

# Denoising Score Matching

Filippo Garagnani

*Mentor:* Prof. Iacopo Masi

Elicsir Foundation – 3° Orthogonal Weekend

March 22, 2026

Can one say: "Where there is no doubt, there is no knowledge either"?

— *Ludwig Wittgenstein*, *On Certainty*

# Introduction

**Task:** Determining output  $y$  from sample  $x$ .

**Task:** Determining output  $y$  from sample  $x$ .

We want a function

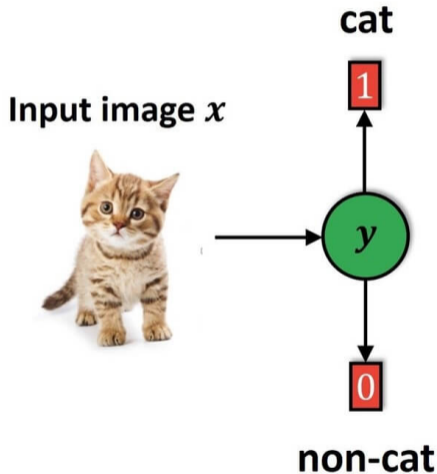
$$f(x) \simeq p(y | x)$$

**Task:** Determining output  $y$  from sample  $x$ .

We want a function

$$f(x) \simeq p(y | x)$$

We do **not** learn anything about how data is structured in space.



**Task:** Model how data is distributed.

**Task:** Model how data is distributed.

We learn

$$p(x)$$

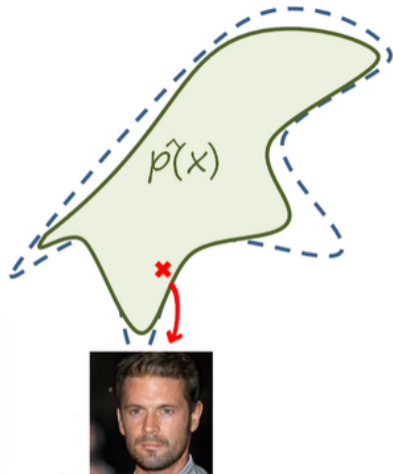
**Task:** Model how data is distributed.

We learn

$$p(x)$$

We could then generate new samples:  $x \sim p(x)$

Evaluate the plausibility of a sample:  $p(\mathbf{x})$



# Learning $p(x)$

We have a specific finite collection of samples:

- **Dataset:**  $\mathcal{D} = \{x_i\}_{i=1}^N$
- **Samples:**  $x_i \in \mathbb{R}^d$

We have a specific finite collection of samples:

- **Dataset:**  $\mathcal{D} = \{x_i\}_{i=1}^N$
- **Samples:**  $x_i \in \mathbb{R}^d$

Assume data samples have been drawn from an unknown underlying distribution:

$$x_i \sim p_{\text{data}}(x)$$

**Problem.**  $p_{\text{data}}$  cannot be computed.

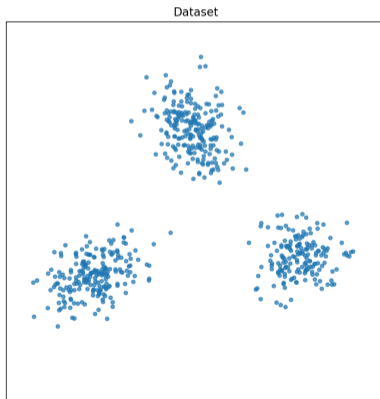
**Problem.**  $p_{\text{data}}$  cannot be computed.

However,

**Problem.**  $p_{\text{data}}$  cannot be computed.

However, we can parametrize a model distribution  $p_{\phi}(x)$  such that:

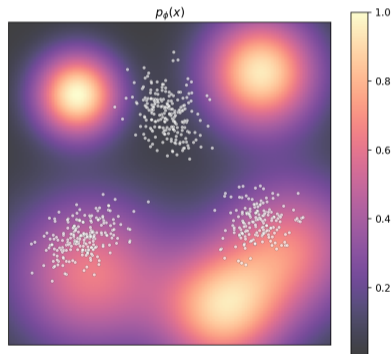
$$p_{\text{data}}(x) \simeq p_{\phi}(x)$$



**Problem.**  $p_{\text{data}}$  cannot be computed.

However, we can parametrize a model distribution  $p_{\phi}(x)$  such that:

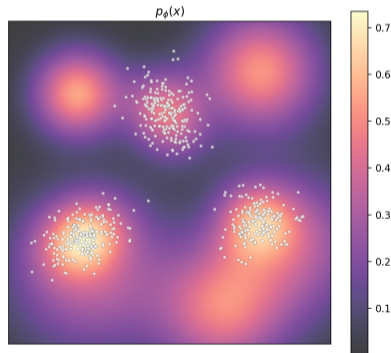
$$p_{\text{data}}(x) \simeq p_{\phi}(x)$$



**Problem.**  $p_{\text{data}}$  cannot be computed.

However, we can parametrize a model distribution  $p_{\phi}(x)$  such that:

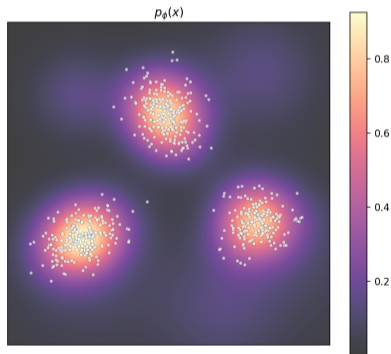
$$p_{\text{data}}(x) \simeq p_{\phi}(x)$$



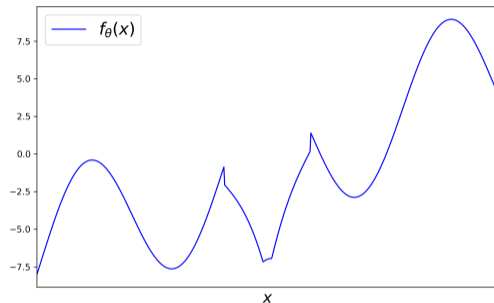
**Problem.**  $p_{\text{data}}$  cannot be computed.

However, we can parametrize a model distribution  $p_{\phi}(x)$  such that:

$$p_{\text{data}}(x) \simeq p_{\phi}(x)$$



Given a generic function  $f_\phi$  (parametrized by  $\phi$ ), for it to be a probability density function it must satisfy:



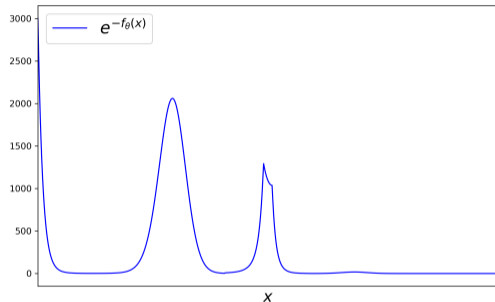
---

<sup>1</sup>[Hinton & Sejnowski, 1985]

<sup>2</sup>[LeCun, Chopra & Hadsell, 2006]

Given a generic function  $f_\phi$  (parametrized by  $\phi$ ), for it to be a probability density function it must satisfy:

**1 Non-Negativity.**  $g_\phi(x) \equiv e^{-f_\phi(x)}$ .



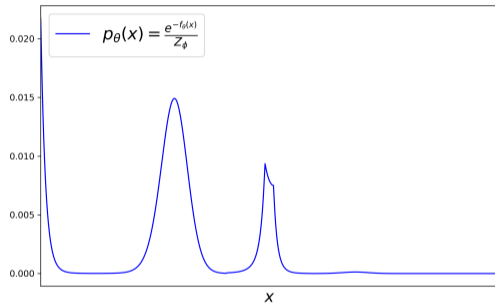
---

<sup>1</sup>[Hinton & Sejnowski, 1985]

<sup>2</sup>[LeCun, Chopra & Hadsell, 2006]

Given a generic function  $f_\phi$  (parametrized by  $\phi$ ), for it to be a probability density function it must satisfy:

- 1 Non-Negativity.**  $g_\phi(x) \equiv e^{-f_\phi(x)}$ .
- 2 Normalization to 1.**  $p_\phi(x) \equiv \frac{1}{Z_\phi} g_\phi(x)$   
( $Z_\phi \equiv \int_{\mathbf{x}} g_\phi(\mathbf{x})$ ).



---

<sup>1</sup>[Hinton & Sejnowski, 1985]

<sup>2</sup>[LeCun, Chopra & Hadsell, 2006]

**Problem.** Usually,  $Z_\phi$  is **intractable**.

**Problem.** Usually,  $Z_\phi$  is **intractable**.  
However,

**Problem.** Usually,  $Z_\phi$  is **intractable**.  
However, we can bypass using it by:

- 1 Transforming it into an **additive term**:  $\log p_\phi(x) = \log g_\phi(x) - \log Z_\phi$ .

**Problem.** Usually,  $Z_\phi$  is **intractable**.

However, we can bypass using it by:

1 Transforming it into an **additive term**:  $\log p_\phi(x) = \log g_\phi(x) - \log Z_\phi$ .

2 Derivating w.r.t. to  $x$ :

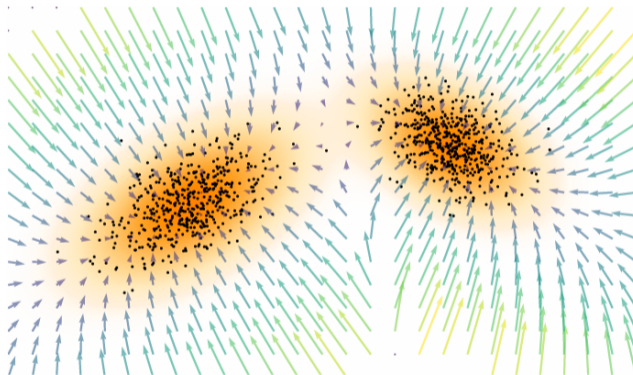
$$\begin{aligned}\nabla_x \log p_\phi(x) &= \nabla_x \log g_\phi(x) - \nabla_x \log Z_\phi \xrightarrow{0} \\ &= -\nabla_x f_\phi(x)\end{aligned}$$

# Score Matching

## Definition

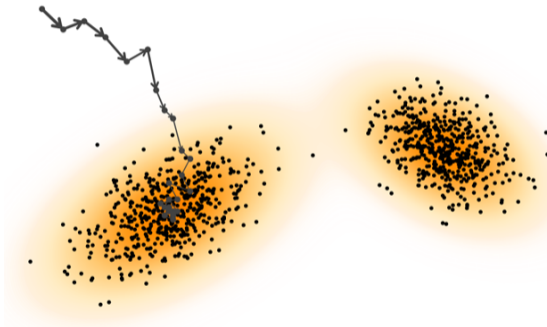
Given a probability density function  $p(x)$ , its **score function**  $s(x)$  is:

$$s(x) \equiv \nabla_x \log p(x)$$

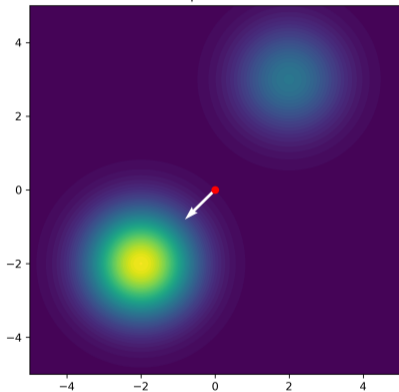


## Langevin Dynamics.

- 1: Initialize  $x_0, \eta$
- 2: **for**  $t \leftarrow 1$  to  $T$  **do**
- 3:      $\mathbf{z}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 4:      $x_t \leftarrow x_{t-1} + \eta \mathbf{s}(x_{t-1}) + \sqrt{2\eta} \mathbf{z}_t$
- 5: **end for**
- 6: **return**  $x_T$

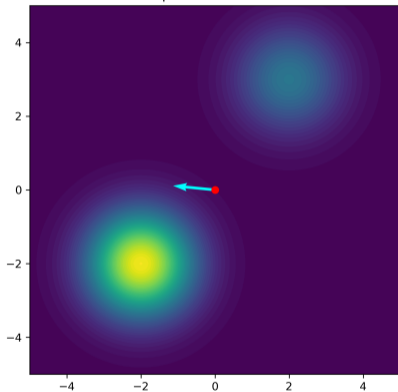


Step 1 - Score



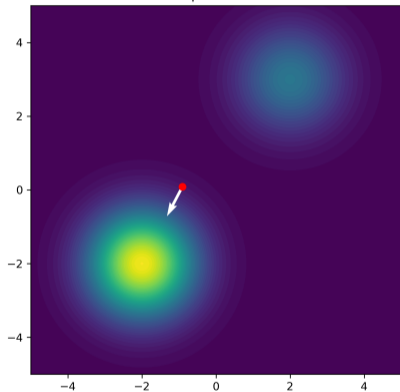
$$\nabla_x \log p(x)$$

Step 1 - Score + Noise



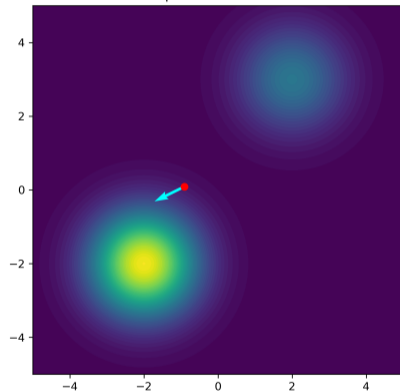
$$\nabla_x \log p(x) + \mathbf{z}_t$$

Step 2 - Score

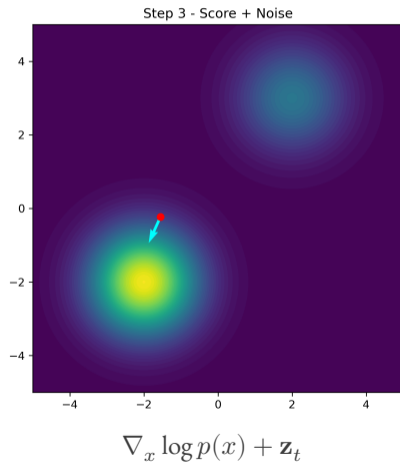
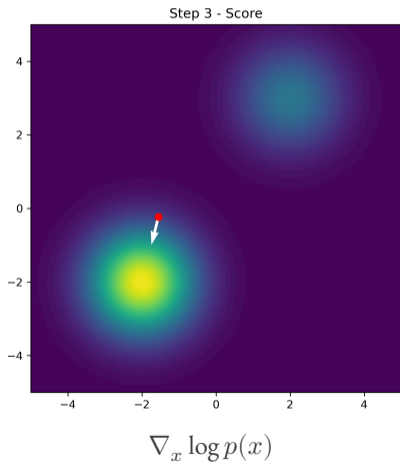


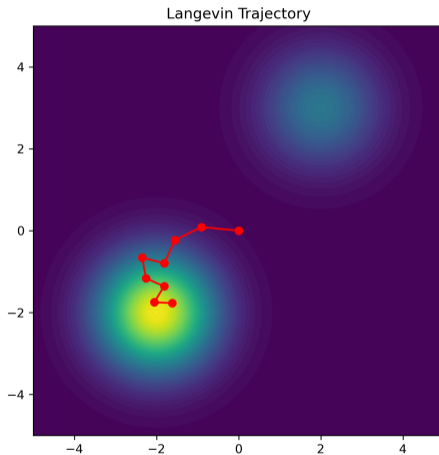
$$\nabla_x \log p(x)$$

Step 2 - Score + Noise



$$\nabla_x \log p(x) + \mathbf{z}_t$$





To approximate  $s_{\text{data}}(x)$ , we can find:

$$\phi^* = \min_{\phi} \mathbb{E}_{x \sim p_{\text{data}}} [\|s_{\phi}(x) - s_{\text{data}}(x)\|^2]$$

To approximate  $s_{\text{data}}(x)$ , we can find:

$$\phi^* = \min_{\phi} \mathbb{E}_{x \sim p_{\text{data}}} [\|s_{\phi}(x) - s_{\text{data}}(x)\|^2]$$

And then employ  $s_{\phi^*}(x)$  to generate data.

**Problem.** To compute  $s_{\text{data}}(x)$ , we would need  $p_{\text{data}}(x)$ .

---

<sup>3</sup>[Hyvärinen, 2005]

**Problem.** To compute  $s_{\text{data}}(x)$ , we would need  $p_{\text{data}}(x)$ .  
However,

---

<sup>3</sup>[Hyvärinen, 2005]

**Problem.** To compute  $s_{\text{data}}(x)$ , we would need  $p_{\text{data}}(x)$ .

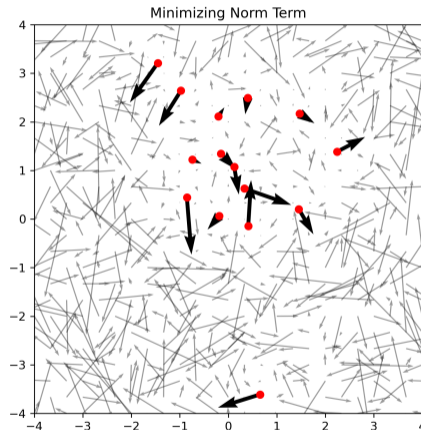
However, it can be shown that  $\phi^*$  is:

$$\begin{aligned}\phi^* &= \min_{\phi} \mathbb{E}_{x \sim p_{\text{data}}} [\|s_{\phi}(x) - s_{\text{data}}(x)\|^2] \\ &= \min_{\phi} \mathbb{E}_{x \sim p_{\text{data}}} [\text{Tr}(\nabla_x s_{\phi}(x)) + \|s_{\phi}(x)\|^2]\end{aligned}$$

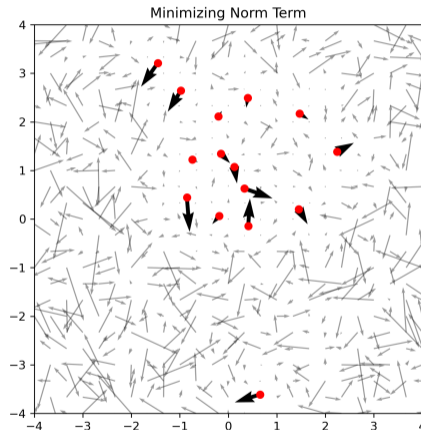
---

<sup>3</sup>[Hyvärinen, 2005]

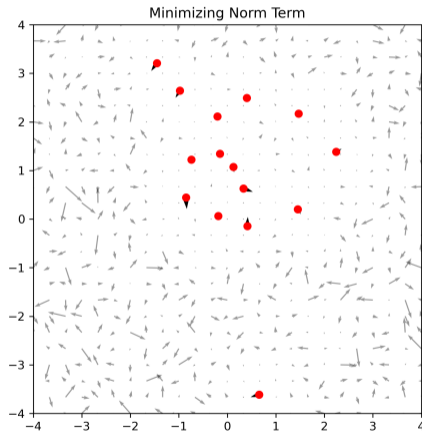
$$\min_{\phi} \mathbb{E}_{x \sim p_{\text{data}}} [\text{Tr}(\nabla_x s_{\phi}(x)) + \|s_{\phi}(x)\|^2]$$



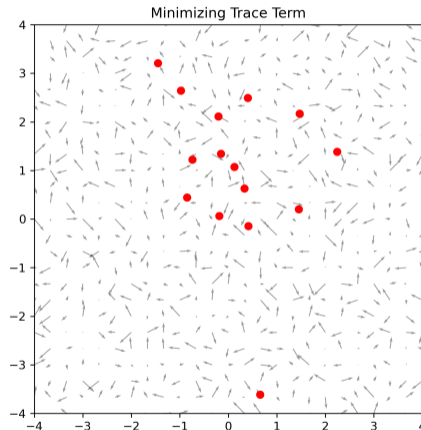
$$\min_{\phi} \mathbb{E}_{x \sim p_{\text{data}}} [\text{Tr}(\nabla_x s_{\phi}(x)) + \|s_{\phi}(x)\|^2]$$



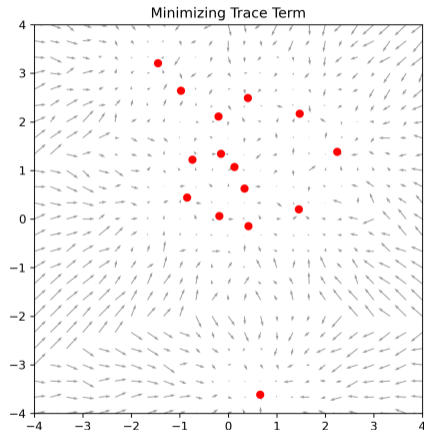
$$\min_{\phi} \mathbb{E}_{x \sim p_{\text{data}}} [\text{Tr}(\nabla_x s_{\phi}(x)) + \|s_{\phi}(x)\|^2]$$



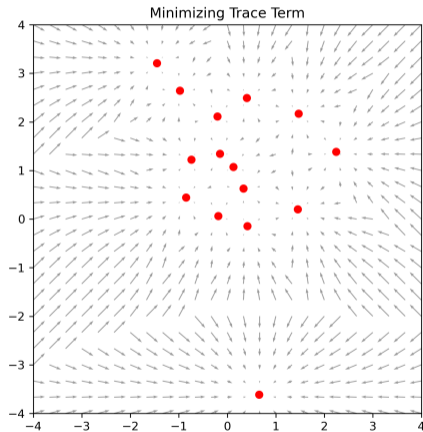
$$\min_{\phi} \mathbb{E}_{x \sim p_{\text{data}}} [\text{Tr}(\nabla_x s_{\phi}(x)) + \|s_{\phi}(x)\|^2]$$



$$\min_{\phi} \mathbb{E}_{x \sim p_{\text{data}}} [\text{Tr}(\nabla_x s_{\phi}(x)) + \|s_{\phi}(x)\|^2]$$



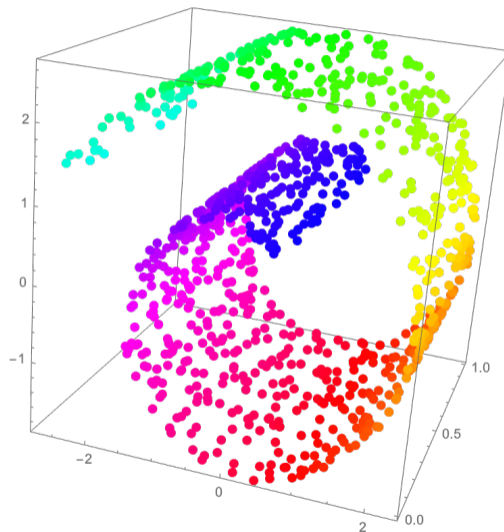
$$\min_{\phi} \mathbb{E}_{x \sim p_{\text{data}}} [\text{Tr}(\nabla_x s_{\phi}(x)) + \|s_{\phi}(x)\|^2]$$



- 1  $\text{Tr}(\nabla_x s_\phi(x))$  is intractable – we would need to compute  $O(n^2)$  derivatives.

- 1  $\text{Tr}(\nabla_x s_\phi(x))$  is intractable – we would need to compute  $O(n^2)$  derivatives.
- 2 Score Matching usually gives accurate results in regions *very close* to the samples.

- 1  $\text{Tr}(\nabla_x s_\phi(x))$  is intractable – we would need to compute  $O(n^2)$  derivatives.
- 2 Score Matching usually gives accurate results in regions *very close* to the samples.
- 3 Manifold Hypothesis – no samples exist outside the manifold.



# Denoising Score Matching

Perturbing data points  $x$  randomly:

$$\hat{x} = x + \sigma\epsilon, \quad \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

---

<sup>4</sup>[Vincent, 2010]

Perturbing data points  $x$  randomly:

$$\hat{x} = x + \sigma\epsilon, \quad \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

Determines a **corruption process**  $q_\sigma(\hat{x}|x) \sim \mathcal{N}(x, \sigma^2\mathbf{I})$ :

$$q_\sigma(\hat{x}) = \int_x q_\sigma(\hat{x}|x) p_{\text{data}}(x) dx$$

---

<sup>4</sup>[Vincent, 2010]

Perturbing data points  $x$  randomly:

$$\hat{x} = x + \sigma\epsilon, \quad \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

Determines a **corruption process**  $q_\sigma(\hat{x}|x) \sim \mathcal{N}(x, \sigma^2\mathbf{I})$ :

$$q_\sigma(\hat{x}) = \int_x q_\sigma(\hat{x}|x) p_{\text{data}}(x) dx$$

It can be proven that:

$$\nabla_{\hat{x}} \log q_\sigma(\hat{x}) \simeq \nabla_x \log p_{\text{data}}(x)$$

---

<sup>4</sup>[Vincent, 2010]

**Problem.** Unfeasible:

$$q_{\sigma}(\hat{x}) = \int_x q_{\sigma}(\hat{x}|x) p_{\text{data}}(x) dx$$

**Problem.** Unfeasible:

$$q_{\sigma}(\hat{x}) = \int_x q_{\sigma}(\hat{x}|x) p_{\text{data}}(x) dx$$

Solution - approximation:

$$q_{\sigma}(\hat{x}) \approx \frac{1}{N} \sum_{i=1}^N q_{\sigma}(\hat{x}|x_i)$$

The objective is:

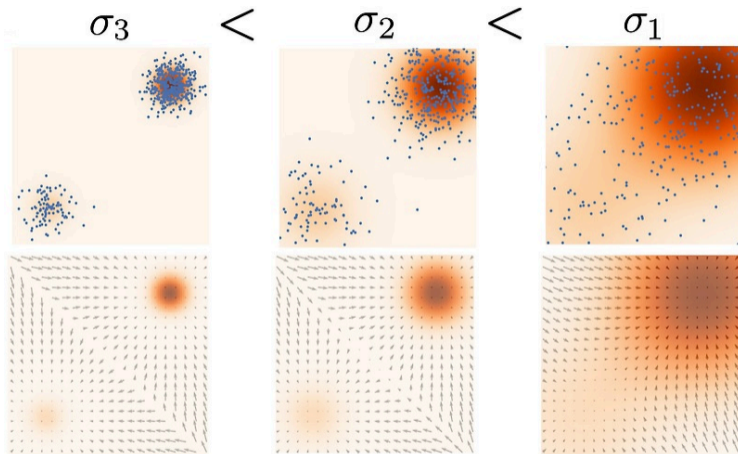
$$\phi^* = \min_{\phi} \mathbb{E}_{x \sim p_{\text{data}}, \hat{x} \sim \mathcal{N}(x, \sigma^2 \mathbf{I})} \left[ \left\| s_{\phi}(\hat{x}) - \nabla_{\hat{x}} \log q_{\sigma}(\hat{x}|x) \right\|^2 \right]$$

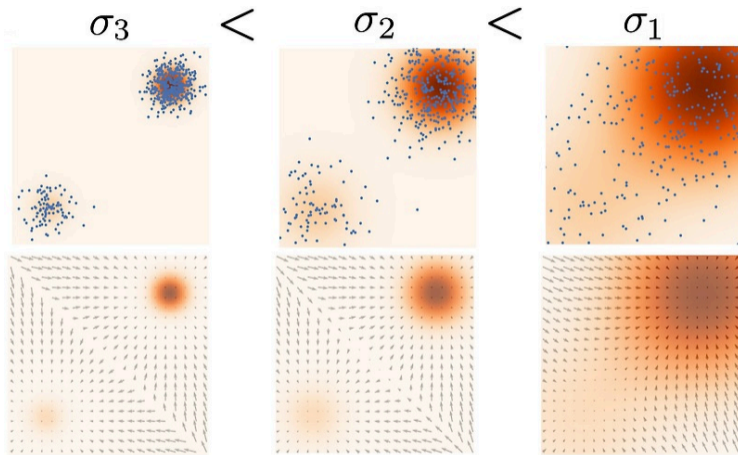
The objective is:

$$\begin{aligned}\phi^* &= \min_{\phi} \mathbb{E}_{x \sim p_{\text{data}}, \hat{x} \sim \mathcal{N}(x, \sigma^2 \mathbf{I})} \left[ \left\| s_{\phi}(\hat{x}) - \nabla_{\hat{x}} \log q_{\sigma}(\hat{x}|x) \right\|^2 \right] \\ &= \min_{\phi} \mathbb{E}_{x \sim p_{\text{data}}, \hat{x} \sim \mathcal{N}(x, \sigma^2 \mathbf{I})} \left[ \left\| s_{\phi}(\hat{x}) - \left( -\frac{\hat{x} - x}{\sigma^2} \right) \right\|^2 \right]\end{aligned}$$

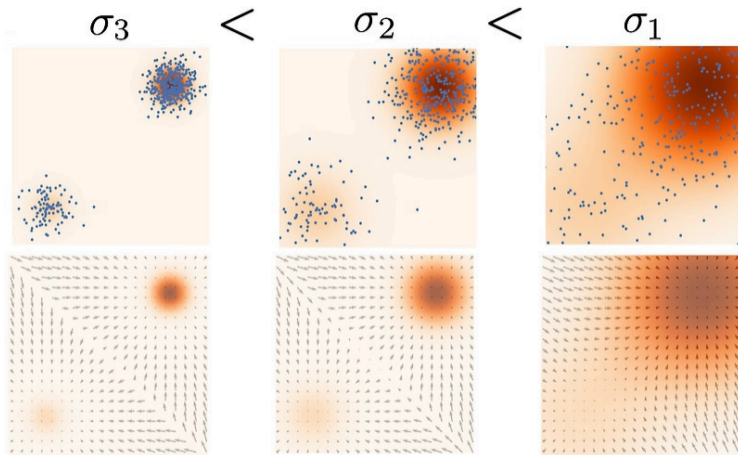
The objective is:

$$\begin{aligned}\phi^* &= \min_{\phi} \mathbb{E}_{x \sim p_{\text{data}}, \hat{x} \sim \mathcal{N}(x, \sigma^2 \mathbf{I})} \left[ \left\| s_{\phi}(\hat{x}) - \nabla_{\hat{x}} \log q_{\sigma}(\hat{x}|x) \right\|^2 \right] \\ &= \min_{\phi} \mathbb{E}_{x \sim p_{\text{data}}, \hat{x} \sim \mathcal{N}(x, \sigma^2 \mathbf{I})} \left[ \left\| s_{\phi}(\hat{x}) - \left( -\frac{\hat{x} - x}{\sigma^2} \right) \right\|^2 \right] \\ &= \min_{\phi} \mathbb{E}_{x \sim p_{\text{data}}, \epsilon \sim \mathcal{N}(0, \mathbf{I})} \left[ \left\| s_{\phi}(x + \sigma\epsilon) + \frac{\epsilon}{\sigma} \right\|^2 \right]\end{aligned}$$



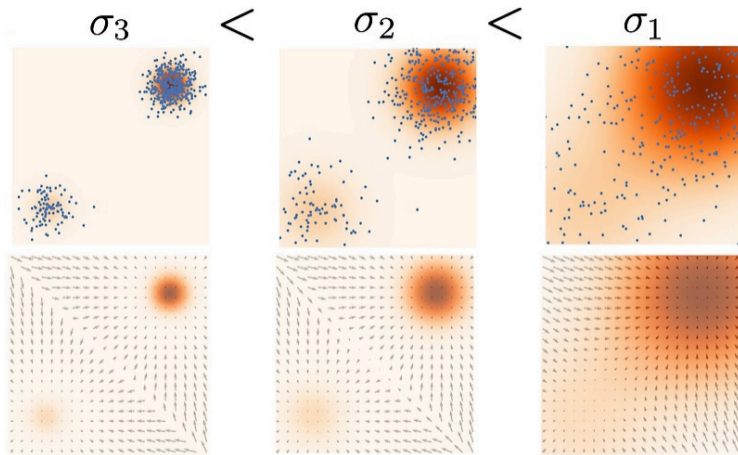


A fixed  $\sigma$  cannot actually yield to an usable model.



A fixed  $\sigma$  cannot actually yield to an usable model.

However,



A fixed  $\sigma$  cannot actually yield to an usable model.  
However, many  $\sigma$  can.

# Noise Conditional Score Networks

Select a sequence of noise levels:  $\{\sigma_i\}_{i=1}^L$ , with  $\sigma_1 > \sigma_2 > \dots > \sigma_L$

---

<sup>4</sup>[Song & Ermon, 2019]

Select a sequence of noise levels:  $\{\sigma_i\}_{i=1}^L$ , with  $\sigma_1 > \sigma_2 > \dots > \sigma_L$

Learn a **score function**  $\mathbf{s}_\phi(\hat{x}, \sigma)$  to approximate  $\nabla_{\hat{x}} \log q_\sigma(\hat{x} | x)$ ,  $\forall \sigma$  and  $\forall x$ .

---

<sup>4</sup>[Song & Ermon, 2019]

Select a sequence of noise levels:  $\{\sigma_i\}_{i=1}^L$ , with  $\sigma_1 > \sigma_2 > \dots > \sigma_L$

Learn a **score function**  $\mathbf{s}_\phi(\hat{x}, \sigma)$  to approximate  $\nabla_{\hat{x}} \log q_\sigma(\hat{x} | x)$ ,  $\forall \sigma$  and  $\forall x$ .



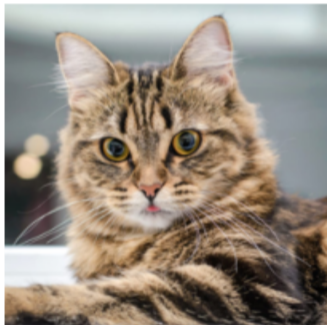
---

<sup>4</sup>[Song & Ermon, 2019]

For a minibatch of size  $N$ , the loss for a noise level  $\sigma_i$  is:

$$\mathcal{L}_i(\phi) = \frac{1}{N} \sum_{n=1}^N \left\| s_\phi(\hat{x}^{(n)}, \sigma_i) + \frac{\hat{x}^{(n)} - x^{(n)}}{\sigma_i^2} \right\|^2$$
$$\mathcal{L}_i(\phi) = \frac{1}{N} \sum_{n=1}^N \left\| s_\phi(x^{(n)} + \sigma_i \epsilon^{(n)}, \sigma_i) + \frac{\epsilon^{(n)}}{\sigma_i} \right\|^2$$

*x*  $\sim$  *D*



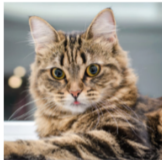
~ D



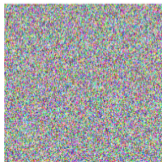
$$\sigma \sim \{\sigma_1, \dots, \sigma_L\}$$

 $\sigma$ 

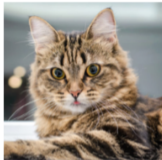
$$\epsilon \sim \mathcal{N}(0, I)$$



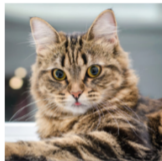
$\sigma$



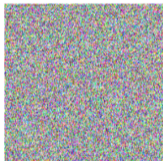
$$\sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

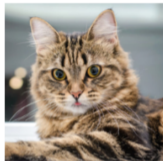


$$\sigma \times \epsilon$$

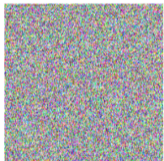


$\sigma \times$





$\sigma \times$

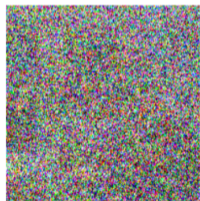


$=$

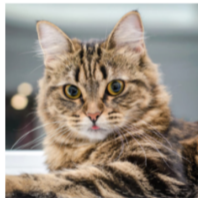


$$\hat{x} = x + \sigma \epsilon$$

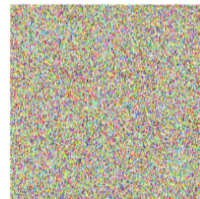
$$\hat{x} = \text{cat} + \text{noise}$$

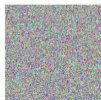



=

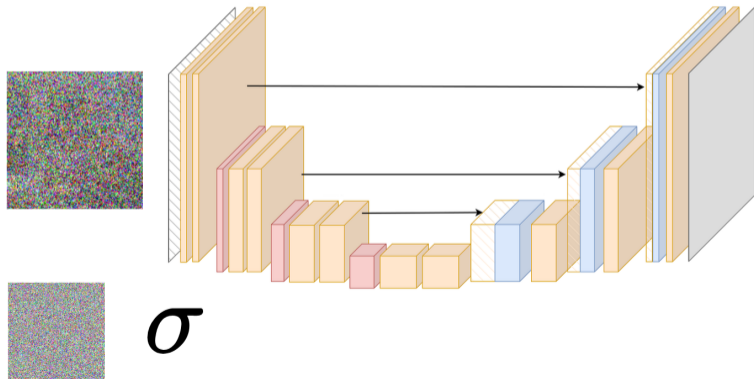


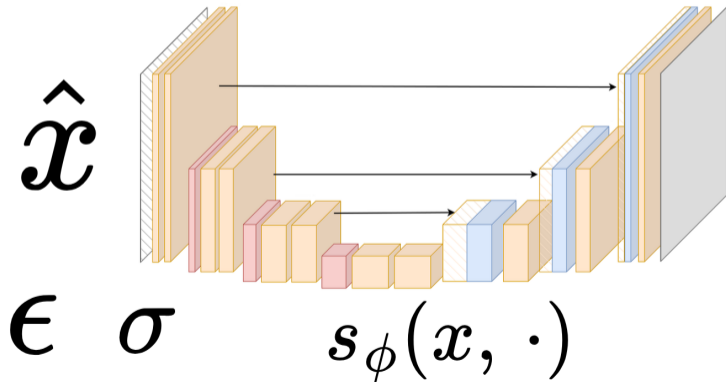
+

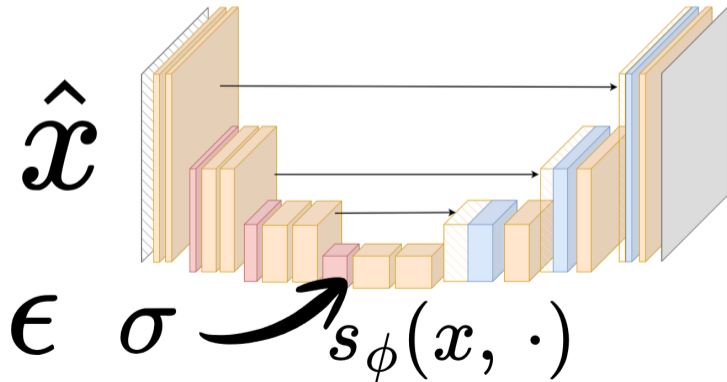


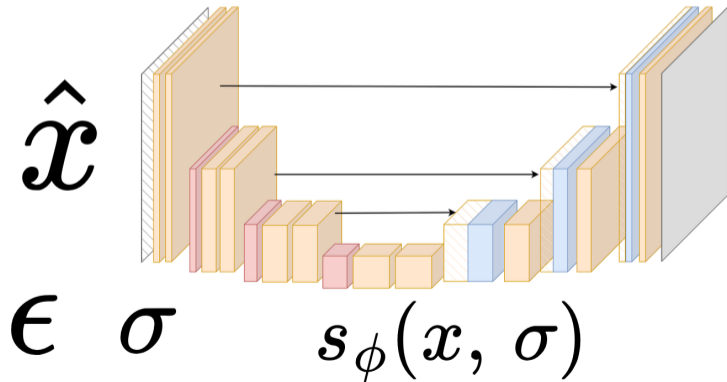


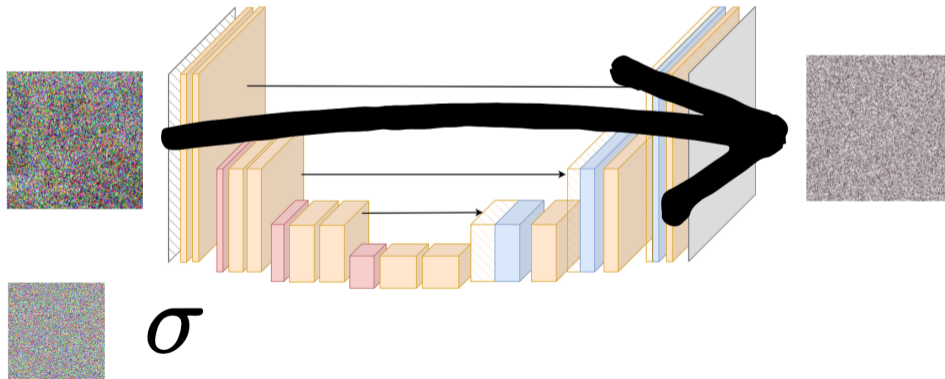
$\sigma$

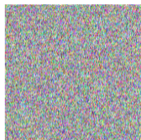






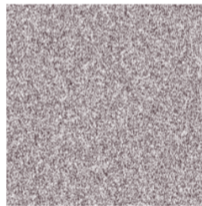






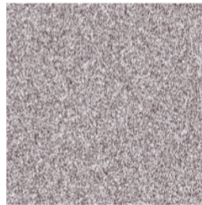
12

12?

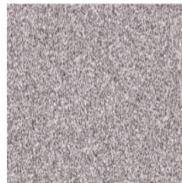
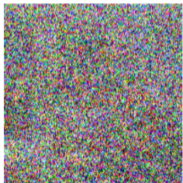


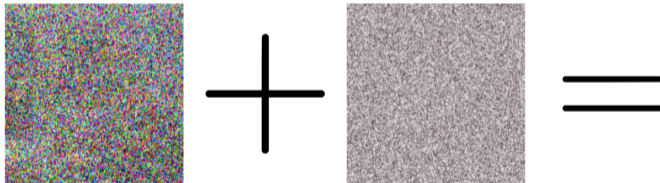
$\epsilon$   
/  
 $\sigma$

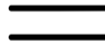
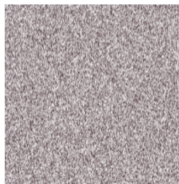
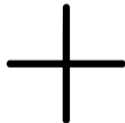
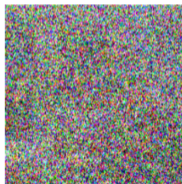
12?











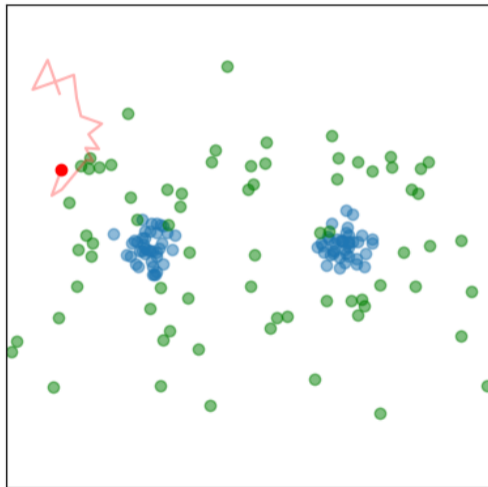
## Annealed Langevin Dynamics.

- 1: Initialize  $x_0$
- 2: **for**  $i \leftarrow 1$  to  $L$  **do**
- 3:      $\alpha_i \leftarrow \epsilon \cdot \sigma_i^2 / \sigma_L^2$
- 4:     **for**  $t \leftarrow 1$  to  $T$  **do**
- 5:          $\mathbf{z}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 6:          $x_t \leftarrow x_{t-1} + \alpha_i \mathbf{s}(x_{t-1}, \sigma_i) + \sqrt{\alpha_i} \mathbf{z}_t$
- 7:     **end for**
- 8:      $x_0 \leftarrow x_T$
- 9: **end for**
- 10: **return**  $x_L$

## Annealed Langevin Dynamics.

- 1: Initialize  $x_0$
- 2: **for**  $i \leftarrow 1$  to  $L$  **do**
- 3:      $\alpha_i \leftarrow \epsilon \cdot \sigma_i^2 / \sigma_L^2$
- 4:     **for**  $t \leftarrow 1$  to  $T$  **do**
- 5:          $\mathbf{z}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 6:          $x_t \leftarrow x_{t-1} + \alpha_i \mathbf{s}(x_{t-1}, \sigma_i) + \sqrt{\alpha_i} \mathbf{z}_t$
- 7:     **end for**
- 8:      $x_0 \leftarrow x_T$
- 9: **end for**
- 10: **return**  $x_L$

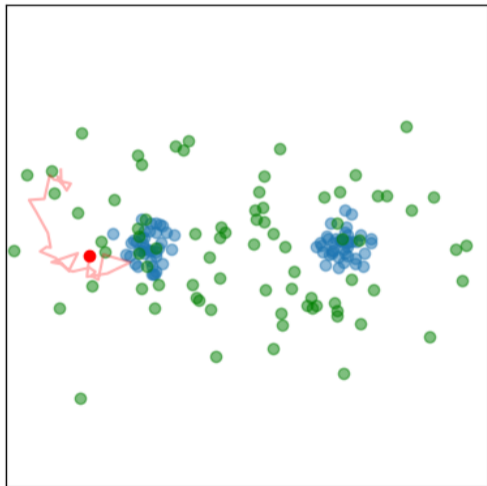
Annealed Langevin Stage (sigma = 2.0)



## Annealed Langevin Dynamics.

- 1: Initialize  $x_0$
- 2: **for**  $i \leftarrow 1$  to  $L$  **do**
- 3:      $\alpha_i \leftarrow \epsilon \cdot \sigma_i^2 / \sigma_L^2$
- 4:     **for**  $t \leftarrow 1$  to  $T$  **do**
- 5:          $\mathbf{z}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 6:          $x_t \leftarrow x_{t-1} + \alpha_i \mathbf{s}(x_{t-1}, \sigma_i) + \sqrt{\alpha_i} \mathbf{z}_t$
- 7:     **end for**
- 8:      $x_0 \leftarrow x_T$
- 9: **end for**
- 10: **return**  $x_L$

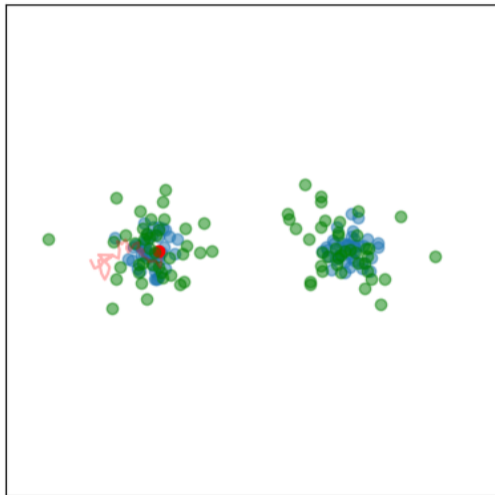
Annealed Langevin Stage (sigma = 1.25)



## Annealed Langevin Dynamics.

- 1: Initialize  $x_0$
- 2: **for**  $i \leftarrow 1$  to  $L$  **do**
- 3:      $\alpha_i \leftarrow \epsilon \cdot \sigma_i^2 / \sigma_L^2$
- 4:     **for**  $t \leftarrow 1$  to  $T$  **do**
- 5:          $\mathbf{z}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 6:          $x_t \leftarrow x_{t-1} + \alpha_i \mathbf{s}(x_{t-1}, \sigma_i) + \sqrt{\alpha_i} \mathbf{z}_t$
- 7:     **end for**
- 8:      $x_0 \leftarrow x_T$
- 9: **end for**
- 10: **return**  $x_L$

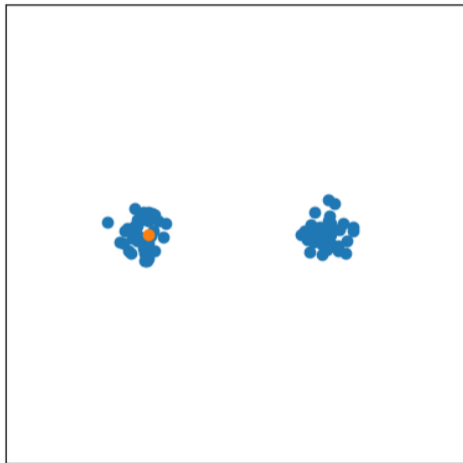
## Annealed Langevin Stage (sigma = 0.5)



## Annealed Langevin Dynamics.

- 1: Initialize  $x_0$
- 2: **for**  $i \leftarrow 1$  to  $L$  **do**
- 3:      $\alpha_i \leftarrow \epsilon \cdot \sigma_i^2 / \sigma_L^2$
- 4:     **for**  $t \leftarrow 1$  to  $T$  **do**
- 5:          $\mathbf{z}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 6:          $x_t \leftarrow x_{t-1} + \alpha_i \mathbf{s}(x_{t-1}, \sigma_i) + \sqrt{\alpha_i} \mathbf{z}_t$
- 7:     **end for**
- 8:      $x_0 \leftarrow x_T$
- 9: **end for**
- 10: **return**  $x_L$

## Final Sample After Annealed Langevin Dynamics



## Training

Input: dataset  $D$ , noise levels  $\sigma_1 > \sigma_2 > \dots > \sigma_L$

Initialize network parameters  $\theta$

while not converged do:

$x \leftarrow$  sample from  $D$

$\sigma \leftarrow$  sample uniformly from  $\sigma_1, \dots, \sigma_L$

$\epsilon \leftarrow$  sample from  $\mathcal{N}(\mathbf{0}, \mathbf{I})$

$\hat{x} \leftarrow x + \sigma\epsilon$

$y \leftarrow -\epsilon/\sigma$

$\mathcal{L} \leftarrow \|s_\theta(\hat{x}, \sigma) - y\|^2$

    update network parameters  $\theta$

end while

**Stefano Ermon**

**Yang Song**





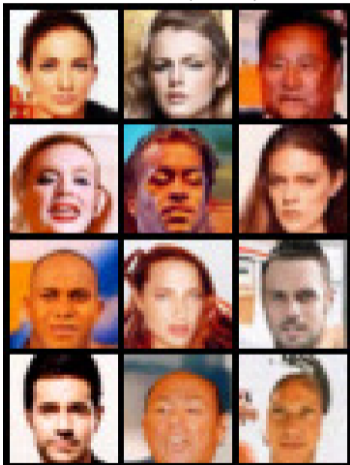
**Stefano Ermon**

**Yang Song**

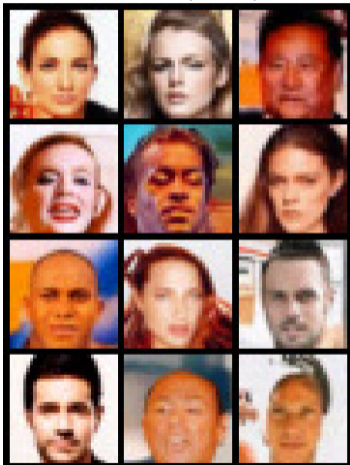
**Gianfranco Bilardi**



## NCSN (2019)



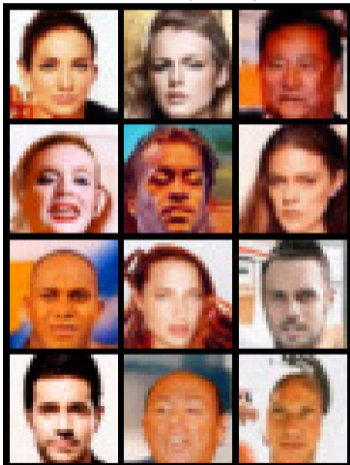
## NCSN (2019)



## DDPM (2020)



## NCSN (2019)



## DDPM (2020)



## Latent Diffusion (2022)



## Foundations

Hinton, G. E., & Sejnowski, T. J. (1985). *Learning and Relearning in Boltzmann Machines*.  
LeCun, Y., Chopra, S., & Hadsell, R. (2006). *A Tutorial on Energy-Based Learning*.

## Score Matching

Hyvärinen, A. (2005). *Estimation of Non-Normalized Statistical Models by Score Matching*.  
Vincent, P. (2010). *A Connection Between Score Matching and Denoising Autoencoders*.

## Noise Conditional Score Networks

Song, Y. & Ermon, S. (2019). *Generative Modeling by Estimating Gradients of the Data Distribution*.

**Thank you!**

$$\phi^* = \min_{\phi} \mathbb{E}_{p_{\text{data}}} [\|s_{\phi}(x) - s_{\text{data}}(x)\|^2]$$

$$\phi^* = \min_{\phi} \mathbb{E}_{p_{\text{data}}} [\|s_{\phi}(x) - s_{\text{data}}(x)\|^2]$$

$$\|s_{\phi}(x) - s_{\text{data}}(x)\|^2 = s_{\phi}(x)^2 + s_{\text{data}}(x)^2 - 2 s_{\phi}(x) s_{\text{data}}(x)$$

$$s_{\text{data}}(x) = \nabla_x \log p_{\text{data}}(x)$$

$$\begin{aligned}\mathbb{E}_{p_{\text{data}}} \left[ -2 s_{\phi}(x) s_{\text{data}}(x) \right] &= -2 \int_{-\infty}^{\infty} s_{\phi}(x) s_{\text{data}}(x) p_{\text{data}}(x) dx \\ &= -2 \int_{-\infty}^{\infty} s_{\phi}(x) \frac{\nabla_x p_{\text{data}}(x)}{p_{\text{data}}(x)} p_{\text{data}}(x) dx \\ &= -2 \int_{-\infty}^{\infty} s_{\phi}(x) \nabla_x p_{\text{data}}(x) dx\end{aligned}$$

$$\begin{aligned} & -2 \int_{-\infty}^{\infty} s_{\phi}(x) \nabla_x p_{\text{data}}(x) dx \\ &= - \int_{-\infty}^{\infty} s_{\phi}(x) \nabla_x p_{\text{data}}(x) dx \\ &= - \left[ s_{\phi}(x) p_{\text{data}}(x) \right] \Big|_{-\infty}^{\infty} + \int_{-\infty}^{\infty} \text{div}(s_{\phi}(x)) p_{\text{data}}(x) dx \\ &= - \cancel{\left[ s_{\phi}(x) p_{\text{data}}(x) \right] \Big|_{-\infty}^{\infty}} + \mathbb{E}_{p_{\text{data}}} [\text{div}(s_{\phi}(x))] \end{aligned}$$