | \vec{x}_t] = \frac{\Delta t}{t} E[\vec{x}_0 | \vec{x}_t] + (1 - \frac{\Delta t}{t}) \vec{x}_t$

Story seems great and complete... So why/how can we sample deterministically??



DDPM-style Sampling tries to walk <sup>random</sup> backwards through the joint distribution. That gives us the right marginals and hence a good  $p_0$ .

In Principle, we only need  $p_0$  to be right

So walking backwards through the same marginals but a different joint distribution ok.

Example  $N(0,1) \xleftarrow{X} N(0,2) \xrightarrow{Y}$

Stochastic a)  $Y = X + N(0,1)$  independent  $\Rightarrow$  Bayes Rule  $P(X|Y=y) = N(y, \frac{1}{2})$   
i.e. sample  $X$  by  $y + \underbrace{N(0,1)}_{\text{Add Noise}}$

Deterministic b)  $X \sim N(0,1) \quad Y = \sqrt{2} \cdot X \quad (\text{forward})$   
 $X = \frac{Y}{\sqrt{2}}$

---

Historical Note: How did people discover this idea originally?

Continuous-time diffusion is understood using stochastic differential equations.  
The probability distributions involved follow ODEs (Actually, PDEs)  
The smoothly changing probability distributions can be connected using transport maps  
that turn one density into another. Deterministically.  
Those transport maps can be ridden back from the end.

---

We will follow a different approach: Guess & Verify.

What should we guess???

$$\hat{X}_{t+\Delta t} = \hat{X}_t + \underbrace{? \cdot (\mu_t(x_t) - x_t)}_{\text{Expected Normal Noise}}$$

What should  $?$  be? Try 0 or 1.  $\Rightarrow$  Doesn't work (HW will make you see)

Try  $\lambda$

General Strategy for such things. Try easiest example. That will give us a candidate.

What's the easiest example?

Scalar case.  $P_0 = \delta_0$  i.e.  $N(0, 0)$  i.e. Always 0.

Now we know  $x_t \sim N(0, \sigma^2)$  and  $\mu_t(z) = E[x_{t+\Delta t} | x_t = z]$

$\xrightarrow{\text{N}(0, (\Delta t)^2 \sigma^2)}$   $\xrightarrow{\text{N}(0, \sigma^2)}$   
 index added noise.

Jointly normal!  $\mu_t(z) = \frac{t-\Delta t}{t} z$

$$\begin{aligned} \text{Plug in } x_{t+\Delta t} &= x_t + \lambda \left( \frac{t-\Delta t}{t} x_t - x_t \right) \\ &= x_t \left( 1 - \lambda \frac{\Delta t}{t} \right) \end{aligned}$$

How to solve for  $\lambda$ ?

Need distribution to match.  $\text{Var}(x_{t+\Delta t}) = \sigma^2 (t-\Delta t)$

$$\cancel{\sigma^2 (t-\Delta t)} = \cancel{\sigma^2} + \left( 1 - \lambda \frac{\Delta t}{t} \right)^2$$

$$\sqrt{\frac{t-\Delta t}{t}} = 1 - \lambda \frac{\Delta t}{t} \Rightarrow \lambda \frac{\Delta t}{t} = 1 - \sqrt{\frac{t-\Delta t}{t}} \Rightarrow \lambda = \frac{t}{\Delta t} \left( 1 - \sqrt{1 - \frac{\Delta t}{t}} \right)$$

Alternate form

$$\begin{aligned}\lambda &= \frac{\Delta t}{\sqrt{t}} \left( 1 - \frac{1}{\sqrt{t}} \sqrt{t - \Delta t} \right) = \frac{\sqrt{\Delta t}}{\sqrt{t}} \left( \sqrt{t} - \sqrt{t - \Delta t} \right) \\ &= \sqrt{\Delta t} \cdot \frac{\sqrt{t} - \sqrt{t - \Delta t}}{\sqrt{t}} \\ &= \frac{\sqrt{\Delta t}}{\sqrt{t} + \sqrt{t - \Delta t}}\end{aligned}$$

When  $t$  is "large"  $\lambda \approx \frac{1}{2}$  i.e. Just go  $\frac{1}{2}$  the way

$\text{At } t = \Delta t \quad \lambda = 1$ $= 2\Delta t \quad \lambda \approx 0.9$ $3\Delta t \quad \lambda \approx 0.55$ $\vdots$ $10\Delta t \quad \lambda \approx 0.51$	$\left. \begin{array}{l} \\ \\ \\ \end{array} \right\}$ Last few steps a bit different
-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------	----------------------------------------------------------------------------------------

But does our guess work in general? Yes!!

See paper.

---

#### Pseudocode 2: DDPM sampling (Code for Algorithm 1)

**Input:** Trained model  $f_\theta$ .  $\cancel{\mu}$   
**Data:** Terminal variance  $\sigma_q$ ; step-size  $\Delta t$ .  
**Output:**  $x_0$

```

1  $x_1 \leftarrow \mathcal{N}(0, \sigma_q^2)$ 
2 for  $t = 1, (1 - \Delta t), (1 - 2\Delta t), \dots, \Delta t$  do
3    $\eta \leftarrow \mathcal{N}(0, \sigma_q^2 \Delta t)$ 
4    $x_{t-\Delta t} \leftarrow f_\theta(x_t, t) + \eta$ 
5 end
6 return  $x_0$ 
```

---



---

#### Pseudocode 3: DDIM sampling (Code for Algorithm 2)

**Input:** Trained model  $f_\theta$ .  $\cancel{\mu}$   
**Data:** Terminal variance  $\sigma_q$ ; step-size  $\Delta t$ .  
**Output:**  $x_0$

```

1  $x_1 \leftarrow \mathcal{N}(0, \sigma_q^2)$ 
2 for  $t = 1, (1 - \Delta t), (1 - 2\Delta t), \dots, \Delta t, 0$  do
3    $\lambda \leftarrow \frac{\sqrt{t}}{\sqrt{t - \Delta t} + \sqrt{\Delta t}}$ 
4    $x_{t-\Delta t} \leftarrow x_t + \lambda(f_\theta(x_t, t) - x_t)$ 
5 end
6 return  $x_0$ 
```

---

ODE Perspective: Rewrite

$$\begin{aligned}\hat{x}_{t-\Delta t} &= \hat{x}_t + \lambda(\mu_t(\hat{x}_t) - \hat{x}_t) \\ &= \hat{x}_t + \Delta t \cdot v_t(\hat{x}_t)\end{aligned}$$

$\cancel{\mu}$   
Velocity field

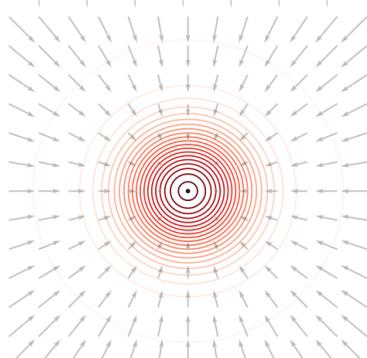


Figure 4: Velocity field  $v_t$  when  $p_0 = \delta_{x_0}$ , overlaid on the Gaussian distribution  $p_t$ .

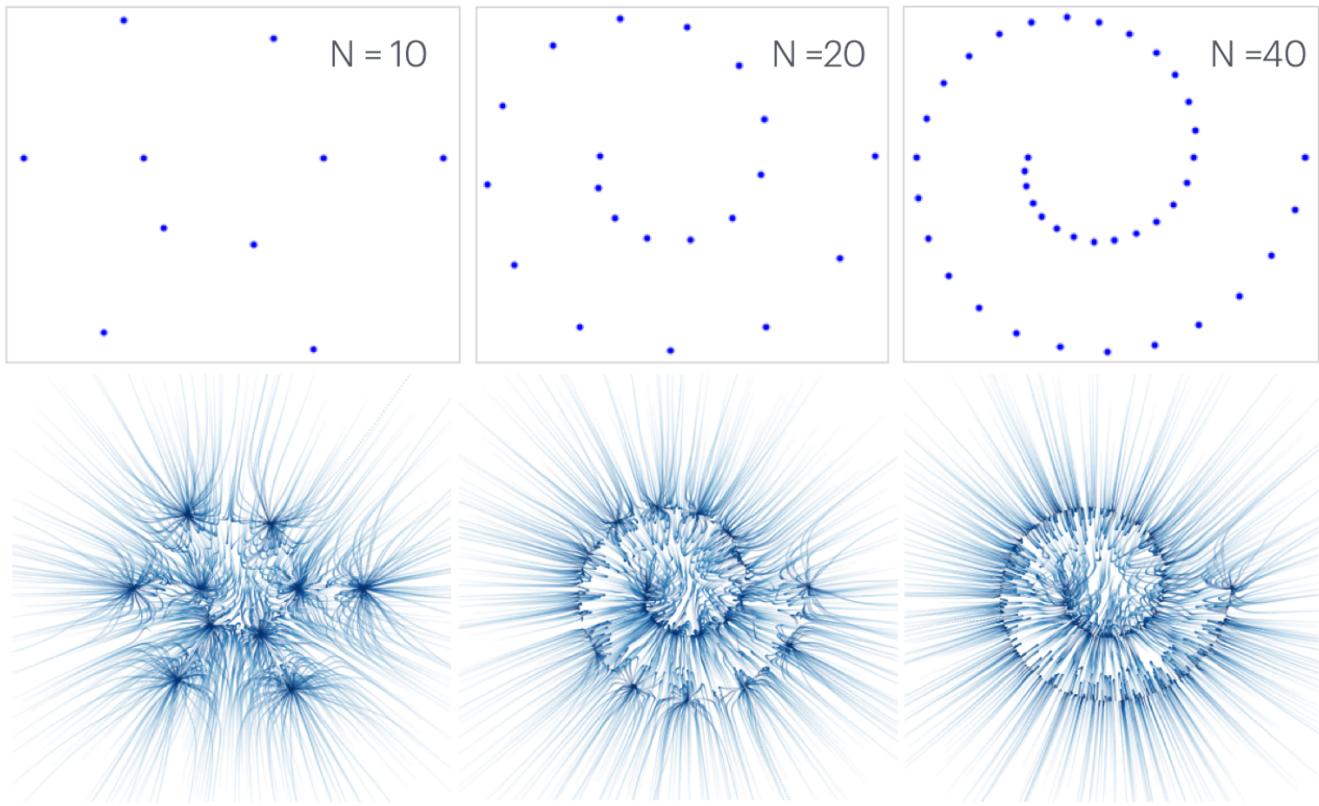


Figure 6: The DDIM trajectories (shaded by timestep  $t$ ) for a spiral dataset. We compare the trajectories with 10, 20, and 40 training samples. Note that as we add more training points (moving left to right) the diffusion algorithm begins to *learn* the underlying spiral and the trajectories look more perpendicular to the underlying manifold. The network used here is a 3 layer ReLU network with 128 neurons per layer.