

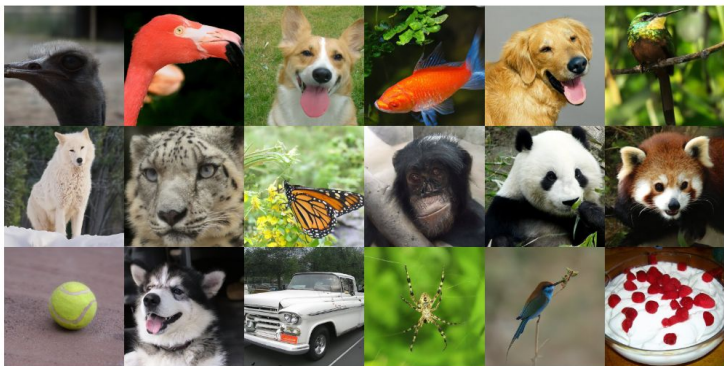


Deep Learning Architectures

Diffusion model



INSTITUT DU
DÉVELOPPEMENT ET DES
RESSOURCES EN
INFORMATIQUE
SCIENTIFIQUE



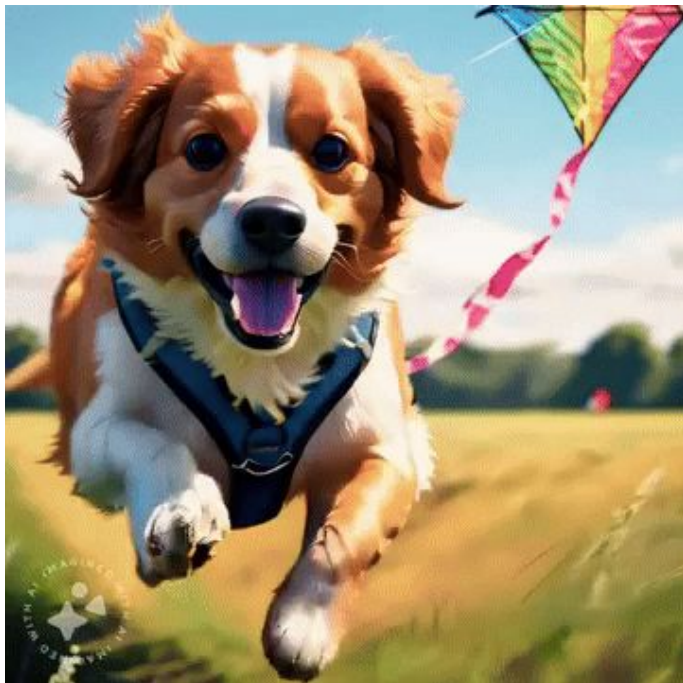
Dhariwal & Nichol, 2021



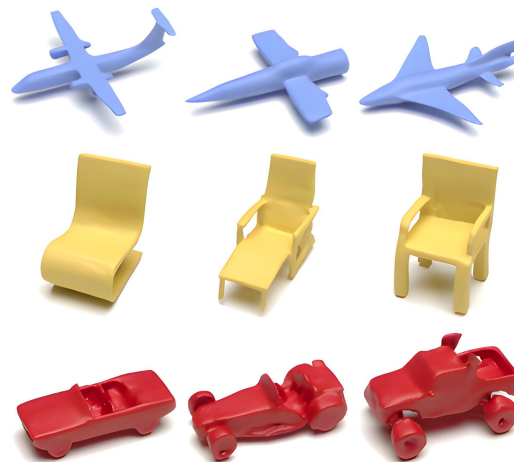
Source : <https://github.com/bentoml/stable-diffusion-bentoml>



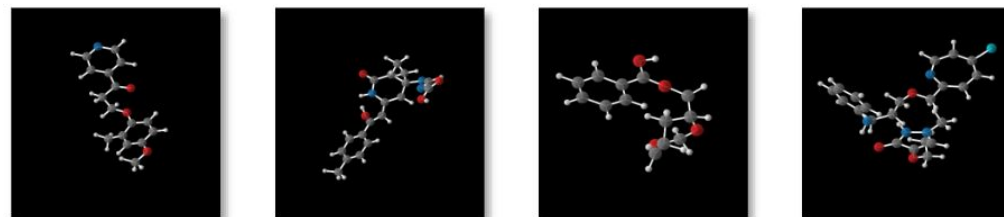
Source : DALL-E 3



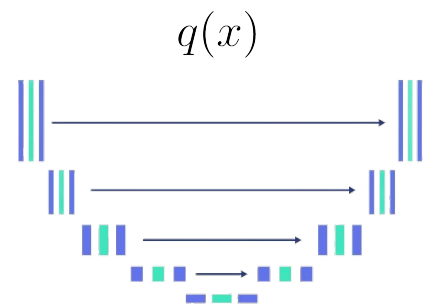
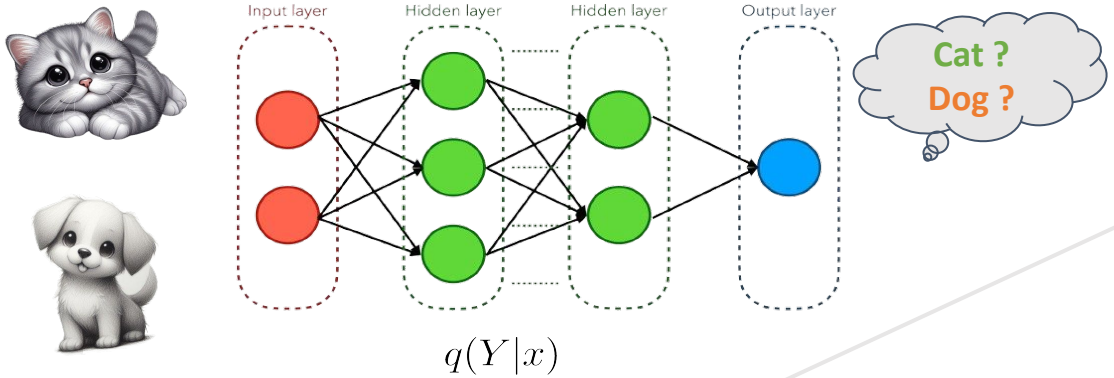
Source : <https://ai.meta.com/blog/emu-text-to-video-generation-image-editing-research/>

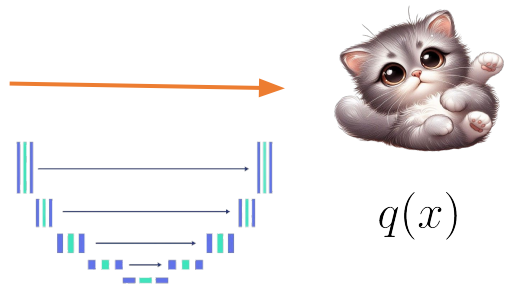
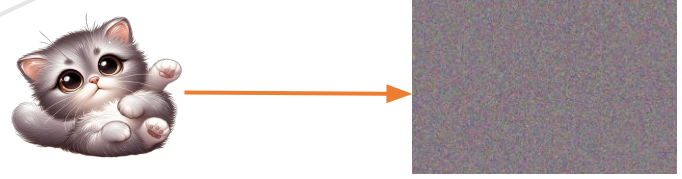
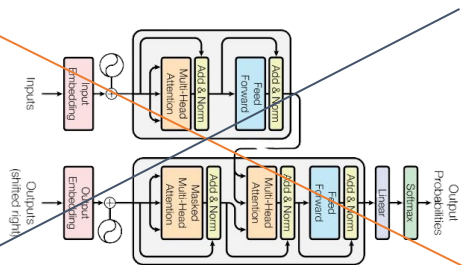
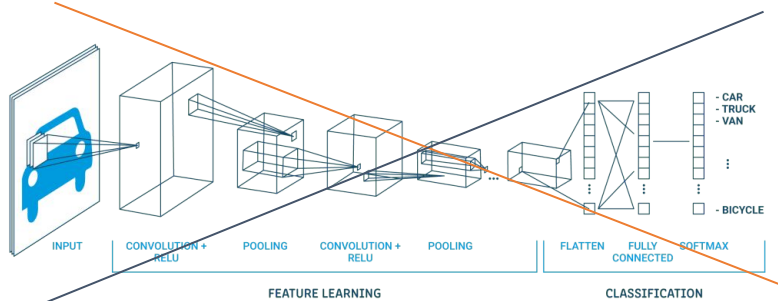


Source : <https://arxiv.org/pdf/2210.06978.pdf>



Source : <https://arxiv.org/pdf/2305.01140.pdf>

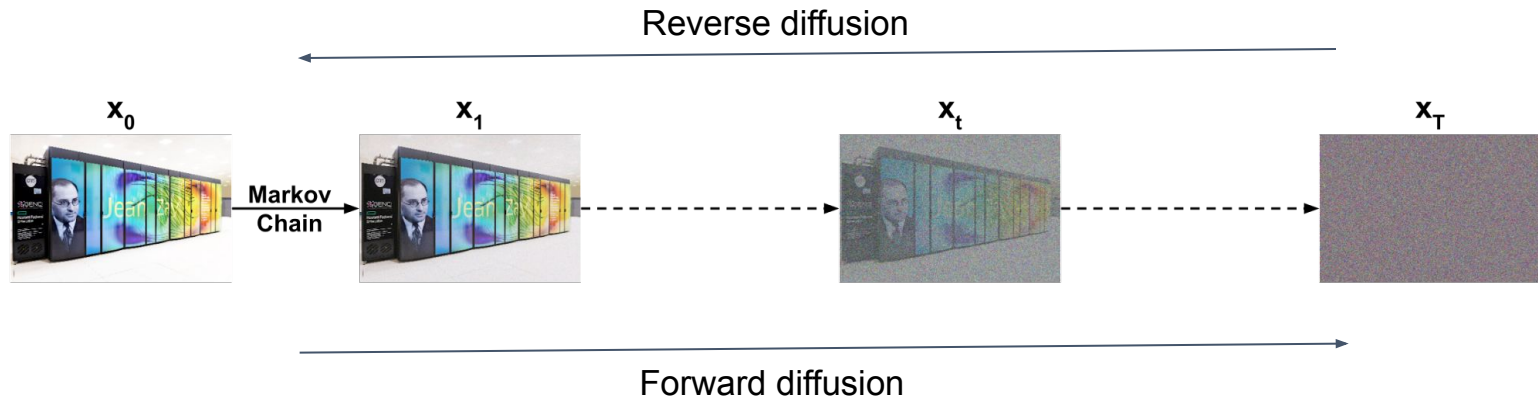




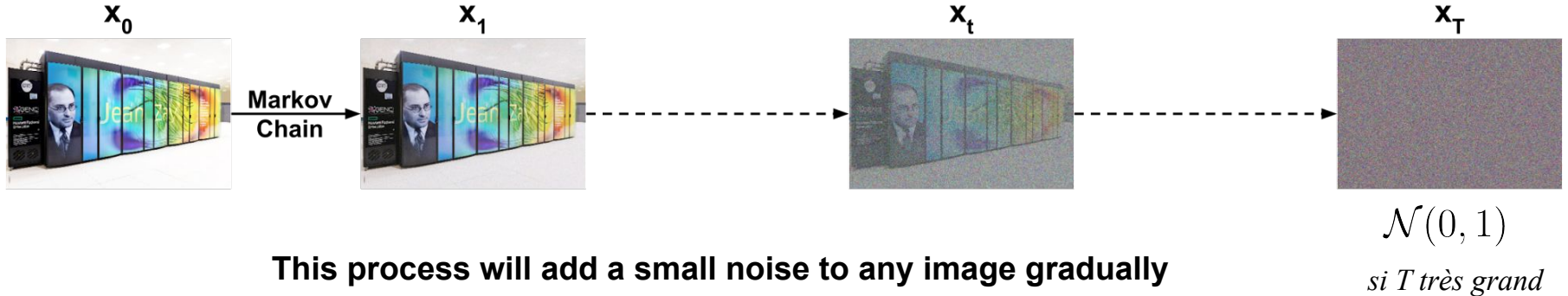
CNN/Transformer/GNN vs Diffusion Model

Denosing diffusion probabilistic models (DDPM) consist of two processes:

- Forward diffusion process that gradually adds noise to input
- Reverse denoising process that learns to generate data by denoising



Forward Diffusion Process



$0 \leq t \leq T$; T is a hyperparameter

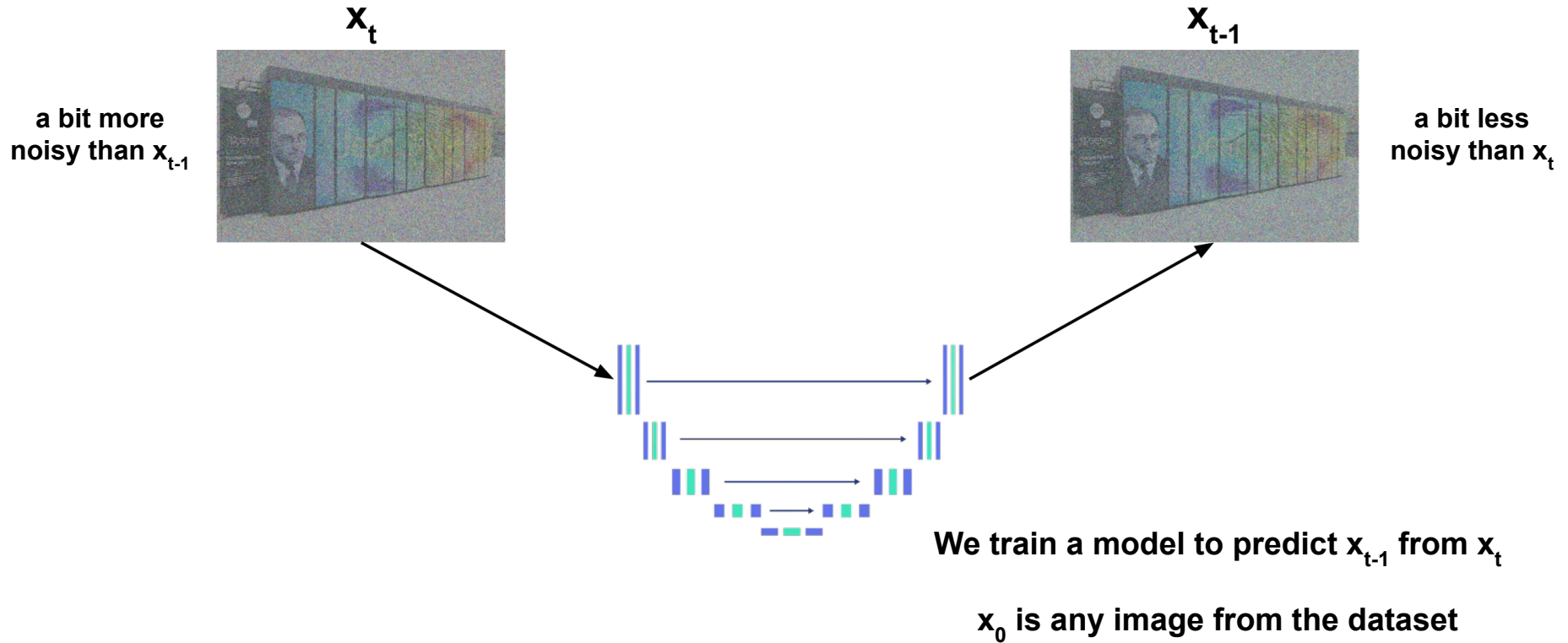
Forward Diffusion Process



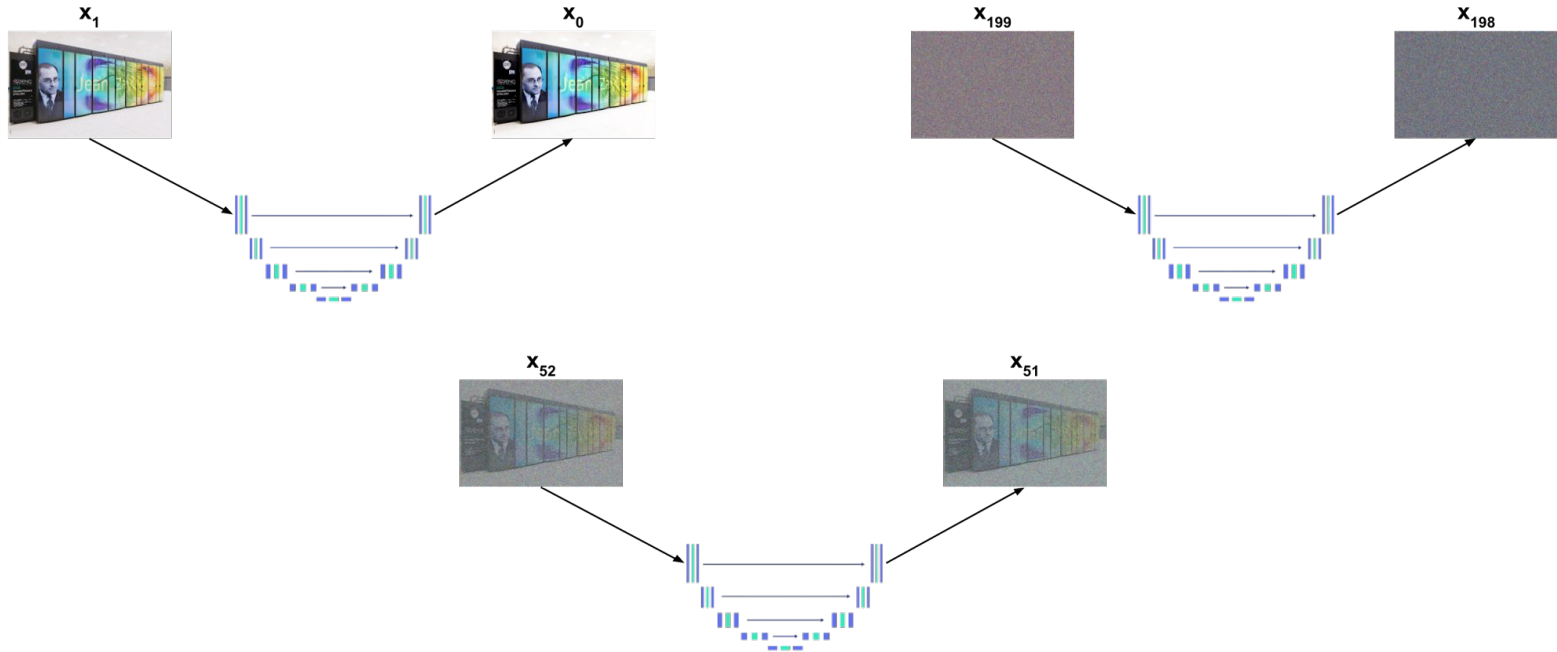
Denoising Diffusion Probabilistic
Models (DDPM)

Here we choose $T=1000$, but it can be different values (it's an hyperparameter)

Reverse Diffusion Process

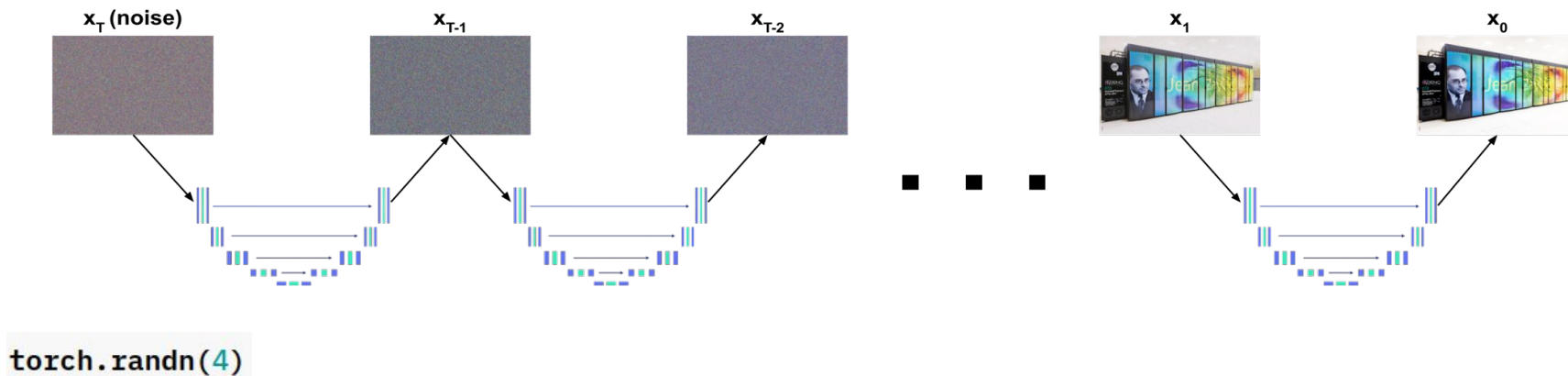


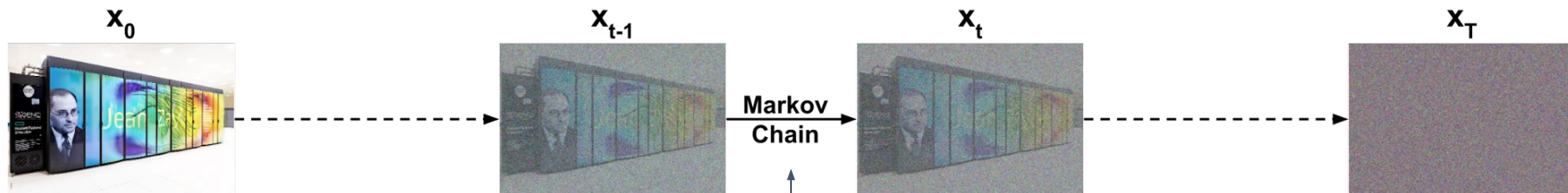
Reverse Diffusion Process



The same model must predict every x_{t-1} from x_t $Model(x_t, t)$

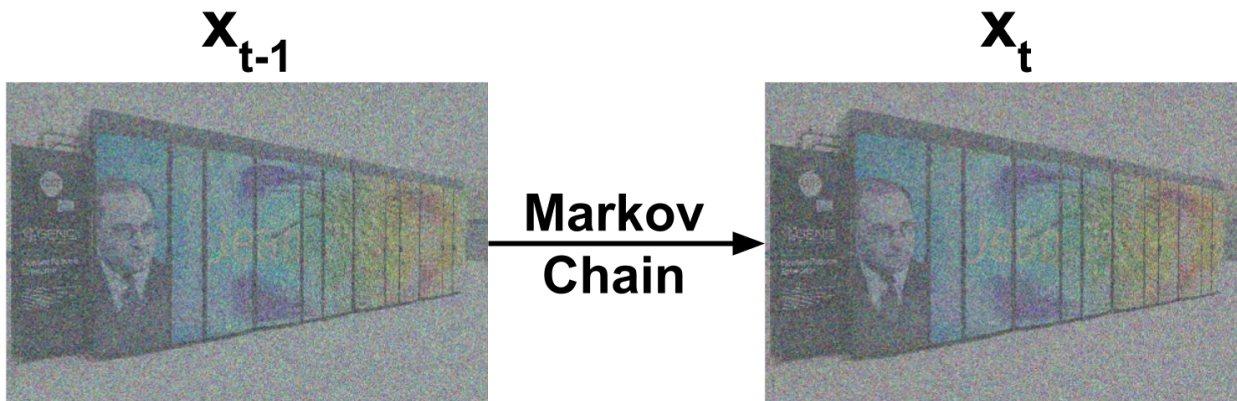
After the training, the model will generate images from Gaussian noise by following a sampling process :





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

(β_t) is a hyperparameter



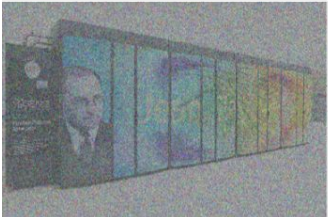
$$\mathcal{N}(x_t; \sqrt{1 - \beta_t}x_{t-1}, \beta_t I)$$

$$x_t = \sqrt{1 - \beta_t}x_{t-1} + \sqrt{\beta_t}z_{t-1}$$

where $z_{t-1} \sim \mathcal{N}(0, I)$

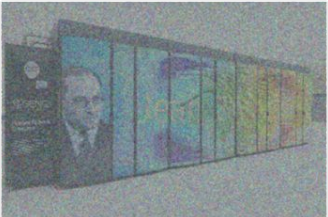
```
torch.randn(4)
```

x_t



$=$


$\sqrt{1 - \beta_t}$

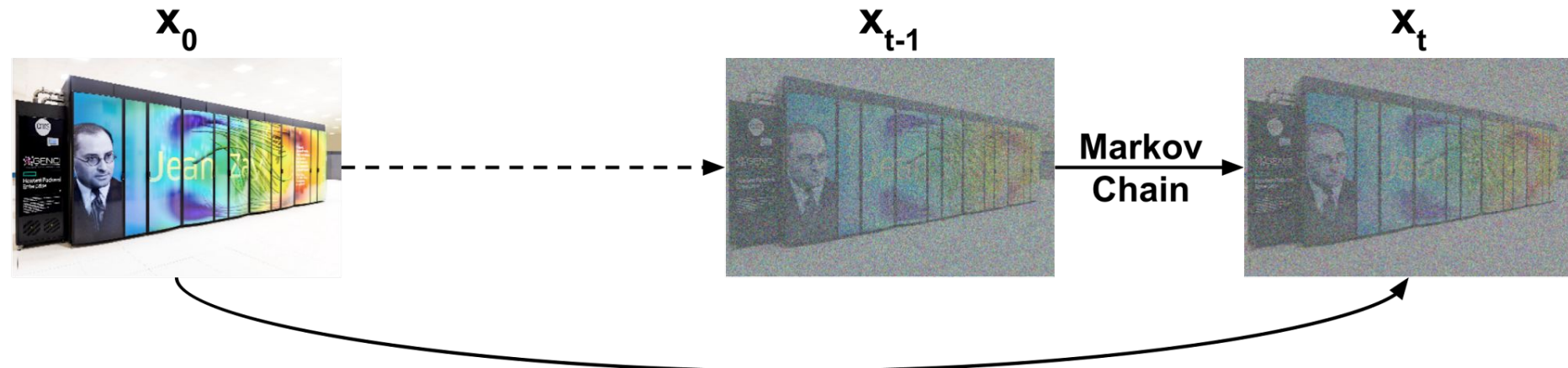


$+$

$\sqrt{\beta_t}$

Gaussian noise





$$q(x_t|x_0) = \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)I)$$

$$\text{where } \bar{\alpha}_t = \prod_{i=1}^T (1 - \beta_i)$$



$$x_t = \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} z$$

So we can sample a noised image at any time step directly from original image

β_t (the noise schedule) such that $q(\mathbf{x}_T|\mathbf{x}_0) \approx \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I})$
 $\beta_t \in (0, 1)$, β_t following a schedule $\beta_1 < \beta_2 < \dots < \beta_T$

Linear
scheduling

$$(\beta_1, \dots, \beta_T) = (\beta_1 + (t-1) * \frac{\beta_T - \beta_1}{T-1})_{t \in \{1 \dots T\}}$$

0.02

0.0001

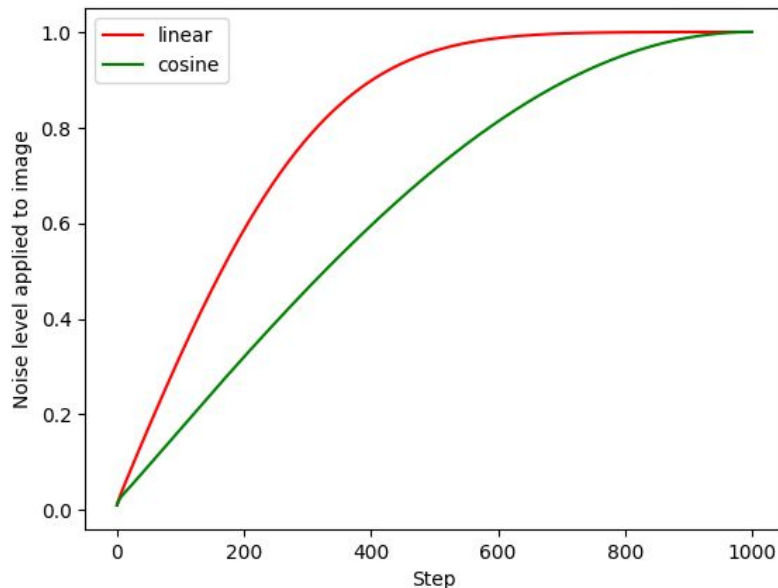
1000

Cosinus
scheduling

$$\bar{\alpha}_t = \frac{f(t)}{f(0)}, \quad f(t) = \cos\left(\frac{t/T + s}{1+s} \cdot \frac{\pi}{2}\right)^2$$

$$\text{where } \bar{\alpha}_t = \prod_{i=1}^T (1 - \beta_i)$$

Linear and Cosine schedules



$$q(x_0)$$

Données

$$q(x_t|x_{t-1})$$

*Données bruitées
par un petit bruit
gaussien*



$$q(x_t|x_0)$$

$$q(x_{t-1}|x_t, x_0)$$



$$q(x_{t-1}|x_t) \quad ???$$

$$q(x_{t-1}|x_t, x_0) = \mathcal{N}(x_{t-1}; \tilde{\mu}_t(x_t, x_0), \tilde{\beta}_t I)$$

$$\tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0) = \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1 - \bar{\alpha}_t}\mathbf{x}_0 + \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t}\mathbf{x}_t \quad \tilde{\beta}_t = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t}\beta_t$$

$$q(x_t|x_0)$$

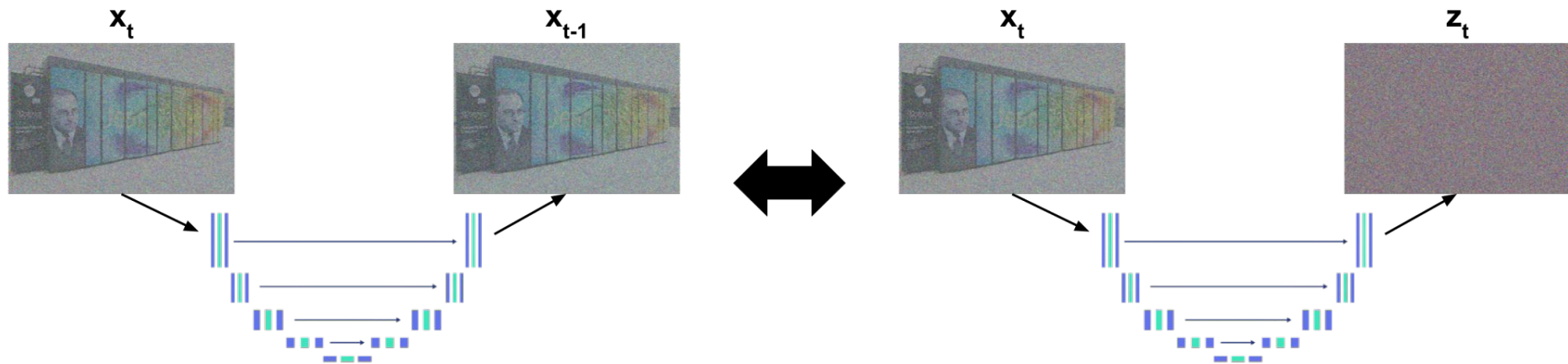
$$x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}z_t$$



$$q(x_{t-1}|x_t, z_t)$$

$$\tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0) = \frac{1}{\sqrt{\alpha_t}}\left(\mathbf{x}_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}}\mathbf{z}_t\right)$$

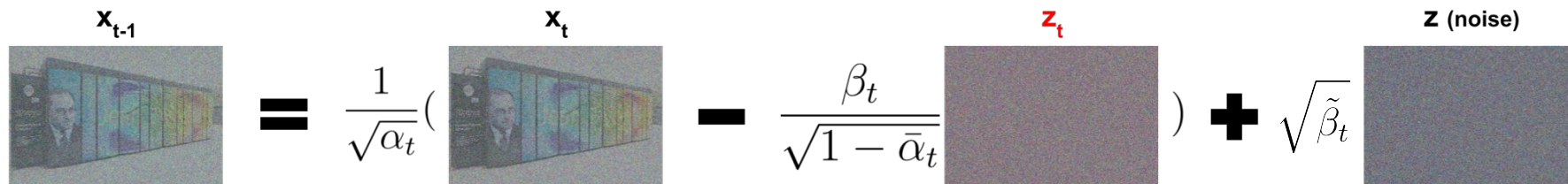
\mathbf{z}_t ???



We can predict x_{t-1} by predicting z_t

$$\tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0) = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \mathbf{z}_t \right) \approx \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \mathbf{z}_\theta(\mathbf{x}_t, \mathbf{t}) \right)$$

Our network



$$\frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \mathbf{z}_\theta(\mathbf{x}_t, \mathbf{t}) \right) \quad x_t = \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} z_t$$

$$L_{\text{VLB}} = \mathbb{E}_{q(\mathbf{x}_{0:T})} \left[\log \frac{q(\mathbf{x}_{1:T} | \mathbf{x}_0)}{p_\theta(\mathbf{x}_{0:T})} \right] \geq -\mathbb{E}_{q(\mathbf{x}_0)} \log p_\theta(\mathbf{x}_0)$$

$$L_{\text{VLB}} = L_T + L_{T-1} + \dots + L_0$$

$$L_t = \mathbb{E}_{\mathbf{x}_0, \epsilon} \left[\frac{(1 - \alpha_t)^2}{2\alpha_t(1 - \bar{\alpha}_t) \|\Sigma_\theta\|_2^2} \|\epsilon_t - \epsilon_\theta(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon_t, t)\|^2 \right]$$

$$L_t^{\text{simple}} = \mathbb{E}_{t \sim [1, T], \mathbf{x}_0, \epsilon_t} \left[\|\epsilon_t - \epsilon_\theta(\mathbf{x}_t, t)\|^2 \right]$$

Algorithm 1 Training

1: **repeat**

2: $\mathbf{x}_0 \sim q(\mathbf{x}_0)$

3: $t \sim \text{Uniform}(\{1, \dots, T\})$

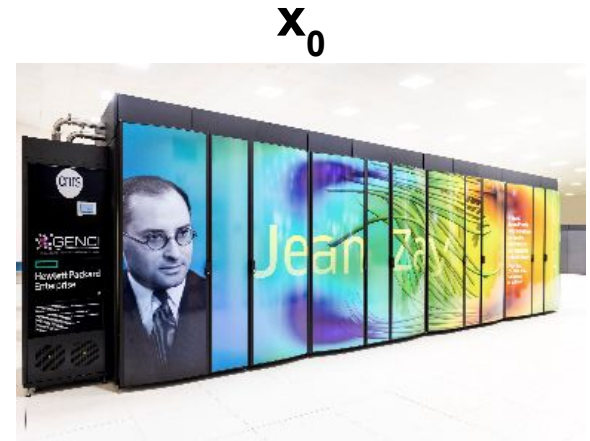
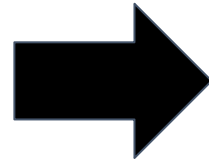
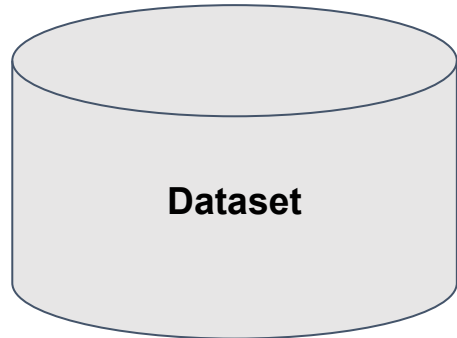
4: $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$

5: Take gradient descent step on

$$\nabla_{\theta} \left\| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta}(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, t) \right\|^2$$

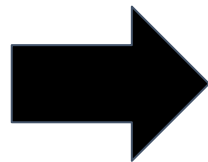
6: **until** converged

<https://arxiv.org/abs/2006.11239>



**Uniform
distribution**

Between 1 and T



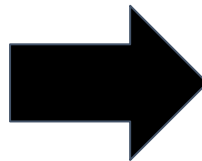
t = 50

x_0

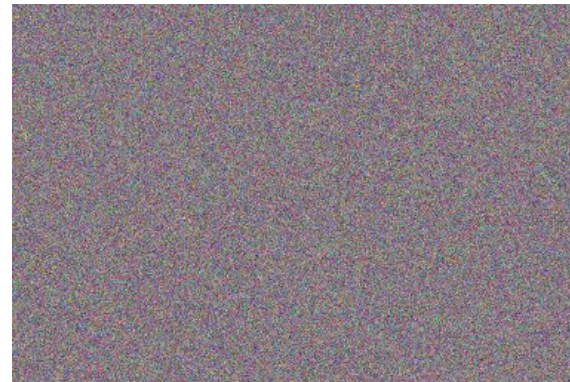


**Gaussian
distribution**

Same shape than x_0



$\mathbf{z}_t (= \epsilon)$



x_0

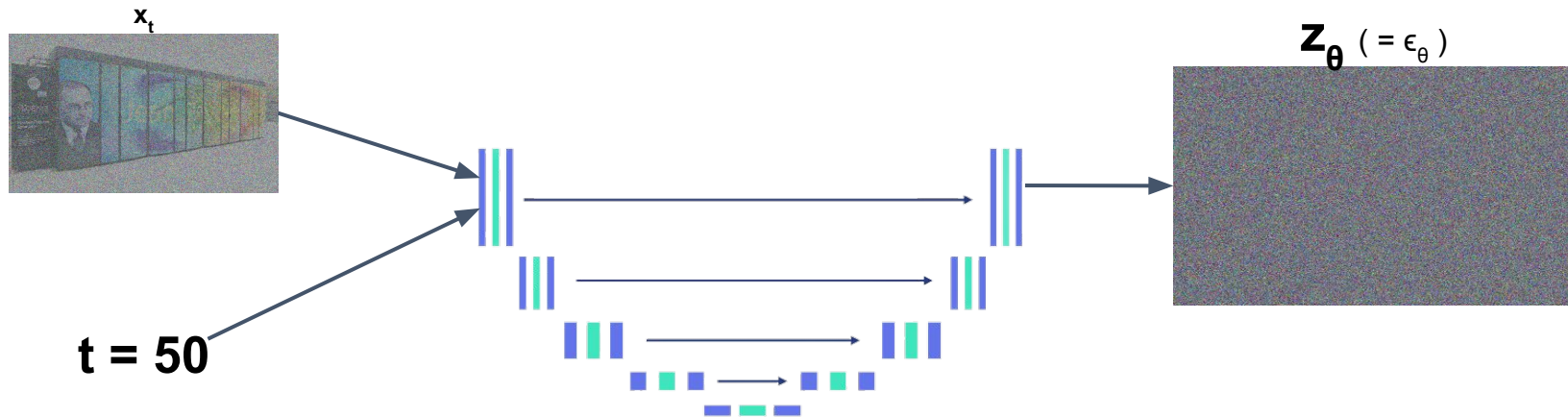


t = 50

$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \mathbf{z}_t$$

$$\mathbf{x}_0 \quad \mathbf{z}_t$$

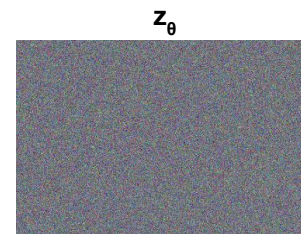
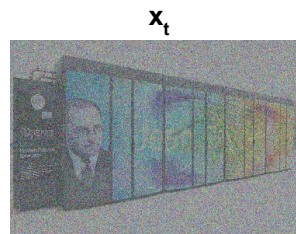
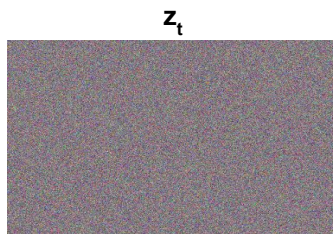
t = 50

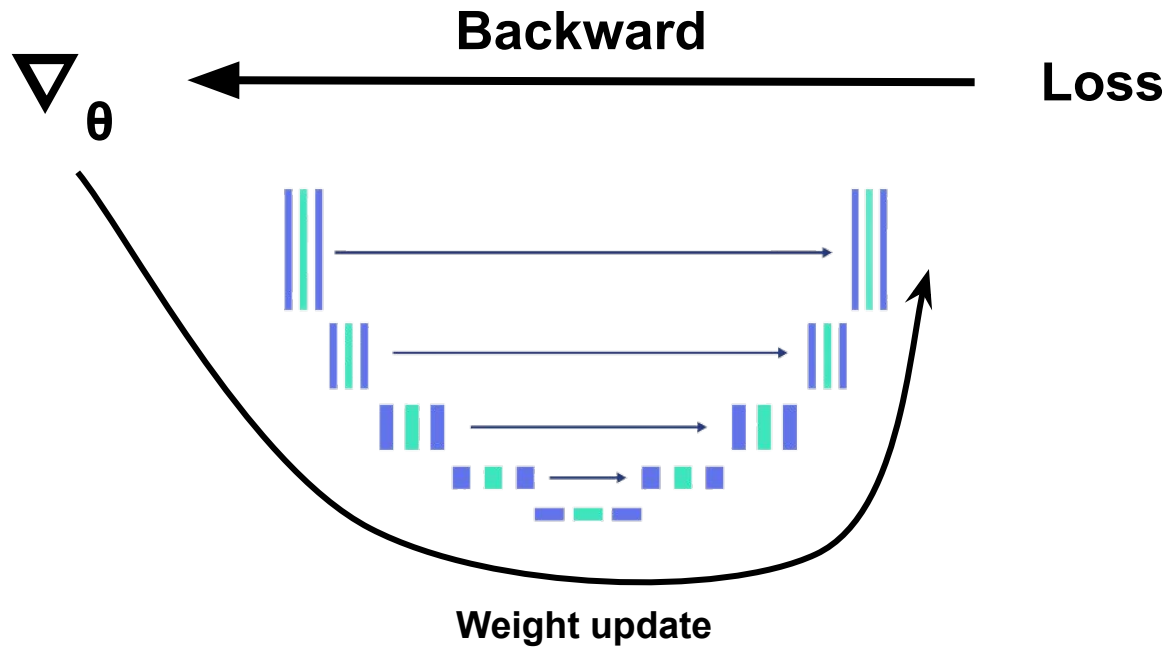


$$\text{Loss} = \left\| z_t - z_\theta \right\|^2$$



$t = 50$





Algorithm 2 Sampling

1: $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$

2: **for** $t = T, \dots, 1$ **do**

3: $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ if $t > 1$, else $\mathbf{z} = \mathbf{0}$

4: $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$

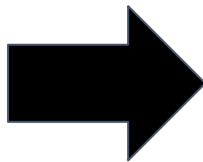
5: **end for**

6: **return** \mathbf{x}_0

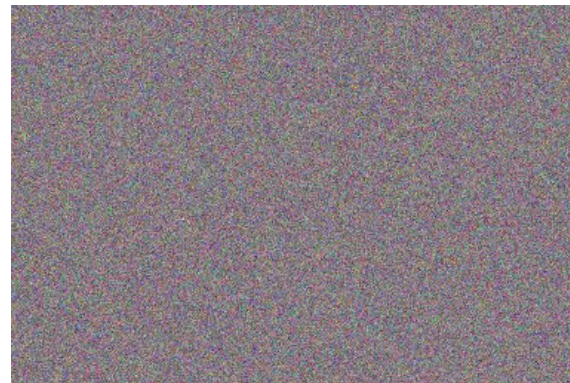
<https://arxiv.org/abs/2006.11239>

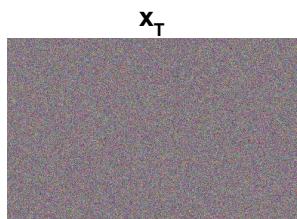
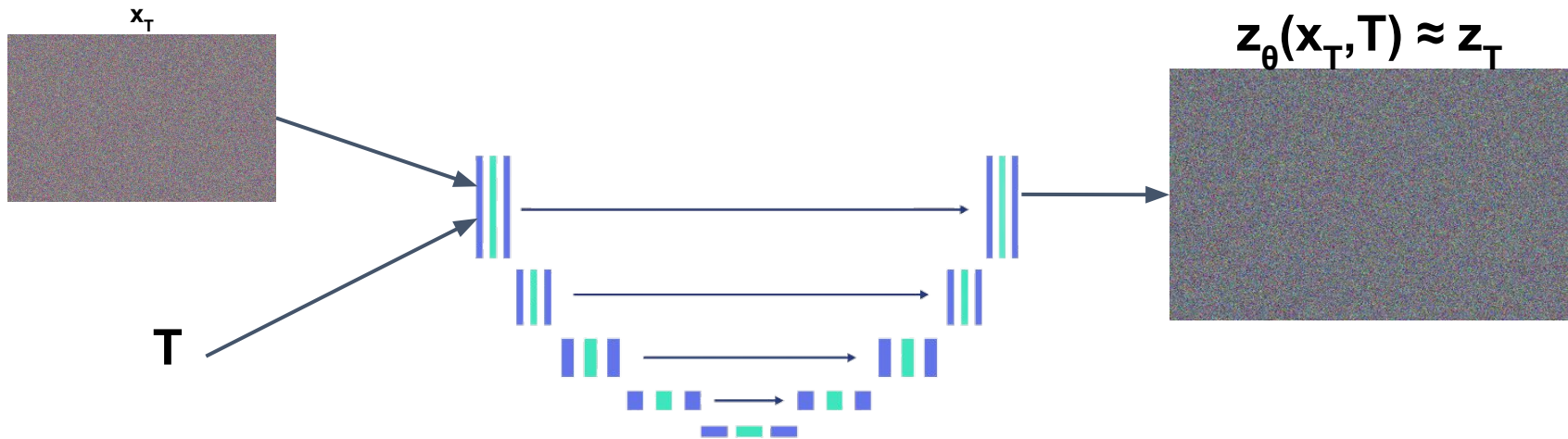
**Gaussian
distribution**

Same shape than training
dataset images



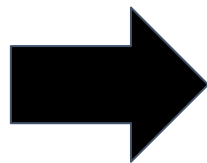
x_T



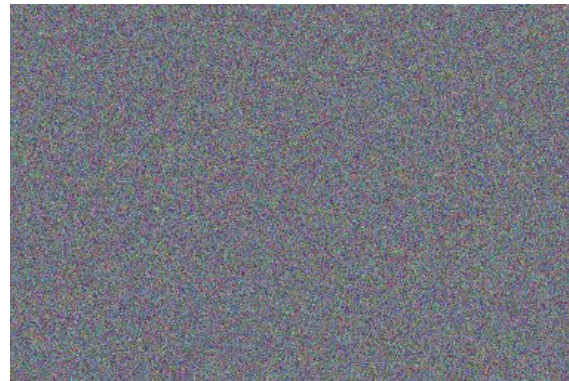


**Gaussian
distribution**

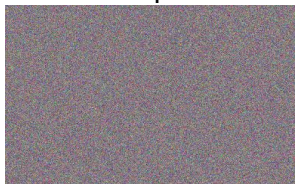
Same shape than training
dataset images



z (noise)



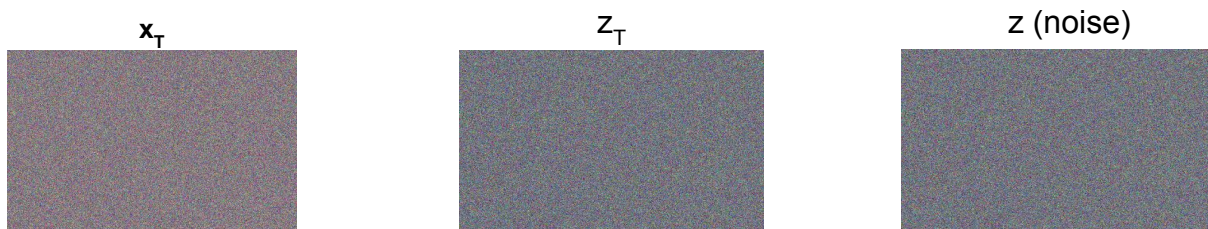
x_T



z_T

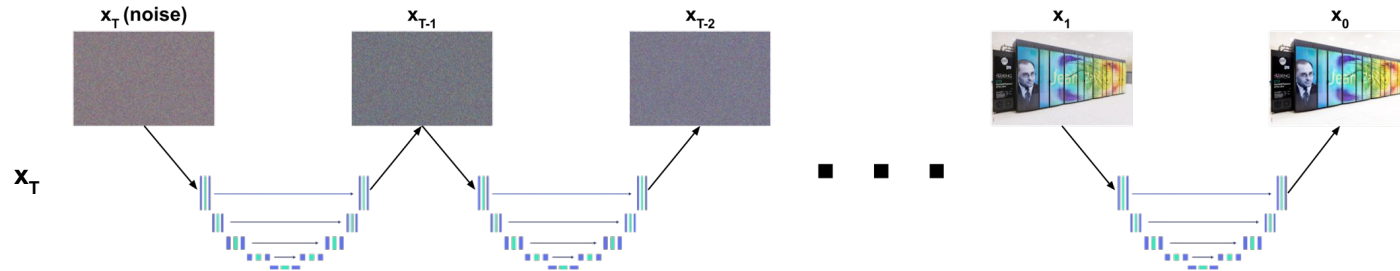


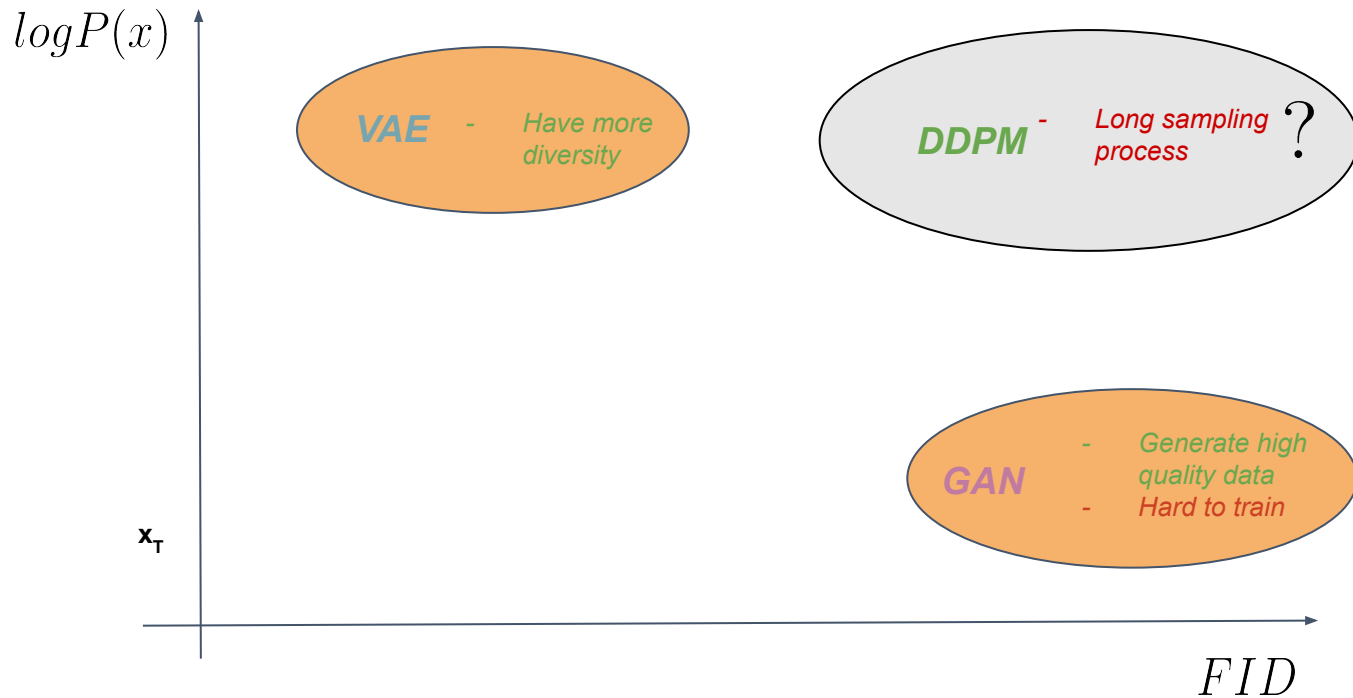
$$\mathbf{x}_{T-1} \approx \frac{1}{\sqrt{\alpha_T}} \left(\mathbf{x}_T - \frac{\beta_T}{\sqrt{1 - \bar{\alpha}_T}} \mathbf{z}_\theta(\mathbf{x}_T, T) \right) + \tilde{\beta}_T \mathbf{z} \text{ (noise)}$$



and repeat !

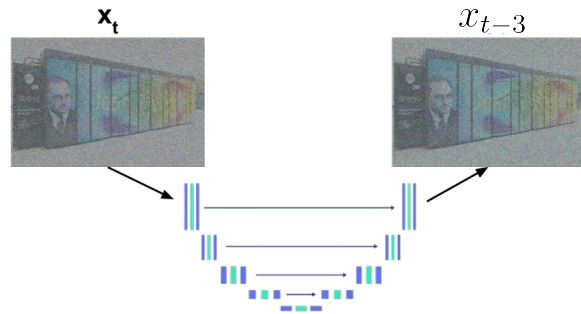
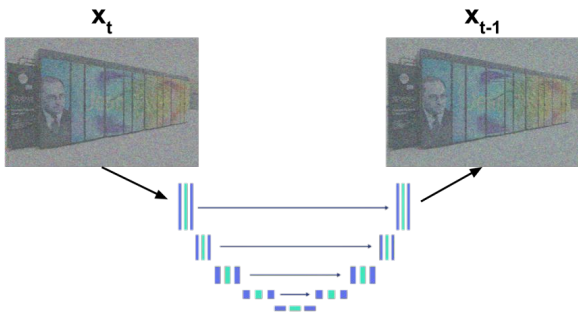
Don't generate x_T , replace it by x_{T-1} and T by T-1... and do it again T time.





Limitations :

“For example, it takes around 20 hours to sample 50k images of size 32 x 32 from a DDPM, but less than a minute to do so from a GAN on a Nvidia 2080 Ti GPU.” (DDIM, 2021)



$$x_t = \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} z_t$$

$$x_{t-3} = \sqrt{\bar{\alpha}_{t-3}} x_0 + \sqrt{1 - \bar{\alpha}_{t-3}} z$$



$$\frac{x_t - \sqrt{1 - \bar{\alpha}_t} z_t}{\sqrt{\bar{\alpha}_t}} = x_0$$

$$x_{t-3} = \sqrt{\bar{\alpha}_{t-3}} * \frac{x_t - \sqrt{1 - \bar{\alpha}_t} z_t}{\sqrt{\bar{\alpha}_t}} + \sqrt{1 - \bar{\alpha}_{t-3}} z$$

Generalisation to a bigger class of inverse process (non-Markovian)

$$q_{\sigma}(x_{t-1}|x_t, x_0) = \mathcal{N}(\sqrt{\alpha_{t-1}^-}x_0 + \sqrt{1 - \alpha_{t-1}^- - \sigma_t^2} \cdot \frac{x_t - \sqrt{\bar{\alpha}_t}x_0}{\sqrt{1 - \bar{\alpha}_t}}, \sigma_t^2 I)$$

$$\sigma_t^2 = \tilde{\beta}_t \quad \Rightarrow \quad DDPM$$

$$q(x_t|x_0) = \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)I)$$



Important :

Same network and training as a DDPM

$$\sigma_t^2 = \tilde{\beta}_t \Rightarrow DDPM$$

$$\sigma_t^2 = \eta \cdot \tilde{\beta}_t \quad \eta \in [0, 1]$$

$$\eta = 0 \Rightarrow DDIM$$

S	CIFAR10 (32 × 32)					CelebA (64 × 64)					
	10	20	50	100	1000	10	20	50	100	1000	
η	0.0	13.36	6.84	4.67	4.16	4.04	17.33	13.73	9.17	6.53	3.51
	0.2	14.04	7.11	4.77	4.25	4.09	17.66	14.11	9.51	6.79	3.64
	0.5	16.66	8.35	5.25	4.46	4.29	19.86	16.06	11.01	8.09	4.28
	1.0	41.07	18.36	8.01	5.78	4.73	33.12	26.03	18.48	13.93	5.98
$\hat{\sigma}$		367.43	133.37	32.72	9.99	3.17	299.71	183.83	71.71	45.20	3.26

FID



DDIM 2

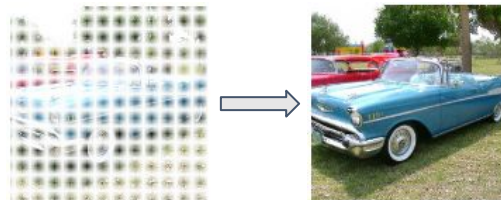
Diffusion models can solve a variety of tasks. We already know about image generation, as well as conditional image generation (for instance with a short paragraph describing the picture)

Other tasks:

→ Inpainting



→ Super-resolution



→ Outpainting



We can solve many of these tasks through the usage of a mask

→ Wide mask



→ Thin mask



→ Outer mask for expanding the image



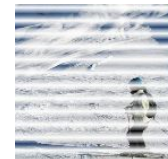
→ Right side mask for halving the image



→ Every second pixel in both directions for super-resolution

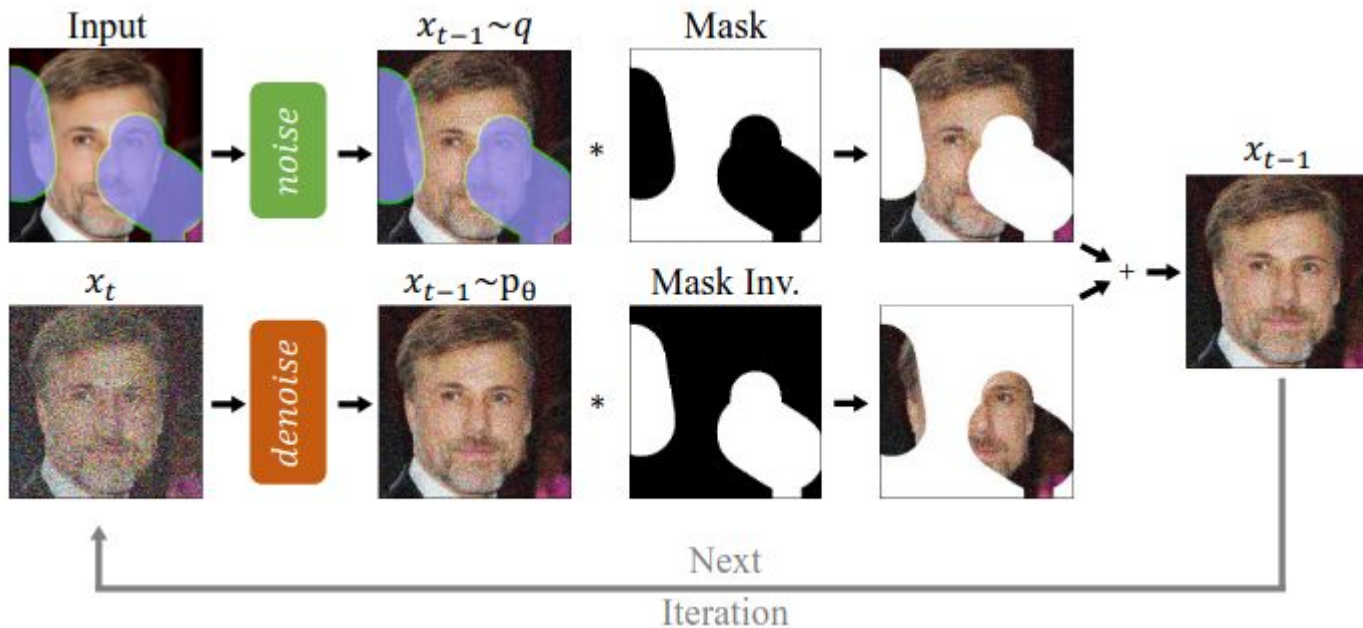


→ Every second row of pixels for alternating lines

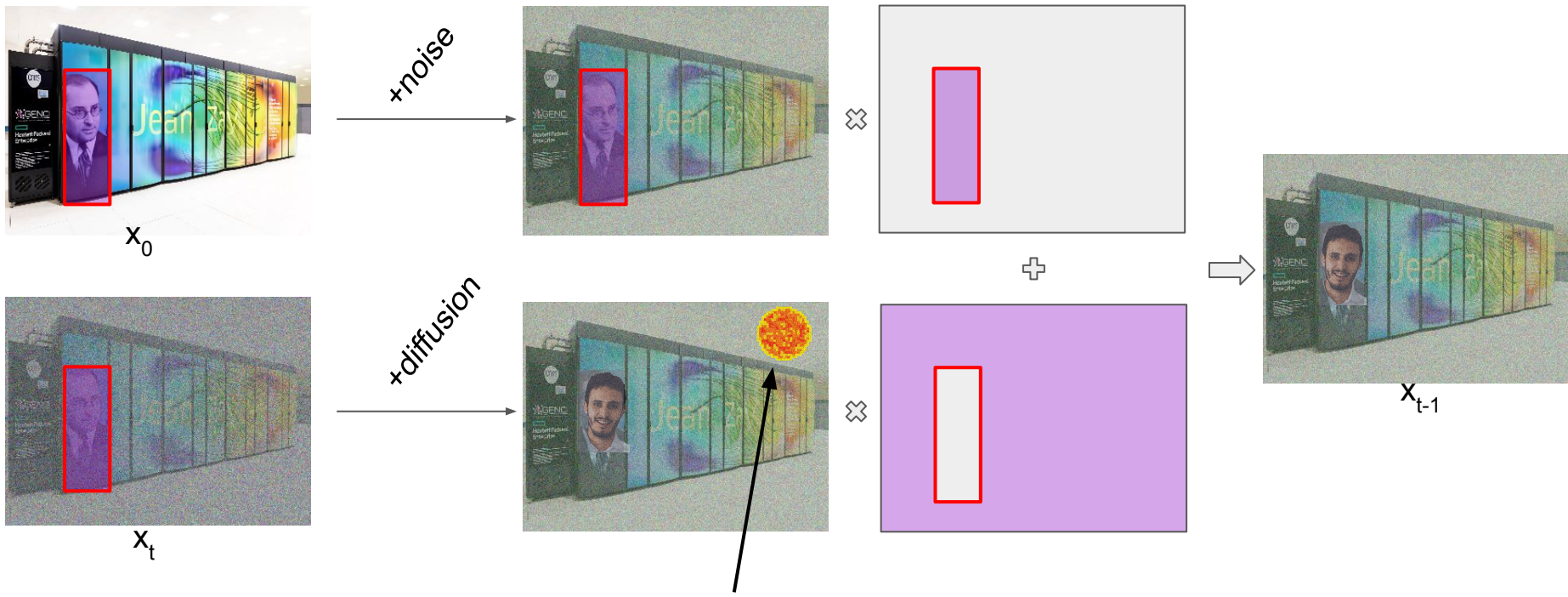


Source Lugmayr, Andreas, et al. "Repaint: Inpainting using denoising diffusion probabilistic models." *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2022.

Step t :



Source Lugmayr, Andreas, et al. "Repaint: Inpainting using denoising diffusion probabilistic models." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2022.



New artifacts added (in this coarse example, our diffusion model drew a sun), so we force the known background again!

Papers:

- Deep Unsupervised Learning using Nonequilibrium Thermodynamics (DPM) (<https://arxiv.org/abs/1503.03585>)
- Denoising Diffusion Probabilistic Models (DDPM) (<https://arxiv.org/abs/2006.11239>)
- Improved Denoising Diffusion Probabilistic Models (IDDP) (<https://arxiv.org/abs/2102.09672>)
- Denoising Diffusion Implicit Models (DDIM) (<https://arxiv.org/abs/2010.02502>)
- Diffusion Models Beat GANs on Image Synthesis (<https://arxiv.org/abs/2105.05233>)
- High-Resolution Image Synthesis with Latent Diffusion Models (LDM) (<https://arxiv.org/abs/2112.10752>)
- Repaint: Inpainting using denoising diffusion probabilistic models (<https://arxiv.org/pdf/2201.09865>)
- Diffusion Models in Vision: A Survey (<https://arxiv.org/abs/2209.04747>)
- Diffusion Models: A Comprehensive Survey of Methods and Applications (<https://arxiv.org/abs/2209.00796>)

Other resources:

- Lilian Weng's article (<https://lilianweng.github.io/posts/2021-07-11-diffusion-models>)
- Yang Song's article (<https://yang-song.net/blog/2021/score>)
- Outlier video (<https://www.youtube.com/watch?v=HoKDTa5jHvg>)