

# Introduction to Score based Generative modelling

Gabriel V. Cardoso

Hi! Paris reading group, 27/02/2024

# Overview

- 1 Background
- 2 Implicit Score Matching
- 3 Denoising score matching
- 4 Noise Conditional Score Networks
- 5 Denoising diffusion implicit models (DDIM)
- 6 Deep image prior
- 7 Conclusion

# Background

---

# Generative Models

- **Task:** generate new samples from a distribution of interest  $q_d$  defined over  $\mathbb{R}^d$ .
- **Context:** We rely only on a dataset  $\mathcal{D}$  of i.i.d samples from  $q_d$ .
- **Examples:** Generative Adversarial Networks (GANs)<sup>1</sup>, Normalizing Flows<sup>2</sup> and *Score-based generative models*<sup>3</sup>.

---

<sup>1</sup>Goodfellow, I. J., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., & Bengio, Y. (2014). Generative adversarial nets. *Proceedings of the 27th International Conference on Neural Information Processing Systems - Volume 2*, 2672–2680.

<sup>2</sup>Rezende, D., & Mohamed, S. (2015, July). Variational inference with normalizing flows. In F. Bach & D. Blei (Eds.), *Proceedings of the 32nd international conference on machine learning* (pp. 1530–1538, Vol. 37). PMLR. <https://proceedings.mlr.press/v37/rezende15.html>

<sup>3</sup>Song, Y., & Ermon, S. (2019). Generative modeling by estimating gradients of the data distribution. *Advances in neural information processing systems*, 32.

# Definitions

- Let  $\mathcal{B}(\mathbb{R}^d)$  be the set of Borelian subsets of  $\mathbb{R}^d$ .

# Definitions

- Let  $\mathcal{B}(\mathbb{R}^d)$  be the set of Borelian subsets of  $\mathbb{R}^d$ .
- Let  $d \in \mathbb{N}_*$  and  $\mathcal{P}_0(\mathbb{R}^d)$  denote the set of (Borelian) probability measures in  $\mathbb{R}^d$ .

# Definitions

- Let  $\mathcal{B}(\mathbb{R}^d)$  be the set of Borelian subsets of  $\mathbb{R}^d$ .
- Let  $d \in \mathbb{N}_*$  and  $\mathcal{P}_0(\mathbb{R}^d)$  denote the set of (Borelian) probability measures in  $\mathbb{R}^d$ .
- $\mathcal{P}_2(\mathbb{R}^d) := \{p \in \mathcal{P}_0(\mathbb{R}^d) \mid \mathbb{E}_{X \sim p} [X^2] < \infty\}$ .

# Definitions

- Let  $\mathcal{B}(\mathbb{R}^d)$  be the set of Borelian subsets of  $\mathbb{R}^d$ .
- Let  $d \in \mathbb{N}_*$  and  $\mathcal{P}_0(\mathbb{R}^d)$  denote the set of (Borelian) probability measures in  $\mathbb{R}^d$ .
- $\mathcal{P}_2(\mathbb{R}^d) := \{p \in \mathcal{P}_0(\mathbb{R}^d) \mid \mathbb{E}_{X \sim p} [X^2] < \infty\}$ .
- For  $q_d \in \mathcal{P}_0(\mathbb{R}^d)$  that admits a density w.r.t the Lebesgue measure, we define the *score* of  $q_d$  as

$$s(x) := \nabla \log q_d(x).$$



# Unadjusted Langevin Algorithm (ULA)

## ULA

$X_0 \sim \mu_0$ , for  $t \in \mathbb{N}_*$

$$X_t := X_{t-1} + \gamma \nabla \log \mathbf{q}_d(X_{t-1}) + (2\gamma)^{1/2} \epsilon_t, \quad (1)$$

where  $\epsilon_t \sim \mathcal{N}(0, I_d)$  and  $\gamma > 0$ .

# ULA guarantees

## Wasserstein 2

For  $(p_1, p_2) \in \mathbf{P}_2(\mathbb{R}^d)^{\otimes 2}$  we define

$$\mathcal{C}(p_1, p_2) := \left\{ \pi \in \mathbf{P}_0(\mathbb{R}^{2d}) \mid \begin{aligned} \pi(A \times \mathbb{R}^d) &= p_1(A); \pi(\mathbb{R}^d \times B) = p_2(B) \\ \text{for } (A, B) &\in \mathcal{B}(\mathbb{R}^d)^{\otimes 2} \end{aligned} \right\}.$$

We define the *Wasserstein 2* distance between  $p_1$  and  $p_2$  as

$$W_2^2(p_1, p_2) := \min_{\pi \in \mathcal{C}(p_1, p_2)} \int \|x - y\|^2 \pi(x, y) dx dy.$$

# ULA guarantees

## ULA guarantees from Durmus et al., 2019, Corollary 10<sup>4</sup>

Assume the score is  $m$ -concave,  $L$  Lipschitz and  $\epsilon > 0$ . Let  $\mu_t := \text{Law}(X_t)$ . If

- $\gamma_\epsilon < \min \{m\epsilon/(4Ld), L^{-1}\}$  and  
 $t_\epsilon > \log(2W_2^2(\mu_0, \mathbf{q}_d)/\epsilon)\gamma_\epsilon^{-1}m^{-1}$

then

$$W_2^2(\mu_{t_\epsilon}, \mathbf{q}_d) < \epsilon.$$

---

<sup>4</sup>Durmus, A., Majewski, S., & Miasojedow, B. (2019). Analysis of langevin monte carlo via convex optimization. *Journal of Machine Learning Research*, 20(73), 1–46. <http://jmlr.org/papers/v20/18-173.html> Corollary 10.

## Implicit Score Matching

---

# Score Matching

Goal: Learn the score of  $q_d$  with a Neural Network  $s_\theta$ , where  $\theta \in \Theta \subset \mathbb{R}^p$ .

# Score Matching

Goal: Learn the score of  $q_d$  with a Neural Network  $s_\theta$ , where  $\theta \in \Theta \subset \mathbb{R}^p$ .

## Score Matching

$$\operatorname{argmin}_\theta \mathbb{E}_{X \sim q_d} [\|s_\theta(X) - s(X)\|^2] . \quad (2)$$

# Learning the score from data

Hyvärinen, 2005<sup>5</sup> shows if  $\lim_{\|x\| \rightarrow \infty} \mathbf{s}_\theta(x) \mathbf{q}_d(x) = 0$ , then the score matching objective (2) is equivalent to

## Implicit score matching loss

$$\operatorname{argmin}_\theta \mathbb{E}_{X \sim \mathbf{q}_d} \left[ \nabla \cdot \mathbf{s}_\theta(X) + 1/2 \|\mathbf{s}_\theta(X)\|^2 \right]. \quad (3)$$

---

<sup>5</sup>Hyvärinen, A. (2005). Estimation of non-normalized statistical models by score matching. *Journal of Machine Learning Research*, 6(24), 695–709.

<http://jmlr.org/papers/v6/hyvarinen05a.html>

## Proof

$$\begin{aligned}\mathbb{E}_{X \sim q_d} [\|\mathbf{s}(X) - \mathbf{s}_\theta(X)\|^2] &= \mathbb{E}_{X \sim q_d} [\|\mathbf{s}(X)\|^2] \\ &\quad - 2\mathbb{E}_{X \sim q_d} [\mathbf{s}(X)^T \mathbf{s}_\theta(X)] + \mathbb{E}_{X \sim q_d} [\|\mathbf{s}_\theta(X)\|^2] .\end{aligned}$$



## Proof

$$\begin{aligned}\mathbb{E}_{X \sim \mathbf{q}_d} [\|\mathbf{s}(X) - \mathbf{s}_\theta(X)\|^2] &= \mathbb{E}_{X \sim \mathbf{q}_d} [\|\mathbf{s}(X)\|^2] \\ &\quad - 2\mathbb{E}_{X \sim \mathbf{q}_d} [\mathbf{s}(X)^T \mathbf{s}_\theta(X)] + \mathbb{E}_{X \sim \mathbf{q}_d} [\|\mathbf{s}_\theta(X)\|^2] .\end{aligned}$$

Note that

$$\mathbf{s}(x)^T \mathbf{s}_\theta(x) = \sum_{i=1}^d \mathbf{s}_{\theta,i}(x) \partial_{x_i} \log \mathbf{q}_d(x) .$$

## Proof

$$\begin{aligned}\mathbb{E}_{X \sim \mathbf{q}_d} [\|\mathbf{s}(X) - \mathbf{s}_\theta(X)\|^2] &= \mathbb{E}_{X \sim \mathbf{q}_d} [\|\mathbf{s}(X)\|^2] \\ &\quad - 2\mathbb{E}_{X \sim \mathbf{q}_d} [\mathbf{s}(X)^T \mathbf{s}_\theta(X)] + \mathbb{E}_{X \sim \mathbf{q}_d} [\|\mathbf{s}_\theta(X)\|^2] .\end{aligned}$$

Note that

$$\mathbf{s}(x)^T \mathbf{s}_\theta(x) = \sum_{i=1}^d \mathbf{s}_{\theta,i}(x) \partial_{x_i} \log \mathbf{q}_d(x) .$$

For  $i \in \llbracket 1, d \rrbracket$ ,

$$\mathbb{E}_{X \sim \mathbf{q}_d} [\mathbf{s}_{\theta,i}(X) \partial_{x_i} \log \mathbf{q}_d(X)] = \int \mathbf{s}_{\theta,i}(x) \partial_{x_i} \log \mathbf{q}_d(x) \mathbf{q}_d(x) dx .$$

## Proof

$$\begin{aligned}\mathbb{E}_{X \sim \mathbf{q}_d} [\|\mathbf{s}(X) - \mathbf{s}_\theta(X)\|^2] &= \mathbb{E}_{X \sim \mathbf{q}_d} [\|\mathbf{s}(X)\|^2] \\ &\quad - 2\mathbb{E}_{X \sim \mathbf{q}_d} [\mathbf{s}(X)^T \mathbf{s}_\theta(X)] + \mathbb{E}_{X \sim \mathbf{q}_d} [\|\mathbf{s}_\theta(X)\|^2] .\end{aligned}$$

Note that

$$\mathbf{s}(x)^T \mathbf{s}_\theta(x) = \sum_{i=1}^d \mathbf{s}_{\theta,i}(x) \partial_{x_i} \log \mathbf{q}_d(x) .$$

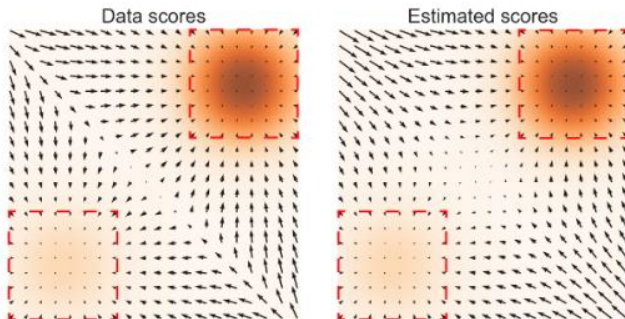
For  $i \in \llbracket 1, d \rrbracket$ ,

$$\mathbb{E}_{X \sim \mathbf{q}_d} [\mathbf{s}_{\theta,i}(X) \partial_{x_i} \log \mathbf{q}_d(X)] = \int \mathbf{s}_{\theta,i}(x) \partial_{x_i} \log \mathbf{q}_d(x) \mathbf{q}_d(x) dx .$$

Since  $\partial_{x_i} \log \mathbf{q}_d(x) \mathbf{q}_d(x) = \partial_{x_i} \mathbf{q}_d(x)$ , we have

$$\int \mathbf{s}_{\theta,i}(x) \partial_{x_i} \log \mathbf{q}_d(x) \mathbf{q}_d(x) dx = - \int \partial_{x_i} \mathbf{s}_{\theta,i}(x) \mathbf{q}_d(x) dx .$$

- $\nabla \cdot s_\theta$  costly in high dimensions.
- Score estimate inaccurate in low data regions, furthermore ULA moves can be stuck in each mode.



**Figure:** Illustration of score on low data regions, from Y. Song and Ermon, 2019<sup>6</sup>.

---

<sup>6</sup>Song, Y., & Ermon, S. (2019). Generative modeling by estimating gradients of the data distribution. *Advances in neural information processing systems*, 32.

Code Break!

https:

//github.com/gabrielvc/tutorial\_ddim

## Denoising score matching

---

## Denosing score matching

Vincent, 2011<sup>7</sup> introduces the idea of learning the score of  $q_\sigma(dx_\sigma) = \int q_\sigma(x_\sigma|x)q_d(dx)$  where  $q_\sigma(\cdot|x) = \mathcal{N}(x, \sigma^2 I)$ .

---

<sup>7</sup>Vincent, P. (2011). A connection between score matching and denosing autoencoders. *Neural Computation*, 23(7), 1661–1674.

[https://doi.org/10.1162/NECO\\_a\\_00142](https://doi.org/10.1162/NECO_a_00142)

# Denosing score matching

Vincent, 2011<sup>7</sup> introduces the idea of learning the score of  $q_\sigma(dx_\sigma) = \int q_\sigma(x_\sigma|x)q_d(dx)$  where  $q_\sigma(\cdot|x) = \mathcal{N}(x, \sigma^2 I)$ .

## Denosing score matching loss

$$\operatorname{argmin}_\theta \mathbb{E}_{X_\sigma \sim q_\sigma(\cdot|X), X \sim q_d} [\|s_{\theta, \sigma}(X_\sigma) - \nabla \log q_\sigma(X_\sigma|X)\|^2] .$$

---

<sup>7</sup>Vincent, P. (2011). A connection between score matching and denoising autoencoders. *Neural Computation*, 23(7), 1661–1674.



## Proof

Let  $q_{0|\sigma}(dx_0|x_\sigma) := q_{\sigma|0}(x_\sigma|x_0)q_d(dx_0)/q_\sigma(x_\sigma)$ .

## Proof

Let  $q_{0|\sigma}(dx_0|x_\sigma) := q_{\sigma|0}(x_\sigma|x_0)\mathbf{q}_d(dx_0)/\mathbf{q}_\sigma(x_\sigma)$ . By Fisher's identity

$$\begin{aligned}\nabla \log \mathbf{q}_\sigma(x_\sigma) &= \frac{\nabla \mathbf{q}_\sigma(x_\sigma)}{\mathbf{q}_\sigma(x_\sigma)} = \mathbb{E}_{X_0 \sim \mathbf{q}_d} \left[ \frac{\nabla q_{\sigma|0}(x_\sigma|X_0)}{\mathbf{q}_\sigma(x_\sigma)} \right] \\ &= \mathbb{E}_{X_0 \sim \mathbf{q}_d} \left[ \nabla \log q_{\sigma|0}(x_\sigma|X_0) \frac{q_{\sigma|0}(x_\sigma|X_0)}{\mathbf{q}_\sigma(x_\sigma)} \right] \\ &= \mathbb{E}_{X_0 \sim q_{0|\sigma}(\cdot|x_\sigma)} \left[ \nabla \log q_{\sigma|0}(x_\sigma|X_0) \right] .\end{aligned}$$

## Proof

Let  $q_{0|\sigma}(dx_0|x_\sigma) := q_{\sigma|0}(x_\sigma|x_0)\mathbf{q}_d(dx_0)/\mathbf{q}_\sigma(x_\sigma)$ . By Fisher's identity

$$\begin{aligned}\nabla \log \mathbf{q}_\sigma(x_\sigma) &= \frac{\nabla \mathbf{q}_\sigma(x_\sigma)}{\mathbf{q}_\sigma(x_\sigma)} = \mathbb{E}_{X_0 \sim \mathbf{q}_d} \left[ \frac{\nabla q_{\sigma|0}(x_\sigma|X_0)}{\mathbf{q}_\sigma(x_\sigma)} \right] \\ &= \mathbb{E}_{X_0 \sim \mathbf{q}_d} \left[ \nabla \log q_{\sigma|0}(x_\sigma|X_0) \frac{q_{\sigma|0}(x_\sigma|X_0)}{\mathbf{q}_\sigma(x_\sigma)} \right] \\ &= \mathbb{E}_{X_0 \sim q_{0|\sigma}(\cdot|x_\sigma)} \left[ \nabla \log q_{\sigma|0}(x_\sigma|X_0) \right].\end{aligned}$$

Thus

$$\begin{aligned}\mathbb{E}_{X_\sigma \sim \mathbf{q}_\sigma} \left[ \|\mathbf{s}_\theta(X_\sigma) - \nabla \log \mathbf{q}_\sigma(X_\sigma)\|^2 \right] \\ = \mathbb{E}_{X_\sigma \sim \mathbf{q}_\sigma} \left[ \left\| \mathbf{s}_\theta(X_\sigma) - \mathbb{E}_{X_0 \sim q_{0|\sigma}(\cdot|X_\sigma)} \left[ \nabla \log q_{\sigma|0}(X_\sigma|X_0) \right] \right\|^2 \right].\end{aligned}$$

By defining

$\mathbf{q}_{\sigma,0}(\mathrm{d}x_\sigma, \mathrm{d}x_0) = q_{\sigma|0}(x_0|\mathrm{d}x_\sigma)\mathbf{q}_d(\mathrm{d}x_0) = q_{0|\sigma}(\mathrm{d}x_0|x_\sigma)\mathbf{q}_\sigma(\mathrm{d}x_\sigma)$ , we have

$$\begin{aligned} & \mathbb{E}_{X_\sigma \sim \mathbf{q}_\sigma} [\|\mathbf{s}_{\theta,\sigma}(X_\sigma) - \nabla \log \mathbf{q}_\sigma(X_\sigma)\|^2] \\ &= \mathbb{E}_{X_\sigma \sim \mathbf{q}_\sigma} \left[ \|\mathbf{s}_\theta(X_\sigma)\|^2 - 2\mathbf{s}_\theta(X_\sigma)^T \mathbb{E}_{X_0 \sim q_{0|\sigma}(\cdot|X_\sigma)} [\nabla \log q_{\sigma|0}(X_\sigma|X_0)] \right] + \\ &= \mathbb{E}_{X_\sigma \sim \mathbf{q}_\sigma} \left[ \mathbb{E}_{X_0 \sim q_{0|\sigma}(\cdot|X_\sigma)} [\|\mathbf{s}_\theta(X_\sigma) - \nabla \log q_{\sigma|0}(X_\sigma|X_0)\|^2] \right] + \tilde{C} \\ &= \mathbb{E}_{(X_\sigma, X_0) \sim \mathbf{q}_{\sigma,0}} [\|\mathbf{s}_\theta(X_\sigma) - \nabla \log q_{\sigma|0}(X_\sigma|X_0)\|^2] + \tilde{C}, \end{aligned}$$

where  $C$  and  $\tilde{C}$  are constants that do not depend on  $\theta$ .

## Denoising score matching

- No need of calculating derivatives of the score network.
- Mixes better and exploit regions of low data density.
- Approaches  $\nabla \log q_d$  only in the limit  $\sigma \rightarrow 0$ .

Code Break!

https:

//github.com/gabrielvc/tutorial\_ddim

## Noise Conditional Score Networks

---

# Noise Conditional Score Networks

Y. Song and Ermon, 2019<sup>8</sup> introduces several noised versions of  $\mathbf{q}_d$ .

## Diffused marginals

For  $t \in \llbracket 1, T \rrbracket$  and  $v_t > 0$ , define  $q_{t|0}(x_t|x_0) = \mathcal{N}(x_t; x_0, v_t^2 \mathbf{I})$  and

$$\mathbf{q}_t(dx_t) := \int q_{t|0}(dx_t|x_0) \mathbf{q}_d(dx_0).$$

---

<sup>8</sup>Song, Y., & Ermon, S. (2019). Generative modeling by estimating gradients of the data distribution. *Advances in neural information processing systems*, 32.



# Noise Conditional Score Networks

Train a neural network  $s_\theta$  to jointly learn the score of  $\{q_t\}_{t=1}^T$ :

## Diffusion score matching

$$\sum_{t=1}^T \gamma_t^2 \mathbb{E}_{X_t \sim q_{t|0}(\cdot|X_0), X_0 \sim q_d} [\|s_\theta(X_t, v_t) - \nabla \log q_{t|0}(X_t|X_0)\|^2] .$$

where  $\{v_t\}_{t=1}^T$  is an increasing sequence of positive values.

# Noise Conditional Score Networks

Generate samples by sequential ULA on  $\{q_t\}_{t=1}^T$ :

---

---

**Data :**  $X_T^0, k, r, \theta$

**Result :**  $X_0^0$

```
1 for  $t \leftarrow T$  to 1 do
2   for  $\ell \leftarrow 1$  to  $k$  do
3     set  $\gamma = rv_t^2/v_T^2$ 
4     draw  $\epsilon_{t,\ell} \sim \mathcal{N}(0, I_d)$ 
5     set  $X_t^\ell = X_t^{\ell-1} + (\gamma/2)\mathfrak{s}_\theta(X_t^{\ell-1}, v_t) + \gamma^{1/2}\epsilon_{t,\ell}$ 
6   set  $X_{t-1}^0 = X_t^\ell$ .
```

---

Code Break!

https:

//github.com/gabrielvc/tutorial\_ddim

## Denoising diffusion implicit models (DDIM)

---

## Score based generative models

Other than sequential ULA, several samplers are available to sample backwards from the sequence of distributions  $\{q_t\}_{t=1}^T$ , based on

- Stochastic differential equations<sup>9</sup>,
- Ordinary differential equations<sup>10</sup>,
- Markov chains<sup>11</sup>.

---

<sup>9</sup>Song, Y., Sohl-Dickstein, J., Kingma, D. P., Kumar, A., Ermon, S., & Poole, B. (2021). Score-based generative modeling through stochastic differential equations. *International Conference on Learning Representations*.

<https://openreview.net/forum?id=PxtIG12RRHS>

<sup>10</sup>Karras, T., Aittala, M., Aila, T., & Laine, S. (2022). Elucidating the design space of diffusion-based generative models. *Proc. NeurIPS*.

<sup>11</sup>Ho, J., Jain, A., & Abbeel, P. (2020). Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33, 6840–6851; Song, J., Meng, C., & Ermon, S. (2021). Denoising diffusion implicit models. *International Conference on Learning Representations*.

<https://openreview.net/forum?id=St1giarCHLP>

Define  $X_0 \sim \mathbf{q}_d$ ,  $X_t = X_{t-1} + (v_t^2 - v_{t-1}^2)^{1/2} \varepsilon_t$  for  $t \in \llbracket 1, T \rrbracket$  with  $\varepsilon_t \sim \mathcal{N}(0, \mathbf{I})$ .

---

<sup>12</sup>Ho, J., Jain, A., & Abbeel, P. (2020). Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33, 6840–6851.

## DDPM<sup>12</sup>

Define  $X_0 \sim \mathbf{q}_d$ ,  $X_t = X_{t-1} + (v_t^2 - v_{t-1}^2)^{1/2} \varepsilon_t$  for  $t \in \llbracket 1, T \rrbracket$  with  $\varepsilon_t \sim \mathcal{N}(0, \mathbf{I})$ .

Then,  $\text{Law}(X_t) = \mathbf{q}_t$  and  $\text{Law}(X_t | X_0 = x_0) = q_{t|0}(\cdot | x_0)$ .

---

<sup>12</sup>Ho, J., Jain, A., & Abbeel, P. (2020). Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33, 6840–6851.

Define  $X_0 \sim \mathbf{q}_d$ ,  $X_t = X_{t-1} + (v_t^2 - v_{t-1}^2)^{1/2} \varepsilon_t$  for  $t \in \llbracket 1, T \rrbracket$  with  $\varepsilon_t \sim \mathcal{N}(0, \mathbf{I})$ .

Then,  $\text{Law}(X_t) = \mathbf{q}_t$  and  $\text{Law}(X_t | X_0 = x_0) = q_{t|0}(\cdot | x_0)$ .

Furthermore, we can write the law of  $X_{t-1}$  conditionally on  $X_t, X_0$ ,

---

<sup>12</sup>Ho, J., Jain, A., & Abbeel, P. (2020). Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33, 6840–6851.



Define  $X_0 \sim \mathbf{q}_d$ ,  $X_t = X_{t-1} + (v_t^2 - v_{t-1}^2)^{1/2} \varepsilon_t$  for  $t \in \llbracket 1, T \rrbracket$  with  $\varepsilon_t \sim \mathcal{N}(0, \mathbf{I})$ .

Then,  $\text{Law}(X_t) = \mathbf{q}_t$  and  $\text{Law}(X_t | X_0 = x_0) = q_{t|0}(\cdot | x_0)$ .

Furthermore, we can write the law of  $X_{t-1}$  conditionally on  $X_t, X_0$ ,

$$q_{t-1|t,0}(x_{t-1} | x_t, x_0) = \mathcal{N} \left( x_0 + \frac{v_{t-1}^2}{v_t^2} (x_t - x_0), (v_t^2 - v_{t-1}^2) \frac{v_{t-1}^2}{v_t^2} \right).$$

---

<sup>12</sup>Ho, J., Jain, A., & Abbeel, P. (2020). Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33, 6840–6851.

## Inference distribution

For  $t \in \llbracket 2, T \rrbracket$  and  $\eta \in (0, v_{t-1})$ , set

$$q_{t-1|t,0}^{\eta}(x_{t-1}|x_t, x_0) := \mathcal{N}(x_{t-1}; \mu_{t-1}(x_0, x_t), \eta^2 \mathbf{I}_d)$$
$$\mu_{t-1}(x_0, x_t) := x_0 + (v_{t-1}^2/v_t^2 - \eta^2/v_t^2)^{1/2}(x_t - x_0).$$

The mean  $\mu_{t-1}$  is chosen to satisfy:

$$q_{t-1|0}(\mathrm{d}x_{t-1}|x_0) = \int q_{t-1|t,0}^{\eta}(\mathrm{d}x_{t-1}|x_t, x_0) q_{t|0}(\mathrm{d}x_t|x_0).$$

---

<sup>13</sup>Song, J., Meng, C., & Ermon, S. (2021). Denoising diffusion implicit models. *International Conference on Learning Representations*.  
<https://openreview.net/forum?id=St1giarCHLP>

## Full inference process

For  $\eta = \{\eta_t \in (0, v_t)\}_{t=1}^T$ , define

$$q_{1:T|0}^{\eta}(\mathrm{d}x_{1:T}|x_0) = q_{T|0}(\mathrm{d}x_T|x_0) \prod_{t=2}^T q_{t-1|t,0}^{\eta_{t-1}}(\mathrm{d}x_{t-1}|x_t, x_0),$$

and  $q_{0:T}^{\eta}(x_{0:T}) = q_{1:T|0}^{\eta}(\mathrm{d}x_{1:T}|x_0)\mathbf{q}_d(\mathrm{d}x_0)$

## Full inference process

For  $\eta = \{\eta_t \in (0, \nu_t)\}_{t=1}^T$ , define

$$q_{1:T|0}^{\eta}(\mathrm{d}x_{1:T}|x_0) = q_{T|0}(\mathrm{d}x_T|x_0) \prod_{t=2}^T q_{t-1|t,0}^{\eta_{t-1}}(\mathrm{d}x_{t-1}|x_t, x_0),$$

and  $q_{0:T}^{\eta}(x_{0:T}) = q_{1:T|0}^{\eta}(\mathrm{d}x_{1:T}|x_0) \mathbf{q}_d(\mathrm{d}x_0)$

The inference process admits the "right" marginals:

$$\mathbf{q}_t(\mathrm{d}x_t) = \int q_{0:T}^{\eta}(\mathrm{d}x_{0:T}). \quad (4)$$

# DDIM backward chain

## DDIM Recursion

$$\mathbf{q}_{t-1}(dx_{t-1}) = \int q_{t-1|t,0}^{\eta}(dx_{t-1}|x_t, x_0) \mathbf{q}_t(dx_t) q_{0|t}(dx_0|x_t).$$

where  $q_{0|t}(dx_0|x_t) = q_{t|0}(x_t|x_0) \mathbf{q}_d(dx_0) / \mathbf{q}_t(x_t)$ .

# DDIM backward chain

## DDIM Recursion

$$\mathbf{q}_{t-1}(dx_{t-1}) = \int q_{t-1|t,0}^{\eta}(dx_{t-1}|x_t, x_0) \mathbf{q}_t(dx_t) q_{0|t}(dx_0|x_t) .$$

where  $q_{0|t}(dx_0|x_t) = q_{t|0}(x_t|x_0) \mathbf{q}_d(dx_0) / \mathbf{q}_t(x_t)$ .

## DDIM Approximation

$$\hat{\mathbf{q}}_{t-1}(dx_{t-1}) = \int q_{t-1|t,0}^{\eta}(dx_{t-1}|x_t, \mu_t(x_t)) \mathbf{q}_t(dx_t) ,$$

where  $\mu_t(x_t) := \mathbb{E}_{X_0 \sim q_{0|t}(\cdot|x_t)} [X_0]$  .

## DDIM Mean approximation

Note that

$$\begin{aligned}v_t^2 \mathbf{s}_\theta(x_t, v_t) &\approx v_t^2 \mathbb{E}_{x_0 \sim q_{0|t}(\cdot|x_t)} [\nabla \log q_{t|0}(x_t|X_0)] \\ &= \mathbb{E}_{X_0 \sim q_{0|t}(\cdot|x_t)} [X_0 - x_t] = \mu_t(x_t) - x_t.\end{aligned}$$

# DDIM Mean approximation

Note that

$$\begin{aligned}v_t^2 \mathbf{s}_\theta(x_t, v_t) &\approx v_t^2 \mathbb{E}_{x_0 \sim q_{0|t}(\cdot|x_t)} [\nabla \log q_{t|0}(x_t|X_0)] \\ &= \mathbb{E}_{X_0 \sim q_{0|t}(\cdot|x_t)} [X_0 - x_t] = \mu_t(x_t) - x_t.\end{aligned}$$

## DDIM Backward Markov chain

Let  $\mu_{t,\theta}(x_t) = x_t + v_t^2 \mathbf{s}_\theta(x_t, v_t)$ ,  $\lambda_T = \mathcal{N}(0, v_T^2 \mathbf{I})$  and  $\eta = \{\eta_t \in (0, v_t)\}_{t=0}^T$ . Define

$$p_{1:T}^\theta(dx_{1:T}) := \lambda_T(dx_T) \prod_{t=1}^T p_{t-1|t}^\theta(dx_{t-1}|x_t).$$

where  $p_{t-1|t}^\theta(dx_{t-1}|x_t) = q_{t-1|t,0}^{\eta_{t-1}}(dx_{t-1}|x_t, \mu_{\theta,t}(x_t))$  for  $t > 1$  and  $p_{0|1}^\theta(x_0|x_1) = \mathcal{N}(x_0; \mu_{1,\theta}(x_1), \eta_0^2 \mathbf{I})$ .



# DDIM as variational inference

## Kullback-Leibner

$$\begin{aligned} & \text{KL}(q_{0:T}^\eta \parallel \mathbf{p}_{0:T}) \\ &= \frac{1}{2} \sum_{t=0}^{T-1} \gamma_t^2 \mathbb{E}_{X_t \sim q_{t|0}(\cdot|X_0), X_0 \sim \mathbf{q}_d} [\|\mu_{\theta,t}(X_t) - X_0\|^2] \\ &+ \frac{1}{2} v_T^{-2} \mathbb{E}_{\mathbf{q}_d} [\|X_0\|^2] + C, \end{aligned}$$

where  $\gamma_t := [v_t - (v_{t-1}^2 - \eta_{t-1}^2)^{1/2}] (\eta_{t-1} v_t)^{-1}$  for  $t > 0$ ,  $\gamma_0 = \eta_0$  and  $C$  is a constant that does not depend on  $\theta$ .

## KL calculation

$$\begin{aligned} & \text{KL}(q_{0:T}^\eta \parallel p_{0:T}^\theta) \\ &= \int \log \left( \frac{\mathbf{q}_d(x_0) q_{T|0}(x_T|x_0) \prod_{t=2}^T q_{t-1|t,0}^\eta(x_{t-1}|x_t, x_0)}{\lambda_T(x_T) \prod_{t=1}^T p_{t-1|t}^\theta(x_{t-1}|x_t)} \right) q_{0:T}^\eta(dx_{0:T}) \\ &= \sum_{t=2}^T \int \text{KL}(q_{t-1|t,0}^\eta(\cdot|x_t, x_0) \parallel p_{t-1|t}^\theta(\cdot|x_t)) \mathbf{q}_{t,0}(dx_t, dx_0) \\ &+ \int \text{KL}(\mathbf{q}_d \parallel p_{0|1}^\theta(\cdot|x_1)) \mathbf{q}_1(dx_1) + \int \text{KL}(q_{T|0}(\cdot|x_0) \parallel \lambda_T) \mathbf{q}_d(dx_0). \end{aligned}$$

## Intermediate KL

$$\begin{aligned} & \text{KL}(q_{t-1|t,0}^\eta(\cdot|x_t, x_0) \parallel p_{t-1|t}^\theta(\cdot|x_t)) \\ &= (2\eta_{t-1}^2)^{-1} \|\mu_{t-1}(x_0, x_t) - \mu_{t-1}(\mu_{t,\theta}(x_t), x_t)\|^2 \\ &= (2\eta_{t-1}^2)^{-1} \left[ 1 + (v_{t-1}^2/v_t^2 - \eta_{t-1}^2/v_t^2)^{1/2} \right]^2 \|x_0 - \mu_{t,\theta}(x_t)\|^2 \\ &= (2\eta_{t-1}^2 v_t^2)^{-1} \left[ v_t + (v_{t-1}^2 - \eta_{t-1}^2)^{1/2} \right]^2 \|x_0 - \mu_{t,\theta}(x_t)\|^2 \\ &= (1/2)\gamma_t^2 \|x_0 - \mu_{t,\theta}(x_t)\|^2. \end{aligned}$$

## Other terms

$$\text{KL}(q_{T|0}(\cdot|x_0) \parallel \lambda_T) = (2\nu_T^2)^{-1}\|x_0\|^2$$

$$\begin{aligned}\text{KL}(\mathbf{q}_d \parallel p_{0|1}^\theta(\cdot|x_1)) &= - \int \log p_{0|1}^\theta(x_0|x_1) \mathbf{q}_d(dx_0) - \mathcal{H}(\mathbf{q}_d) \\ &= (2\eta_0^2)^{-1}\|x_0 - \mu_{1,\theta}(x_1)\|^2 + (d/2) \log(2\pi\eta_0) \\ &\quad - \mathcal{H}(\mathbf{q}_d).\end{aligned}$$

Code Break! [https://github.com/gabrielvc/tutorial\\_ddim](https://github.com/gabrielvc/tutorial_ddim)

## Deep image prior

---

# Deep image prior

Let  $\tilde{x}$  be a corrupted version of an image (inpainting, denoising) and  $\mathbf{m} \in \{0, 1\}^d$  the associated mask. Ulyanov et al., 2018<sup>14</sup> proposes solving the reconstruction task by

$$\operatorname{argmin}_{\theta} \|\mu_{\theta}(z) \odot \mathbf{m} - \tilde{x} \odot \mathbf{m}\|^2,$$

where  $z$  is a fixed seed.

---

<sup>14</sup>Ulyanov, D., Vedaldi, A., & Lempitsky, V. (2018). Deep image prior. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.

# Deep image prior



Figure: Inpainting example from Ulyanov et al., 2018<sup>15</sup>

---

<sup>15</sup>Ulyanov, D., Vedaldi, A., & Lempitsky, V. (2018). Deep image prior. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.

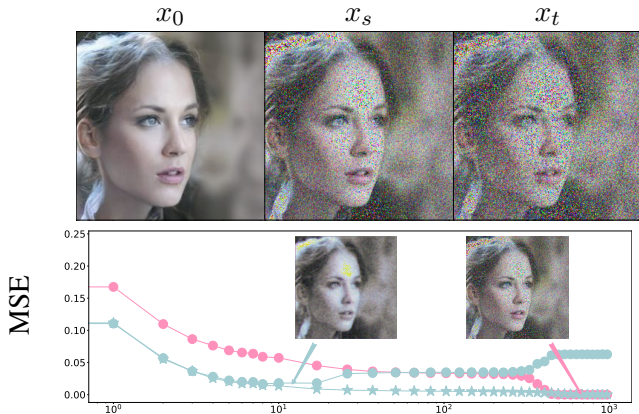


# Denoising tasks

Suppose  $\mu_{t,\theta}(x_t)$  is a UNet. Consider the following losses:

$$L_{t|s}(\theta) := \|\mu_{t,\theta}(x_t) - x_s\|^2 \quad \text{and} \quad L_{t|0}(\theta) := \|\mu_{t,\theta}(x_t) - x_0\|^2.$$

We train  $\mu_{t,\theta}(x_t)$  to minimize  $L_{t|s}(\theta)$ .



## Conclusion

---

## Interesting papers

- Going further on Diffusion models: Karras et al., 2022<sup>16</sup>, Y. Song, Sohl-Dickstein, et al., 2021<sup>17</sup>, Y. Song, Durkan, et al., 2021.<sup>18</sup>

---

<sup>16</sup>Karras, T., Aittala, M., Aila, T., & Laine, S. (2022). Elucidating the design space of diffusion-based generative models. *Proc. NeurIPS*.

<sup>17</sup>Song, Y., Sohl-Dickstein, J., Kingma, D. P., Kumar, A., Ermon, S., & Poole, B. (2021). Score-based generative modeling through stochastic differential equations. *International Conference on Learning Representations*.  
<https://openreview.net/forum?id=PXTIG12RRHS>

<sup>18</sup>Song, Y., Durkan, C., Murray, I., & Ermon, S. (2021). Maximum likelihood training of score-based diffusion models. *Advances in Neural Information Processing Systems*, 34, 1415–1428.

## Interesting papers

- Diffusion models as priors for inverse problems: Chung et al., 2023<sup>19</sup>, Cardoso et al., 2023<sup>20</sup>, Wu et al., 2023<sup>21</sup>.

---

<sup>19</sup>Chung, H., Kim, J., Mccann, M. T., Klasky, M. L., & Ye, J. C. (2023). Diffusion posterior sampling for general noisy inverse problems. *The Eleventh International Conference on Learning Representations*.

<https://openreview.net/forum?id=OnD9zGAGT0k>

<sup>20</sup>Cardoso, G., Idrissi, Y. J. E., Corff, S. L., & Moulines, E. (2023). Monte carlo guided diffusion for bayesian linear inverse problems.

<sup>21</sup>Wu, L., Trippe, B. L., Naeseth, C. A., Blei, D. M., & Cunningham, J. P. (2023). Practical and asymptotically exact conditional sampling in diffusion models.

## Interesting papers

- Developpements on diffusion models: Rombach et al., 2022<sup>22</sup>, Y. Song et al., 2023<sup>23</sup>.

---

<sup>22</sup>Rombach, R., Blattmann, A., Lorenz, D., Esser, P., & Ommer, B. (2022). High-resolution image synthesis with latent diffusion models. *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 10684–10695.

<sup>23</sup>Song, Y., Dhariwal, P., Chen, M., & Sutskever, I. (2023). Consistency models.

## **Temporary page!**

$\text{\LaTeX}$  was unable to guess the total number of pages correctly. As there was some unprocessed data that should have been added to the final page this extra page has been added to receive it.

If you rerun the document (without altering it) this surplus page will go away, because  $\text{\LaTeX}$  now knows how many pages to expect for the document.