



# Privacy Issues in DGMs: How to detect & mitigate

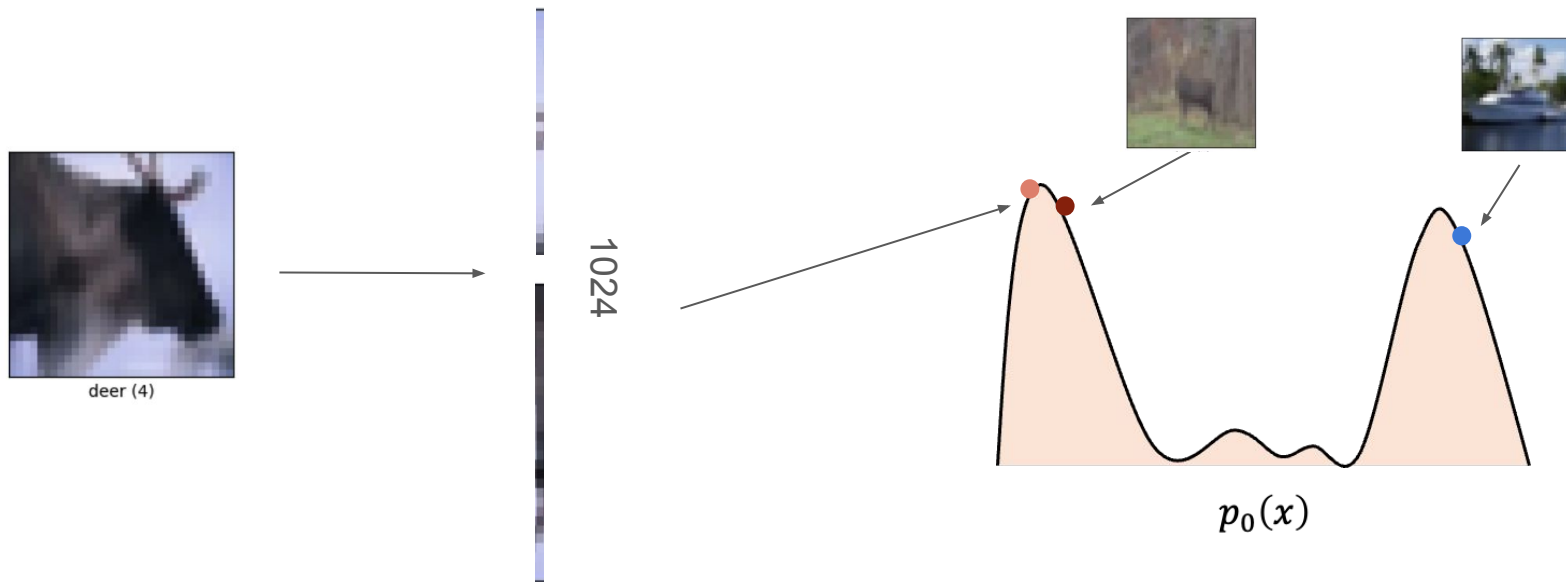
---

Dongjae Jeon

Paper: [lcml link](#)

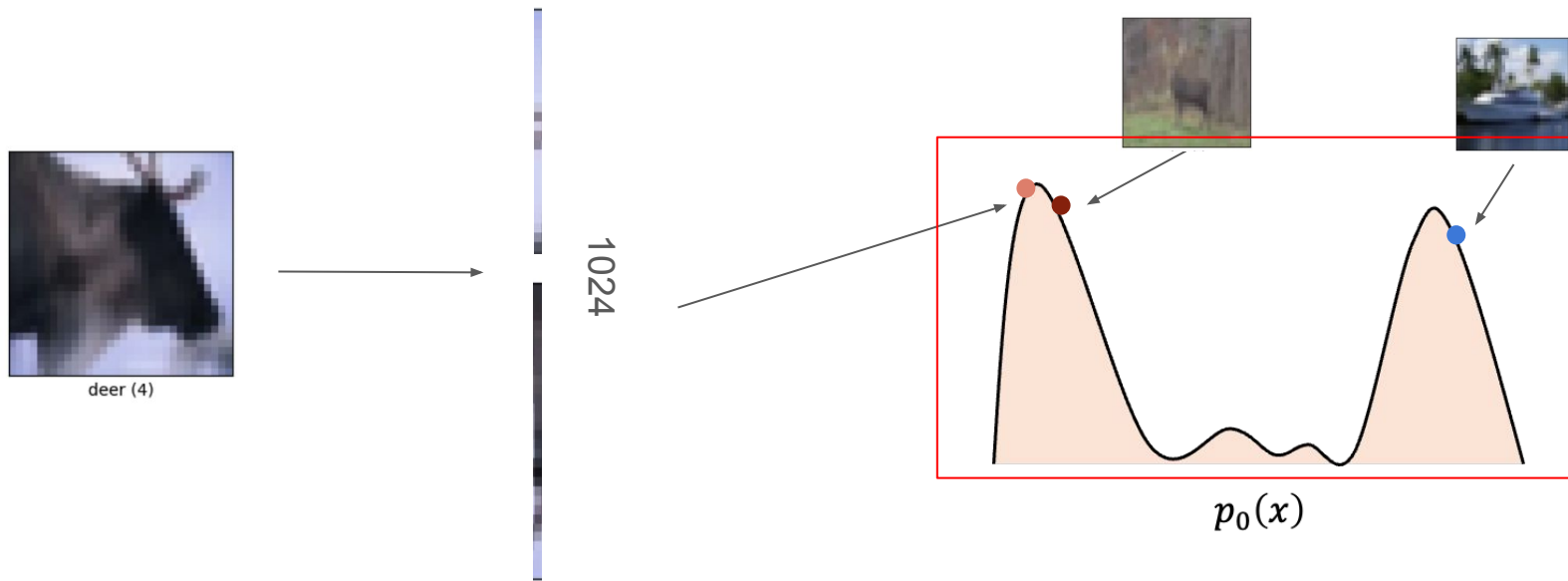
# 1. Detailed Background on Diffusion Model

Image = 1024 sized vector = lives in 1024  $\mathbf{p}(\mathbf{x})$  (we don't know)

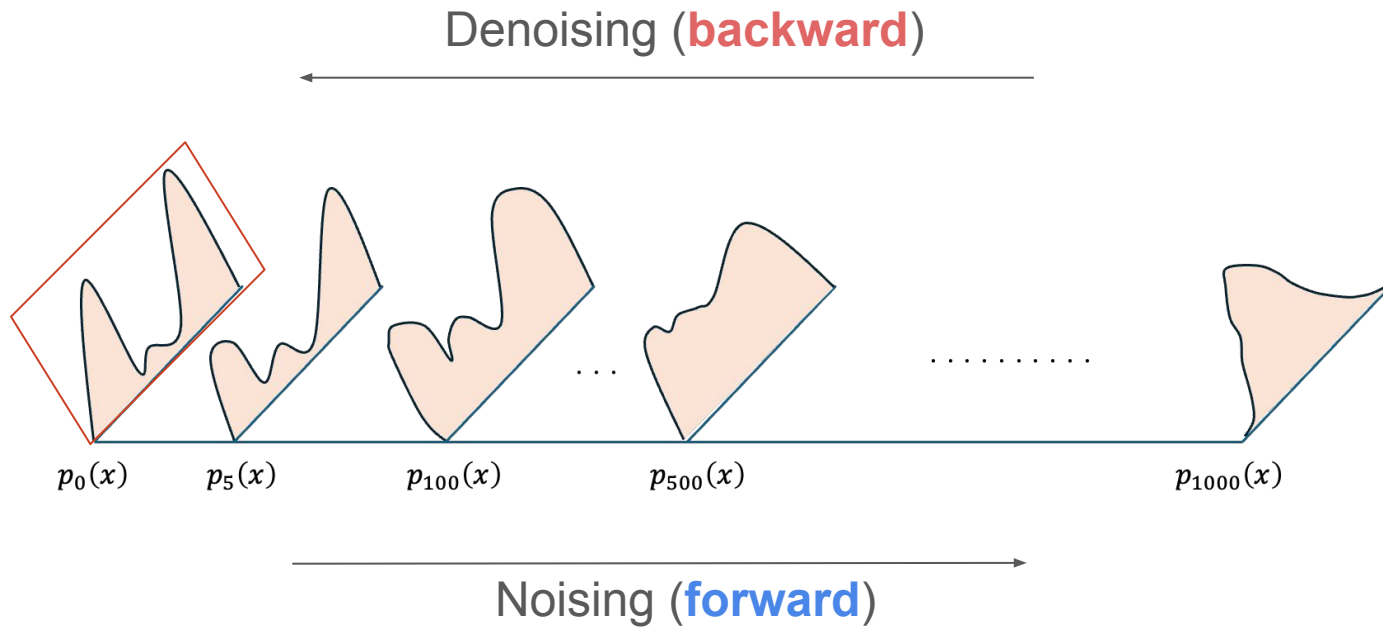


# 1. Detailed Background on Diffusion Model

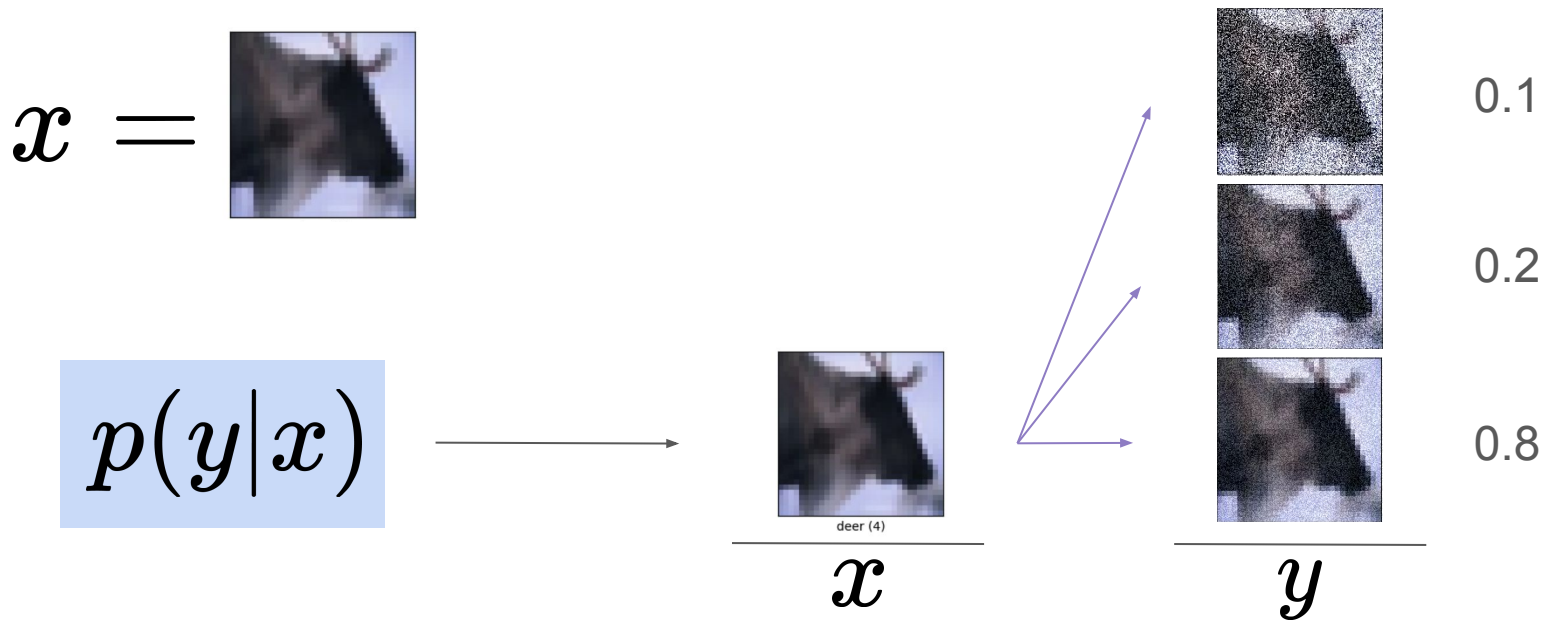
Image = 1024 sized vector = lives in 1024  $\mathbf{p}(\mathbf{x})$  (we don't know)



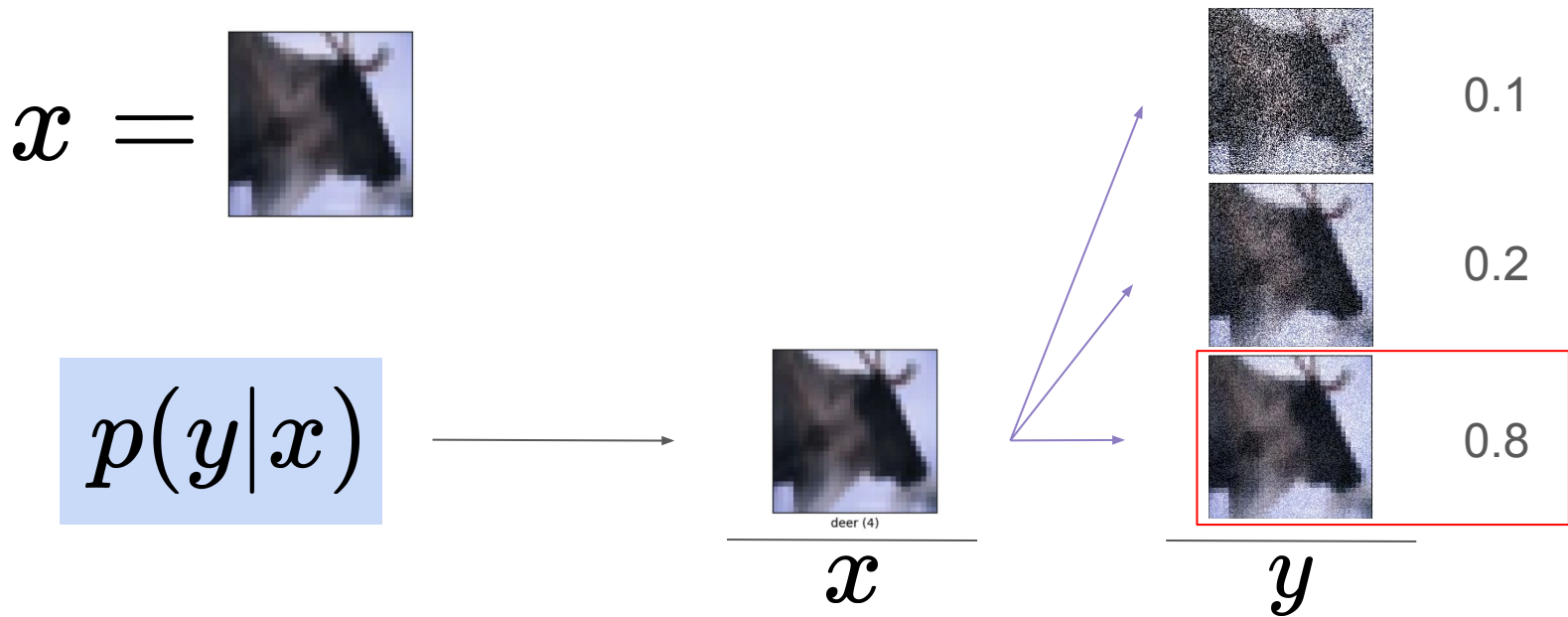
# 1. Detailed Background on Diffusion Model (cont'd)



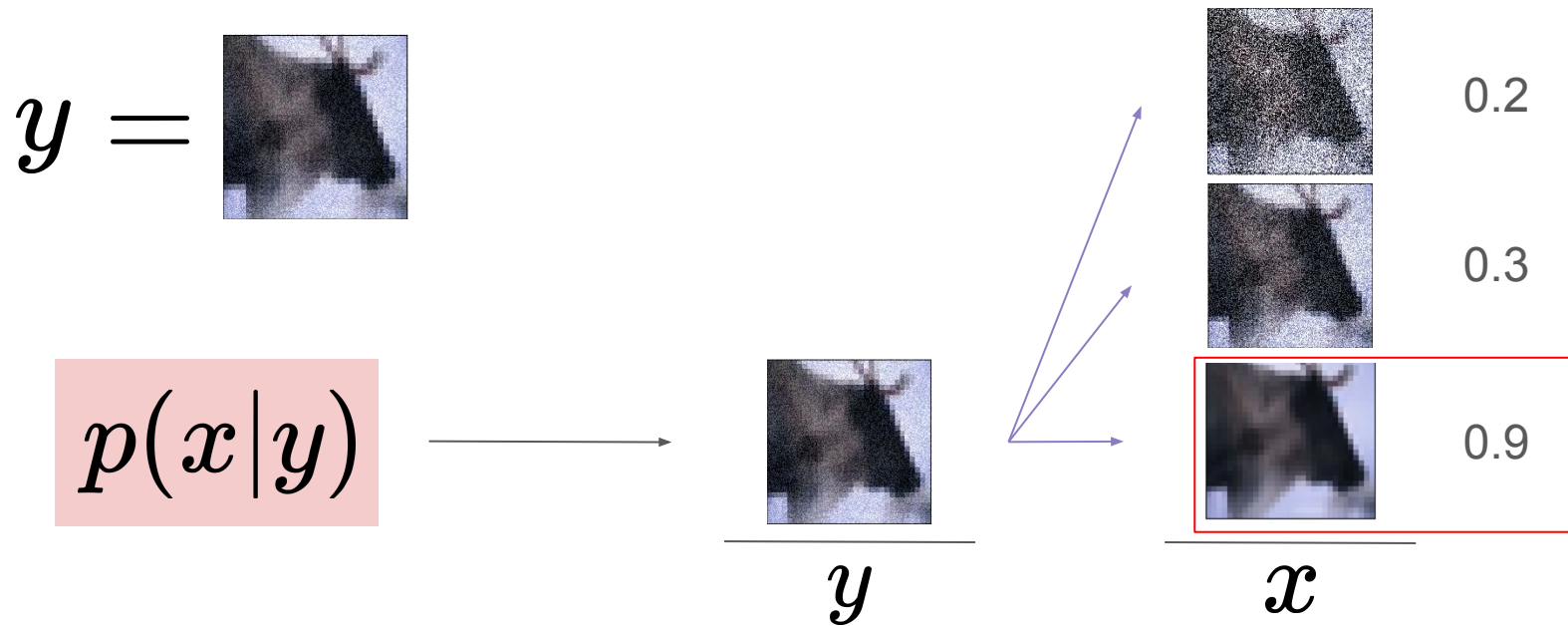
# 1. Detailed Background on Diffusion Model (cont'd)



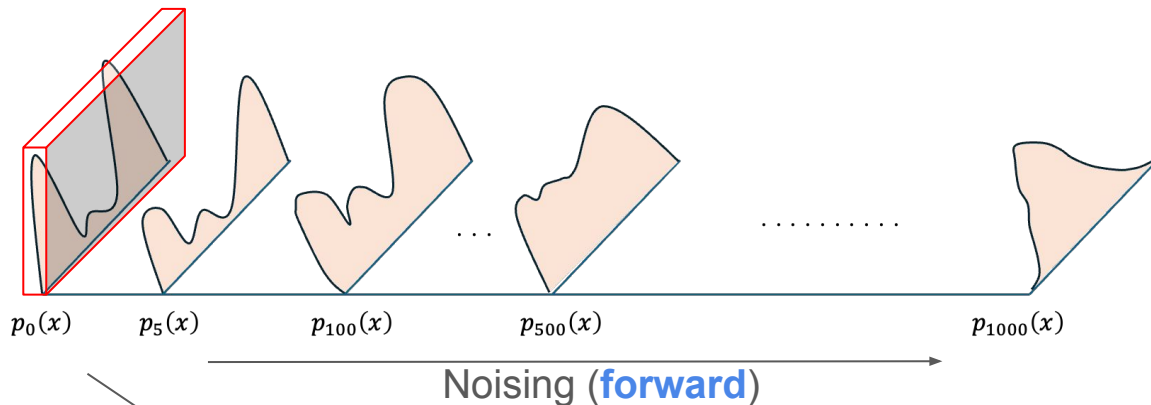
# 1. Detailed Background on Diffusion Model (cont'd)



# 1. Detailed Background on Diffusion Model (cont'd)



# 1. Detailed Background on Diffusion Model (cont'd)



$$y = x + z, \quad z \sim \mathcal{N}(0, I)$$

$$y \sim \mathcal{N}(x, I)$$

$$p_{1|0}(y|x) = \mathcal{N}(y; x, I)$$

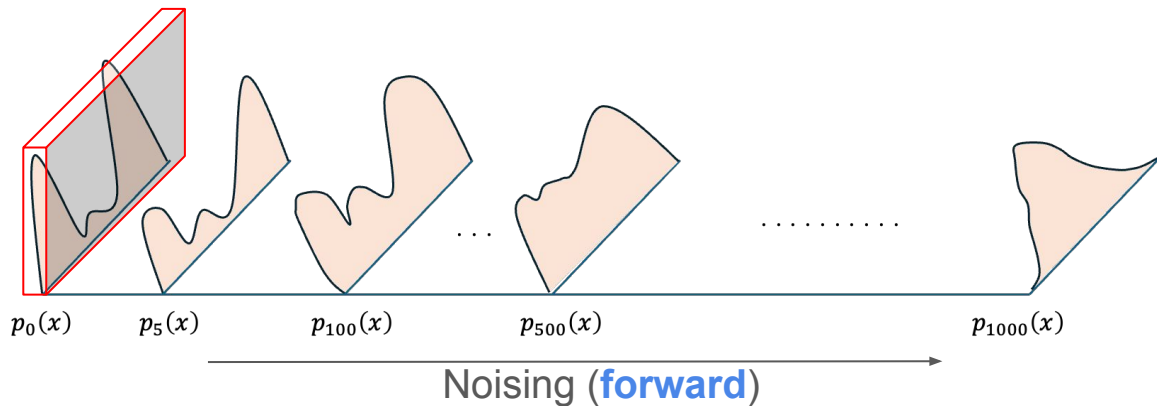
Do this in every time t

$$p_{t|t-1}(x_t | x_{t-1}) = \mathcal{N}(x_t; x_{t-1}, I)$$

Gaussian convolution demo: <https://phiresky.github.io/convolution-demo/>



# 1. Detailed Background on Diffusion Model (cont'd)



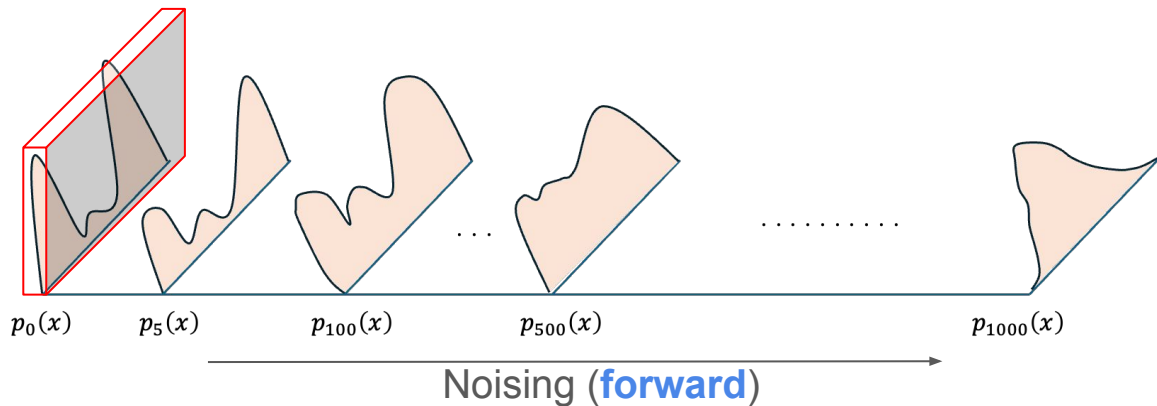
$p_{1000}(x) \neq N(0, I)$   
Then, how to  
start sampling?

$$p_{t|t-1}(x_t \mid x_{t-1}) = \mathcal{N}(x_t; x_{t-1}, I) \quad \text{Variance Exploding}$$

$$p_{t|t-1}(x_t \mid x_{t-1}) = \mathcal{N}\left(x_t; \sqrt{1 - \beta_t} x_{t-1}, \beta_t I\right) \quad \text{Variance Preserving}$$

Suppress to prevent exploding variance.

# 1. Detailed Background on Diffusion Model (cont'd)



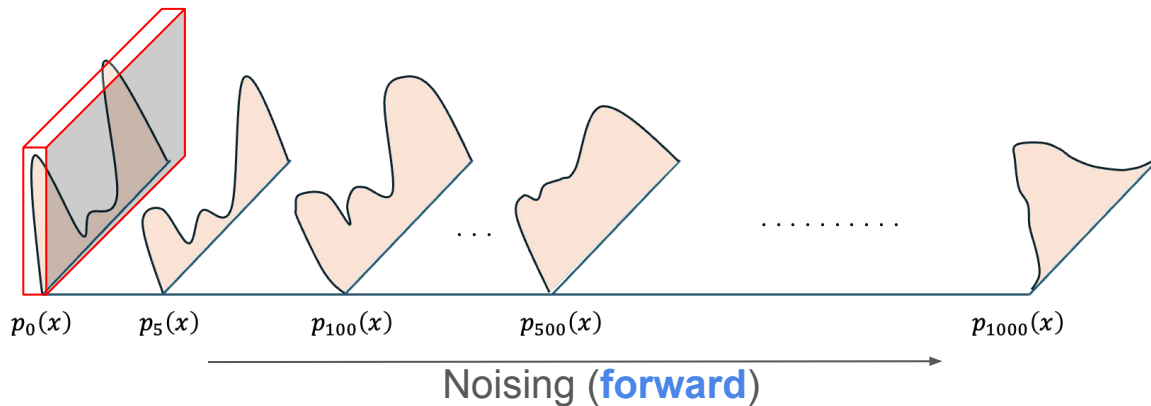
$p_{1000}(x) \neq N(0, I)$   
Then, how to  
start sampling?

$$p_{t|t-1}(x_t | x_{t-1}) = \mathcal{N}(x_t; x_{t-1}, I) \quad \text{Variance Exploding}$$

$$p_{t|t-1}(x_t | x_{t-1}) = \mathcal{N}\left(x_t; \sqrt{1 - \beta_t} x_{t-1}, \beta_t I\right) \quad \text{Variance Preserving}$$

Suppress to prevent exploding variance.

# 1. Detailed Background on Diffusion Model (cont'd)

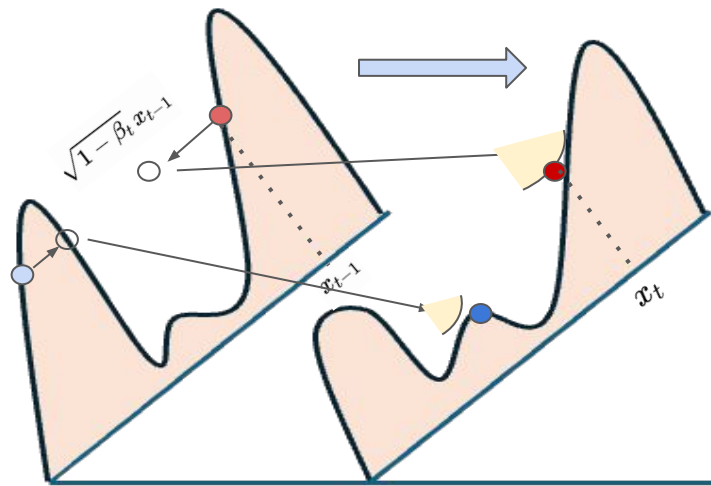


$$p_{t|t-1}(x_t | x_{t-1}) = \mathcal{N}(x_t; x_{t-1}, I) \quad \text{Variance Exploding}$$

$$p_{t|t-1}(x_t | x_{t-1}) = \mathcal{N}\left(x_t; \sqrt{1 - \beta_t} x_{t-1}, \beta_t I\right) \quad \text{Variance Preserving}$$

Our interest

# 1. Detailed Background on Diffusion Model (cont'd)



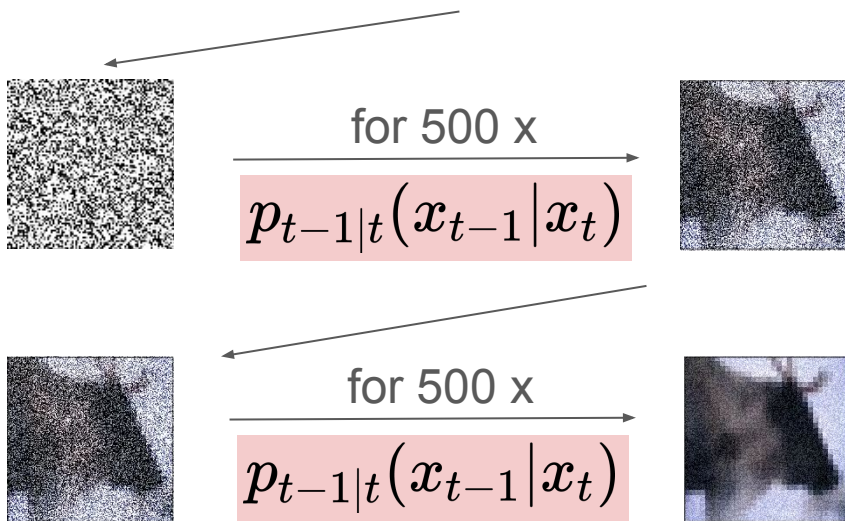
$$p_{t|t-1}(x_t \mid x_{t-1}) = \mathcal{N}(x_t; \sqrt{1 - \beta_t} x_{t-1}, \beta_t I)$$

$$p_{1000}(x) = \mathcal{N}(0, I) \quad \text{No matter what } p_0(x) \text{ is.}$$

# 1. Detailed Background on Diffusion Model (cont'd)

$$p_{1000}(x) = \mathcal{N}(0, I)$$

Sample  $k \sim \mathcal{N}(0, I)$  (What we can sample.)



# 1. Detailed Background on Diffusion Model (cont'd)

$$p_{t|t-1}(x_t | x_{t-1}) = \mathcal{N}(x_t; \sqrt{1 - \beta_t} x_{t-1}, \beta_t I)$$

$$p_{t-1|t}(x_{t-1} | x_t) \approx \mathcal{N}(x_{t-1}; BLANK, \tilde{\beta}_t I)$$

**Tweedie's formula**

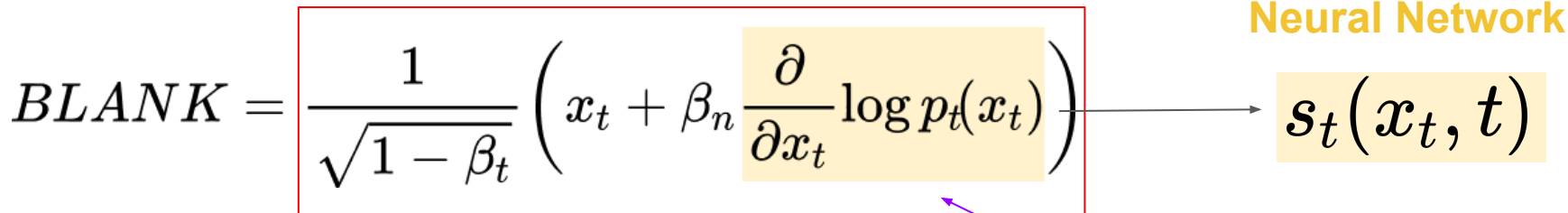
$$BLANK = \frac{1}{\sqrt{1 - \beta_t}} \left( x_t + \beta_t \frac{\partial}{\partial x_t} \log p_t(x_t) \right)$$

# 1. Detailed Background on Diffusion Model (cont'd)

$$p_{t|t-1}(x_t | x_{t-1}) = \mathcal{N}(x_t; \sqrt{1 - \beta_t} x_{t-1}, \beta_t I)$$

$$p_{t-1|t}(x_{t-1} | x_t) \approx \mathcal{N}(x_{t-1}; BLANK, \tilde{\beta}_t I)$$

**Tweedie's formula**

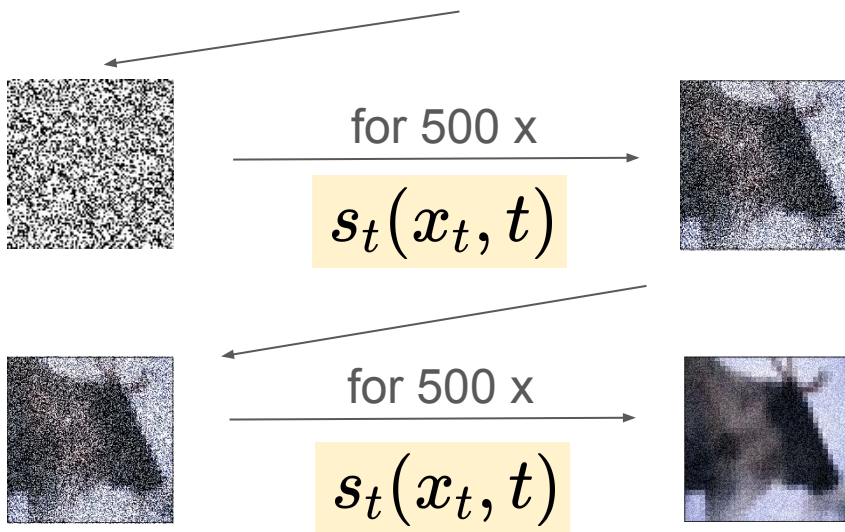
$$BLANK = \frac{1}{\sqrt{1 - \beta_t}} \left( x_t + \beta_n \frac{\partial}{\partial x_t} \log p_t(x_t) \right) \longrightarrow s_t(x_t, t)$$


also called "Score"

# 1. Detailed Background on Diffusion Model (cont'd)

$$p_{1000}(x) = \mathcal{N}(0, I)$$

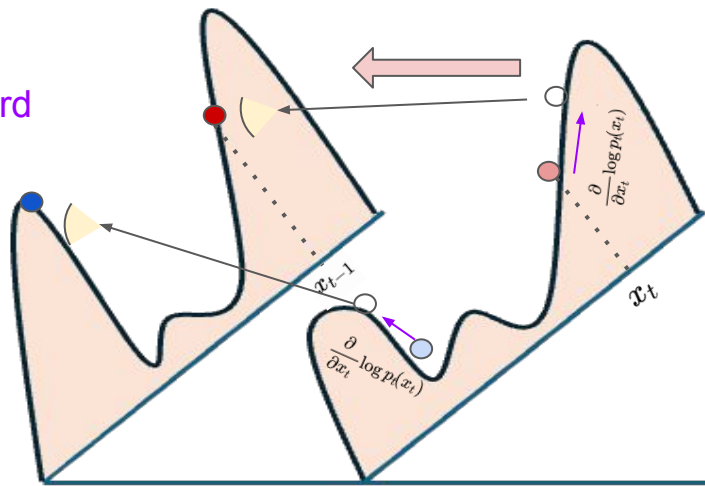
Sample  $k \sim \mathcal{N}(0, I)$  (What we can easily sample)





# 1. Detailed Background on Diffusion Model (cont'd)

Gradually moving samples toward higher likelihood region using “gradient” information.



$$p_{t-1|t}(x_{t-1}|x_t) \approx \mathcal{N}(x_{t-1}; BLANK, \tilde{\beta}_t I)$$

**Tweedie's formula**

$$BLANK = \frac{1}{\sqrt{1 - \beta_t}} \left( x_t + \beta_n \frac{\partial}{\partial x_t} \log p_t(x_t) \right) \xrightarrow{\text{Neural Network}} s_t(x_t, t)$$

# 1. Detailed Background on Diffusion Model (cont'd)

## Takeaway:

- 1) Diffusion Models learn gradient of  $p_t(x)$ :  $\frac{\partial}{\partial x_t} \log p_t(x_t)$
- 2) We can not sample from data distribution directly, but, we can sample from **Gaussian**, and **gradually pushing** it as a “real-like” image.

## 2. Memorization in Diffusion Models

Exact mem.

Training Image



Generated Image

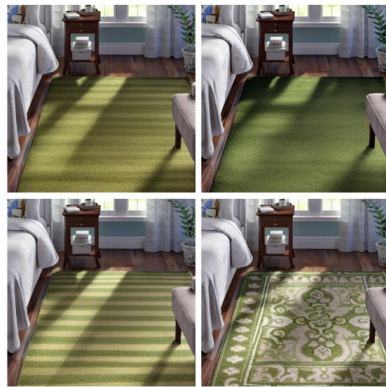


Partial mem.

Training Image



Generated Image



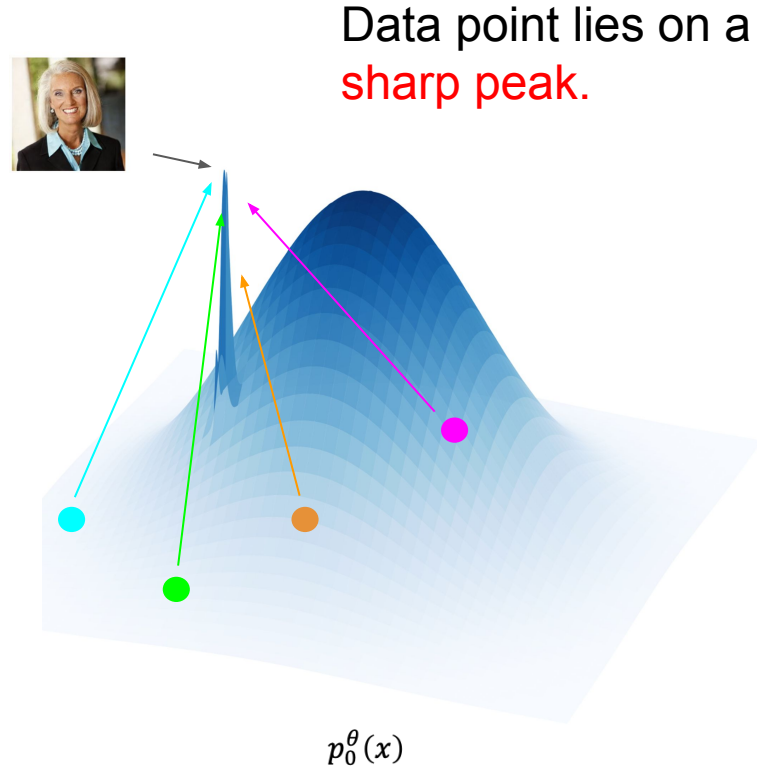
“Living in the Light with Ann Graham Lotz”

“Plattville Green Area Rug by Andover Mills”

Image credit: <https://arxiv.org/pdf/2407.21720>

### 3. What does it mean to be memorized?

Generated Image



## 4. How can we detect it?

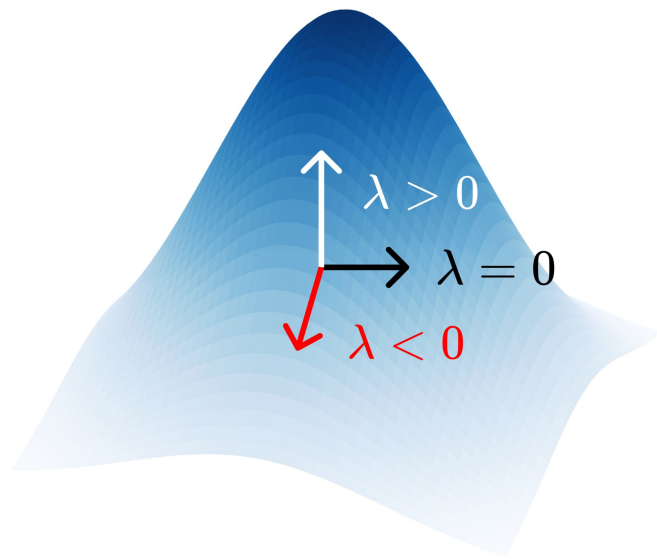
---

$$s_t(x_t, t) \text{ --- } \frac{\partial}{\partial x_t} \log p_t(x_t)$$

$$\frac{\partial}{\partial x_t} s_t(x_t, t) \text{ --- } \frac{\partial^2}{\partial x_t^2} \log p_t(x_t)$$

Hessian Eigenvalues tell Curvature:

- $\lambda \geq 0$ : Concave downward or Flat
- $\lambda < 0$ : Concave upward (Key for finding peaks)



## 4. How can we detect it? (cont'd)

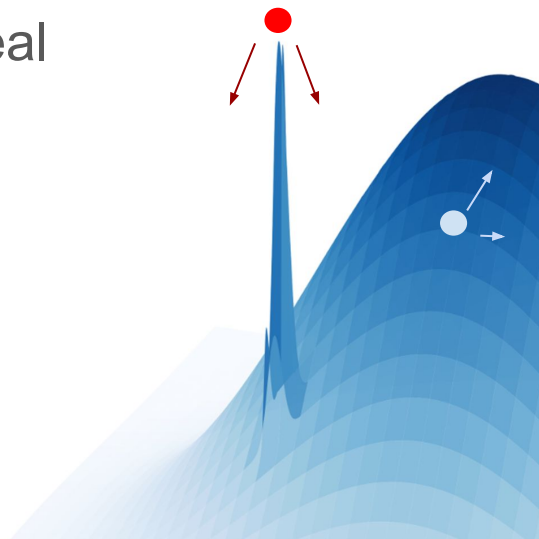
---

Hessian Eigenvalues tell Curvature:

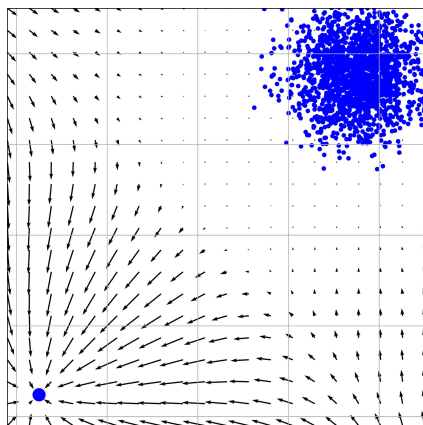
- $\lambda \geq 0$ : Concave downward or Flat
- $\lambda < 0$ : Concave upward (Key for finding peaks)

Memorized sample should reveal  
**large negative eigenvalues,**

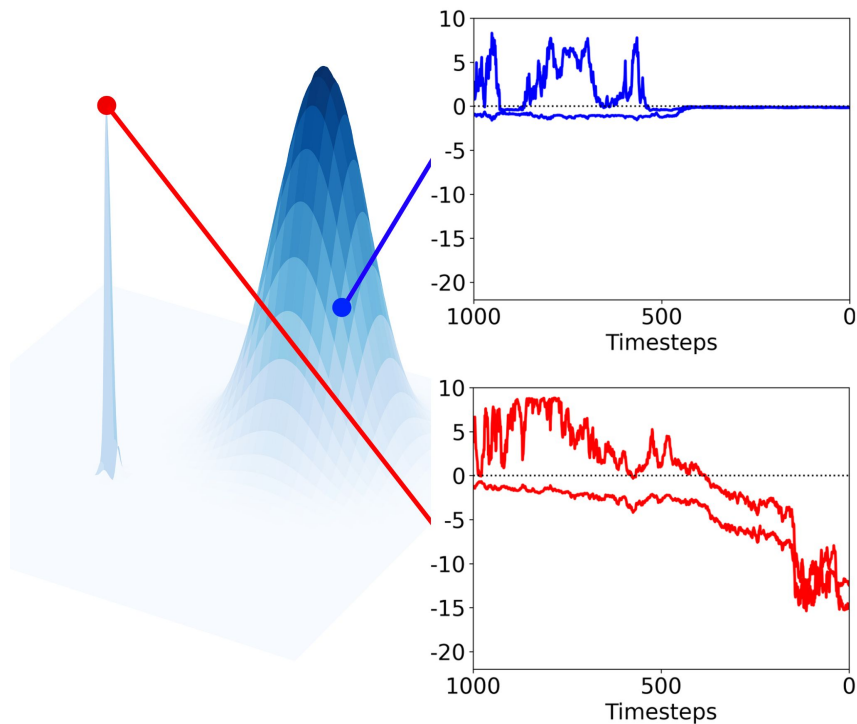
while non-memorized show  
**positive eigenvalues**



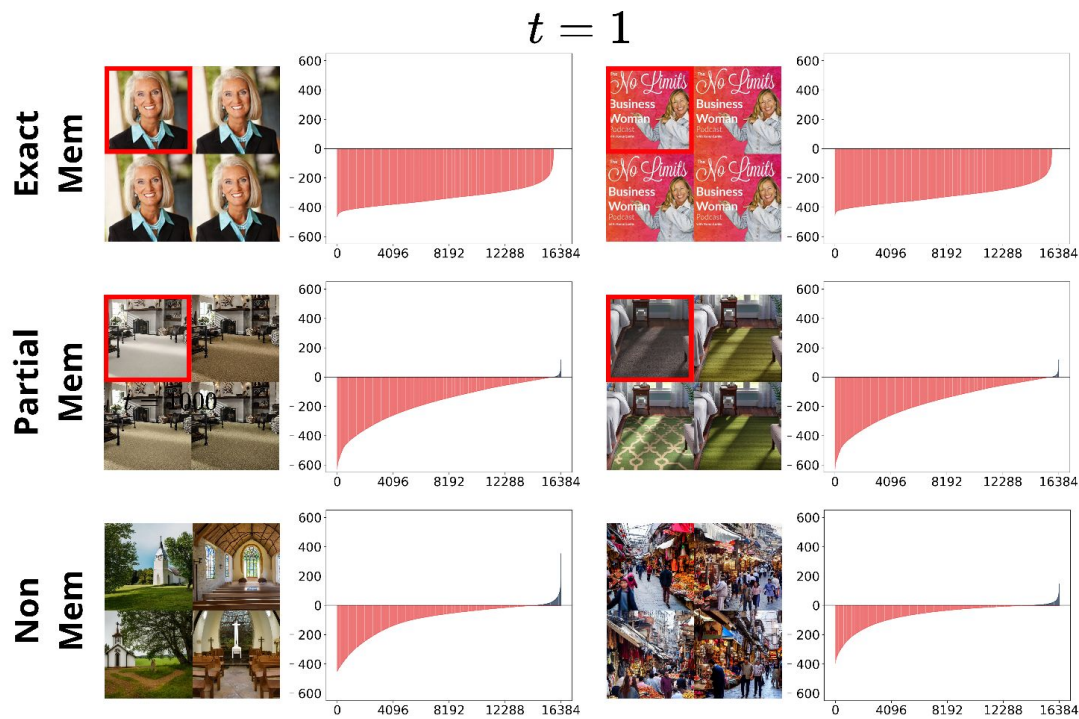
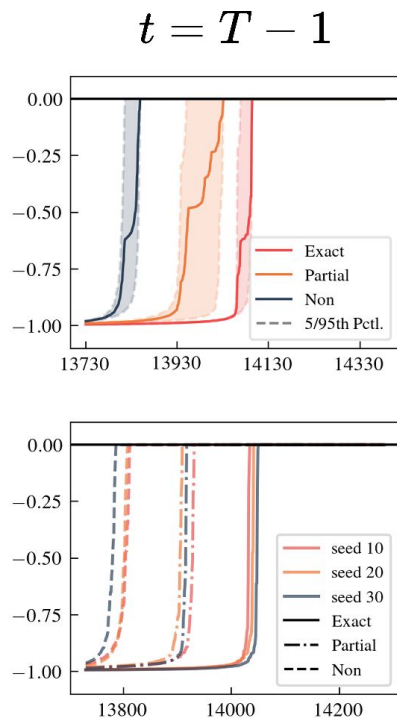
## 4. How can we detect it? (cont'd)



Train data



# 4. How can we detect it? (cont'd)



Eigenvalues in Stable Diffusion



## 4. How can we detect it? (cont'd)

---

But, doing **backpropagation** in Stable Diffusion is nonsense  
We use the **sum of eigenvalues as a proxy!**

Very cheap to compute.

$$\mathbb{E}[\|s_t(x_t, t)\|^2] = -\text{Tr}(H_t(x_t, t)) = -\sum_{i=1}^d \lambda_i,$$

Under gaussian assumption,

$$\mathbb{E}[\|H_t(x_t, t) s_t(x_t, t)\|^2] = -\text{Tr}(H_t(x_t, t)^3) = -\sum_{i=1}^d \lambda_i^3.$$

## 4. How can we detect it? (cont'd)

Method	Steps	$n$	SD v1.4		SD v2.0	
			AUC	TPR@1%FPR	AUC	TPR@1%FPR
Tiled $\ell_2$ (Carlini et al., 2023)	50	4	0.908	0.088	0.792	0.114
		16	0.94	0.232	0.907	0.114
LE (Ren et al., 2024)	1	1	0.846	0.116	0.848	0
		4	0.839	0.13	0.853	0
		16	0.832	0.124	0.851	0
AE (Ren et al., 2024)	50	1	0.606	0	0.809	0
		4	0.628	0	0.82	0
		16	0.598	0	0.817	0
BE (Chen et al., 2024)	50	1	0.986	0.95	0.983	0.908
		4	0.997	0.98	0.99	0.945
		16	0.997	0.982	0.99	0.949
$\ s_\theta^\Delta(\mathbf{x}_t)\ $ (Wen et al., 2024)	1	1	0.976	0.896	0.948	0.739
		4	0.992	0.944	0.98	0.876
		16	0.99	0.928	0.983	0.881
	5	1	0.991	0.932	0.969	0.885
		4	0.997	0.978	0.984	0.917
		16	<b>0.998</b>	<b>0.982</b>	0.987	0.931
	50	1	0.983	0.948	0.982	0.904
		4	0.996	<b>0.982</b>	0.99	<b>0.949</b>
		16	<b>0.998</b>	0.98	<b>0.991</b>	0.945
$\ H_\theta^\Delta(\mathbf{x}_T)s_\theta^\Delta(\mathbf{x}_T)\ ^2$ (Ours)	1	1	0.987	0.908	0.959	0.74
		4	<b>0.998</b>	<b>0.982</b>	<b>0.991</b>	0.895

## 5. How can we mitigate it?

---

Previous approaches,

- [1] Change text prompts
- [2] Put random tokens between prompts
- [3] Weaken text-conditioning during sampling
- .....

**Degrade user utility and image quality!!**

## 5. How can we mitigate it? (cont'd)



## 5. How can we mitigate it? (cont'd)

---

ODE samplers have **1 to 1** relationship between (Xt, Image)  
**Memorization is revealed even at the first timestep!**

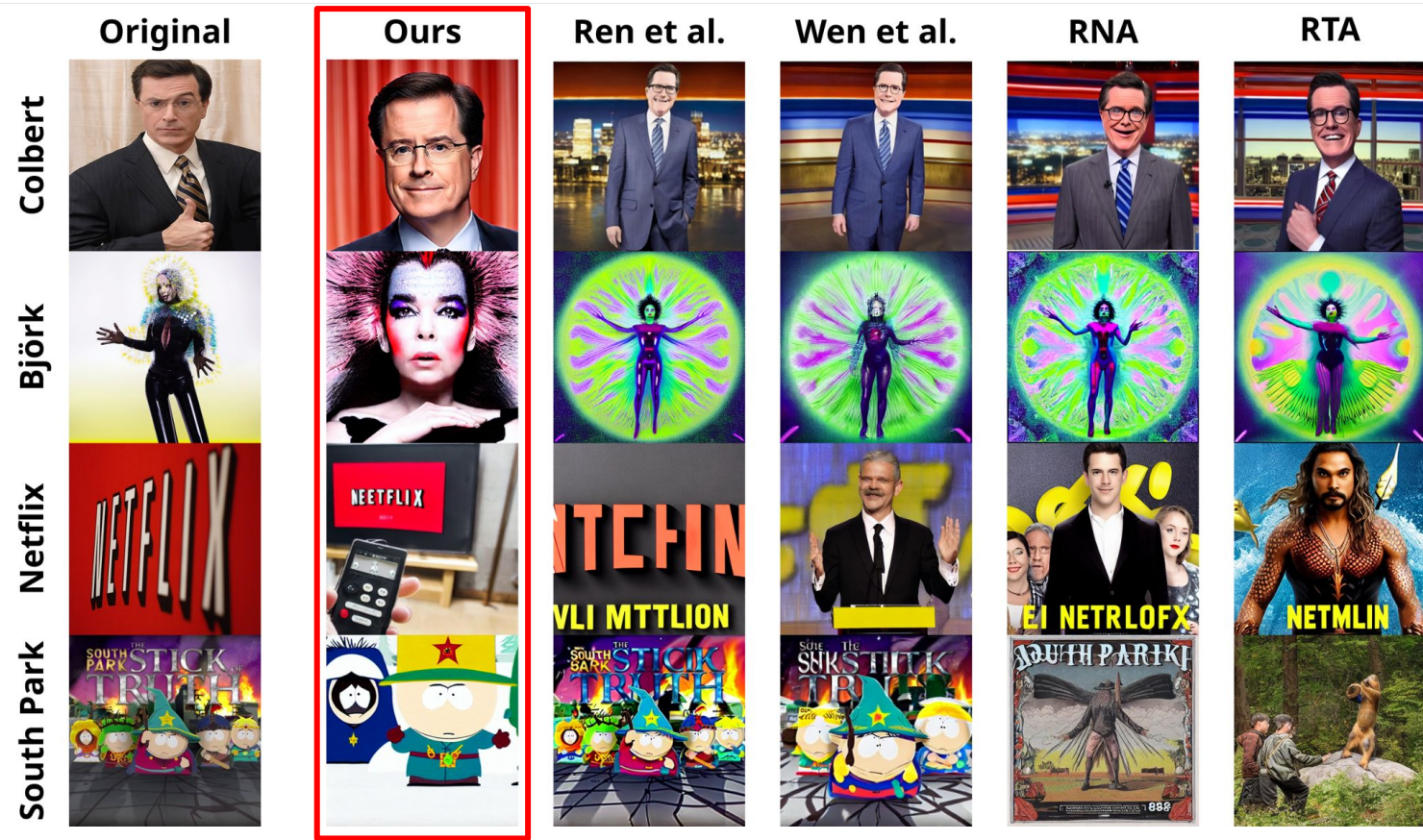
Why don't we just start sampling from  
Gaussian latent on less sharper landscape?  
(a.k.a Seed sampling)

$$\|H_{\Delta\theta}(x_T) s_{\Delta\theta}(x_T)\|^2 - \alpha \log p_G(x_T)$$

Sharpness measure












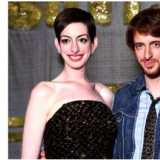



Gaussian regularization

## 5. How can we mitigate it? (cont'd)

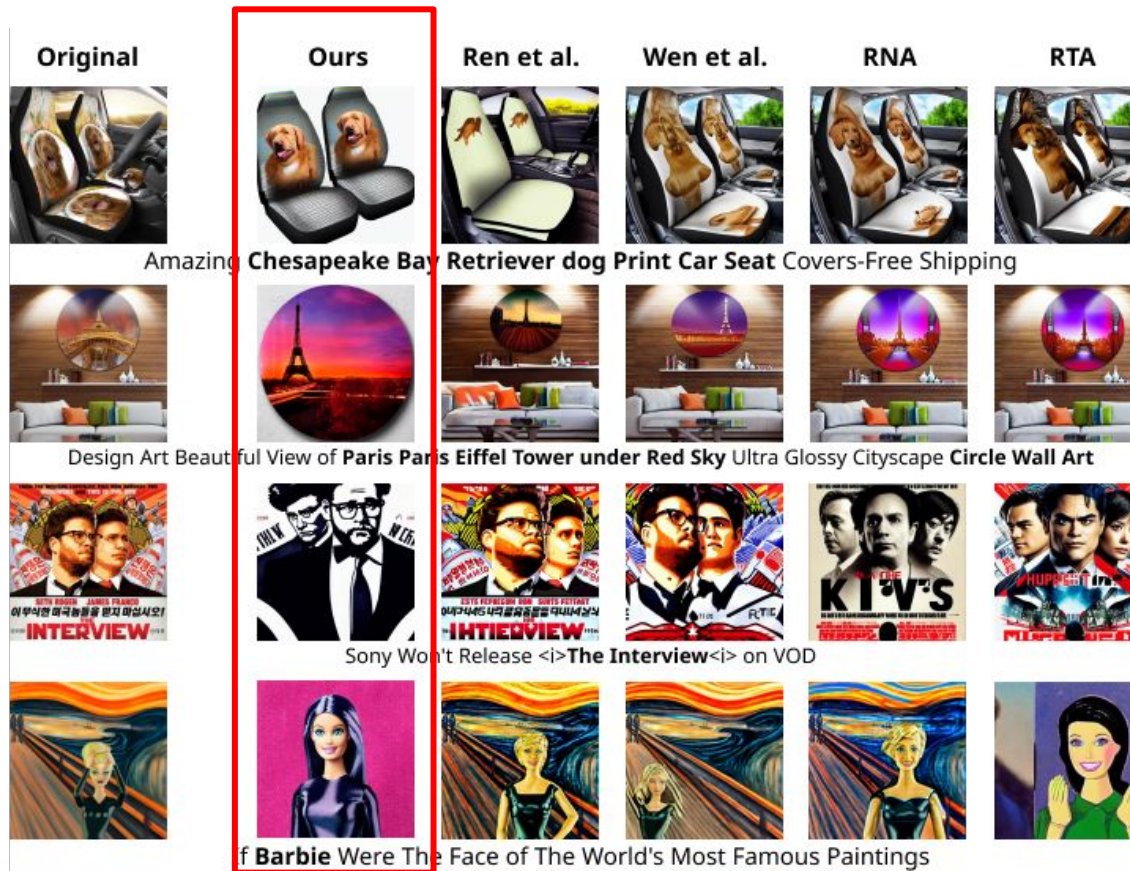




## 5. How can we mitigate it? (cont'd)

Original	Ours	Ren et al.	Wen et al.	RNA	RTA
					
<i>Breaking Bad</i> Fans Get a Chance to Call Saul with Albuquerque Billboard					
					
35 Possible Titles for the <i>Mrs. Doubtfire</i> Sequel.					
					
Baby Shower Turned Meteor Shower: Anne Hathaway Fights Off Aliens in Sci-Fi Comedy <i>The Shower</i>					
					
Listen to Ricky Gervais Perform "Slough" as David Brent					

## 5. How can we mitigate it? (cont'd)





## 6. Advertisement

---

**Visit:**

[https://github.com/Dongjae0324/sharpness\\_memorization\\_diffusion](https://github.com/Dongjae0324/sharpness_memorization_diffusion)

and push “**STAR**”!

