# Deep Learning Architectures

## Diffusion model

*Dhariwal & Nichol, 2021*



*Source : https://github.com/bentoml/stable-diffusion-bentoml*



*Source : Dall-E 3*

# DPM Landscape

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

# Generation of videos, of molecules, of …

Source : Dall-E 3

$$p(x)$$ = **Probabilité que le point x appartient à notre base de données**

**Distribution de données P(x)**

$p(y|x)$

$\sim$

$q(y|x)$

$y$

*Classification*

*Generation*

$p(x)$

$\sim$

$q(x)$

**Generation vs Classification**

Classification
with 3 classes

Deep neural network

a discrete
probability
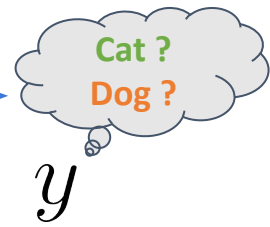
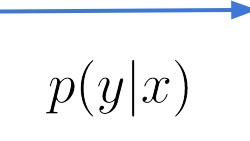$= \begin{bmatrix} 0.1 \\ 0.1 \\ 0.8 \end{bmatrix}$ *: chien*
*: chat*
*: humain*

$= \begin{bmatrix} 0.0 \\ 0.0 \\ 1.0 \end{bmatrix}$

$q(y|x)$

???

$= \mathcal{N}(output, \sigma^2)$

$= output + z * \sigma, \text{ where } z \sim \mathcal{N}(0,1)$

*Generation of one continuous value*

**Output of models**

$\mathcal{N}(0,1)$

$p(x)$

$q(x|\mathcal{N}(0,1))$

$\sim$

$$q(x|\mathcal{N}(0,1)) * \mathcal{N}(0,1) = q(x)$$

**CNN/Transformer/GNN vs Diffusion Model**

**Denoising 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

Reverse diffusion

$x_0$             $x_1$             $x_t$             $x_T$

**Markov Chain**

Forward diffusion

# DDPM Principle

# Forward Diffusion Process



| t=0 | t=5 | t=10 | t=50 | t=100 | t=T=1000 |

Denoising Diffusion Probabilistic Models (DDPM)

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

## DDPM Principle 3

# Reverse Diffusion Process

**x$_t$**

**x$_{t-1}$**

**a bit more noisy than x$_{t-1}$**

**a bit less noisy than x$_t$**

**We train a model to predict x$_{t-1}$ from x$_t$**

**x$_0$ is any image from the dataset**

**DDPM Principle 4**

# Reverse Diffusion Process



**Important :**

The same model must predict every $x_{t-1}$ from $x_t$

**DDPM Principle 5**

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



$x_T$ (noise)  $x_{T-1}$  $x_{T-2}$  $x_1$  $x_0$

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

$$\mathrm{PE}(t) = \left[\sin\left(\frac{t}{10000^{\frac{2i}{d}}}\right), \cos\left(\frac{t}{10000^{\frac{2i}{d}}}\right)\right]_{i=0}^{d/2}$$

**Unet: a good denoiser + t-embedding**

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

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

`torch.randn(4)`

**Forward diffusion 2 - Sampling a gaussian distribution**

$$\mathbf{x_t} = \sqrt{1-\beta_t} \quad \mathbf{x_{t-1}} \quad + \quad \sqrt{\beta_t} \quad \text{Gaussian noise}$$

**Forward diffusion 3**

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

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

**Forward diffusion 4**

$\beta_t$ (the noise schedule) such that $q(\mathbf{x}_T|\mathbf{x}_0) \approx \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I})$

$\beta_t \epsilon (0,1),\ \beta_t\ following\ a\ schedule\ \beta_1 < \beta_2 < ... < \beta_T$

*Linear scheduling*

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

0.02    0.0001    1000

*Cosinus scheduling*

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

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

Linear and Cosine schedules



## Beta scheduling

$q(x_0)$

*Données*

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

*Données bruitées
par un petit bruit
gaussien*

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

*Données bruitées
à l'étape t*

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

**Computation 0**

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

**Chaîne de Markov**

$$q(x_{t-1}|x_t, x_0) = \frac{q(x_t|x_{t-1}, x_0)q(x_{t-1}|x_0)}{q(x_t|x_0)} = \boxed{\frac{q(x_t|x_{t-1})q(x_{t-1}|x_0)}{q(x_t|x_0)}}$$

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

**Computation 1 - Gaussian reverse diffusion**

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

$$\tilde{\mu}_t(\mathbf{x_t}, \mathbf{x_0}) = \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1 - \bar{\alpha}_t}\mathbf{x_0} + \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t}\mathbf{x_t} \qquad\qquad \tilde{\beta}_t = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t}\beta_t$$

---

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

$$x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}z_t \qquad\longrightarrow\qquad \tilde{\mu}_t(\mathbf{x_t}, \mathbf{x_0}) = \frac{1}{\sqrt{\alpha_t}}(\mathbf{x_t} - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}}\mathbf{z_t})$$

$$\mathbf{z_t} \quad ???$$

**Computation 2 - Mean and Variance of the reverse diffusion**

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

**Reverse diffusion 1**
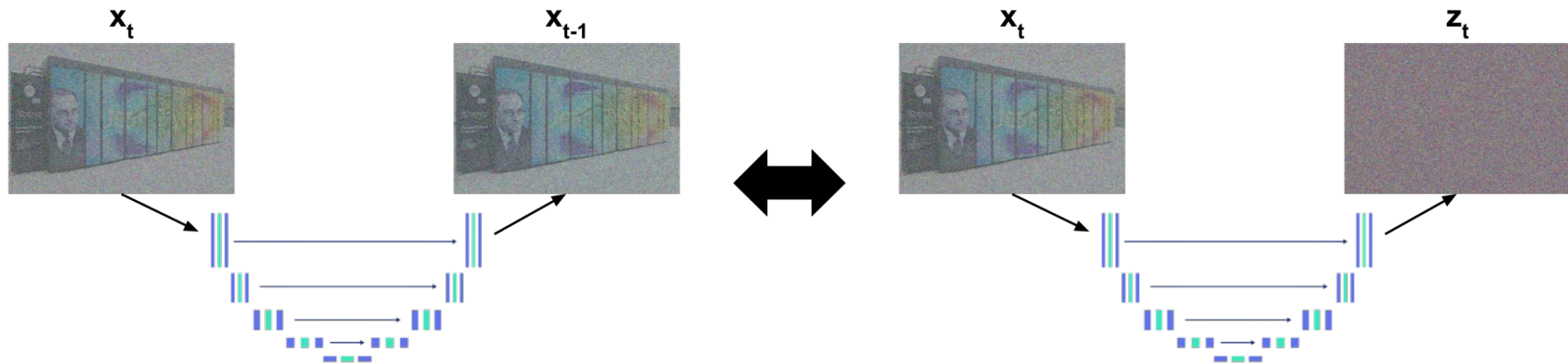
$$\tilde{\mu}_t(\mathbf{x_t}, \mathbf{x_0}) = \frac{1}{\sqrt{\alpha_t}}\left(\mathbf{x_t} - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}}\mathbf{z_t}\right) \approx \frac{1}{\sqrt{\alpha_t}}\left(x_t - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}}\mathbf{z}_\theta(\mathbf{x_t}, \mathbf{t})\right)$$
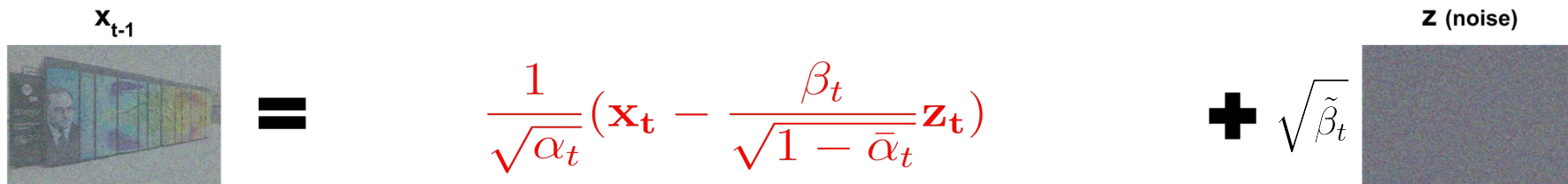
*Our network*



$\mathbf{x_{t-1}}$

$\mathbf{z}$ (noise)

$$= \frac{1}{\sqrt{\alpha_t}}\left(\mathbf{x_t} - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}}\mathbf{z_t}\right) \quad + \sqrt{\tilde{\beta}_t}$$

**Reverse diffusion 2**

**Algorithm 1** Training

1: **repeat**
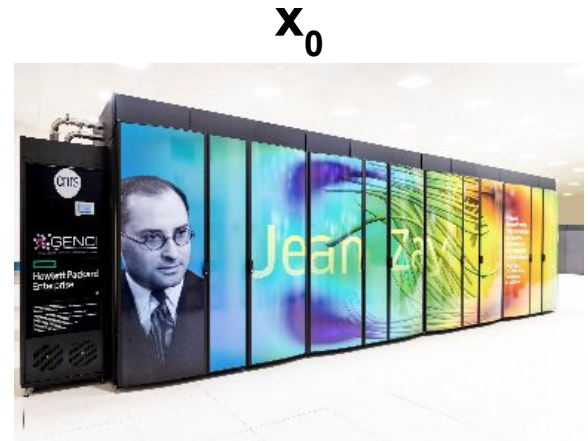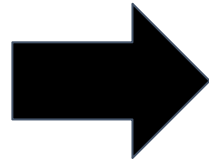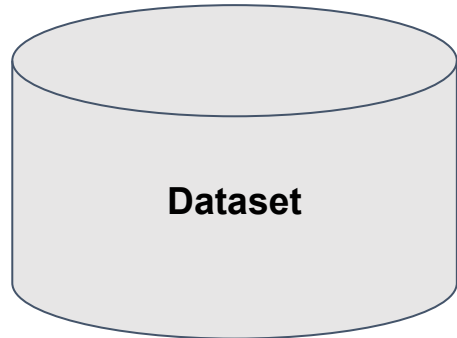2: $\quad \mathbf{x}_0 \sim q(\mathbf{x}_0)$
3: $\quad t \sim \text{Uniform}(\{1, \ldots, T\})$
4: $\quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
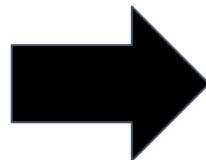5: $\quad$ Take gradient descent step on
$$\nabla_\theta \left\| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta(\sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\epsilon}, t) \right\|^2$$
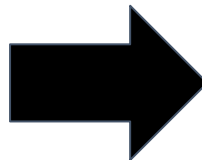6: **until** converged

**Training 1**

$x_0$

Dataset

**Uniform distribution**

Between 1 and T

**t = 50**

$x_0$

$$\mathbf{z_t} \ (= \epsilon)$$

Gaussian distribution

Same shape than $x_0$

$x_0$

t = 50

$$\mathbf{x_t} \qquad \mathbf{x_0} \qquad \mathbf{z_t}$$



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

$\mathbf{x_0}$

$\mathbf{z_t}$



**t = 50**

$x_t$

$z_\theta \ (= \epsilon_\theta)$

t = 50

$x_0$

t = 50

$z_t$

$x_t$

**Training 6**

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

$x_0$    $t = 50$    $z_t$    $x_t$    $z_\theta$

**Backward**

**Loss**

$\nabla_\theta$

**Weight update**

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

$$L_{\mathrm{VLB}} = L_T + L_{T-1} + \cdots + L_0$$

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

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

**A more complicated loss …**

## Algorithm 2 Sampling

1: $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
2: **for** $t = T, \ldots, 1$ **do**
3:    $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ if $t > 1$, else $\mathbf{z} = \mathbf{0}$
4:    $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$
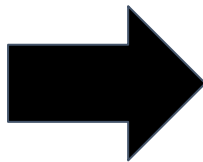5: **end for**
6: **return** $\mathbf{x}_0$

**Sampling 1**

$\mathbf{x_T}$



**Gaussian distribution**

Same shape than training dataset images

**Sampling 2**

$x_T$

$z_\theta(x_T,T) \approx z_T$

T

$x_T$

**z (noise)**



**Gaussian distribution**

Same shape than training dataset images

$x_T$



$z_T$



**Sampling 3**

**x**$_{T-1}$ ≈ $\frac{1}{\sqrt{\alpha_T}}$ ( **x**$_T$ — $\frac{\beta_T}{\sqrt{1-\bar{\alpha}_T}}$ **z**$_\theta$(**x**$_T$,**T**) ) ➕ $\tilde{\beta}_T$ **z** (noise)

**x**$_T$  $z_\theta(x_T, T)$  **z ( random 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.**

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

**New mathematical view :**

$$\nabla_{\mathbf{x_t}} \log(\mathbf{q}(\mathbf{x_t})) = -\frac{\mathbf{z}_\theta(\mathbf{x_t}, \mathbf{t})}{\sqrt{1 - \bar{\alpha}_t}}$$

$$\tilde{\mu}_t(x_t, x_0) \approx \frac{1}{\sqrt{\alpha_t}}(x_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \mathbf{z}_\theta(\mathbf{x_t}, \mathbf{t})) = \frac{1}{\sqrt{\alpha_t}}(x_t + \beta_t \times \nabla_{\mathbf{x_t}} \log(\mathbf{q}(\mathbf{x_t})))$$

$$\nabla_{\mathbf{x_t}} \log(\mathbf{q}(\mathbf{x_t})) \longrightarrow \nabla_{\mathbf{x_t}} \log(\mathbf{q}(\mathbf{x_t}|\mathbf{y}))$$

**Unconditional to Conditional Diffusion Model**

$$\nabla_{\mathbf{x_t}} \log(\mathbf{q}(\mathbf{x_t}|\mathbf{y})) = \nabla_{\mathbf{x_t}} \log(\mathbf{q}(\mathbf{x_t})) + \nabla_{\mathbf{x_t}} \log(\mathbf{q}(\mathbf{y}|\mathbf{x_t}))$$

$$\tilde{\mu}_t(x_t, x_0, y) \approx \frac{1}{\sqrt{\alpha_t}} (x_t + \beta_t \times [\nabla_{\mathbf{x_t}} \log(\mathbf{q}(\mathbf{x_t})) + \lambda \times \nabla_{\mathbf{x_t}} \log(\mathbf{q}(\mathbf{y}|\mathbf{x_t}))])$$

**Strength of guidance**

**Unconditional pretrained model**

**Classifier output**

Input layer   Hidden layer   Hidden layer   Output layer

Noisy Jean Zay ?
A cat ?

**Classifier Guidance**

Principle of the Classifier Guidance method

$$\frac{1}{\sqrt{\alpha_t}}(x_t + \beta_t \times [\nabla_{\mathbf{x_t}}\log(\mathbf{q}(\mathbf{x_t})) + \lambda \times \nabla_{\mathbf{x_t}}\log(\mathbf{q}(\mathbf{y}|\mathbf{x_t}))])$$

$$\nabla_{\mathbf{x_t}}\log(\mathbf{q}(\mathbf{y}|\mathbf{x_t})) = \nabla_{\mathbf{x_t}}\log(\mathbf{q}(\mathbf{x_t}|\mathbf{y})) - \nabla_{\mathbf{x_t}}\log(\mathbf{q}(\mathbf{x_t}))$$

**x_t**

**t = 50**

**y = "Supercomputer"**

**Classifier-Free Guidance**

$logP(x)$

VAE - *Have more diversity*

DDPM - *Long sampling process* ?

GAN
- *Generate high quality data*
- *Hard to train*

$x_T$

$FID$

**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)

**DDPM vs GAN vs VAE**

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

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

$$\longrightarrow$$

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

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

**Reduce the sampling step**

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

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

$$\sigma_t^2 = \tilde{\beta_t} \implies DDPM$$

**Important :**
Same network and
training as a DDPM

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

**DDIM 1**

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

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

$$\eta = 0 \qquad \Rightarrow \text{DDIM}$$

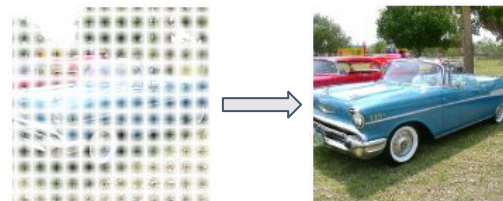|  | | CIFAR10 (32 × 32) | | | | | CelebA (64 × 64) | | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| $S$ | | 10 | 20 | 50 | 100 | 1000 | 10 | 20 | 50 | 100 | 1000 |
| | 0.0 | **13.36** | **6.84** | **4.67** | **4.16** | 4.04 | **17.33** | **13.73** | **9.17** | **6.53** | 3.51 |
| | 0.2 | 14.04 | 7.11 | 4.77 | 4.25 | 4.09 | 17.66 | 14.11 | 9.51 | 6.79 | 3.64 |
| $\eta$ | 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:

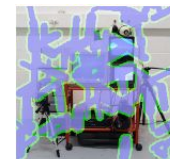➔ Inpainting



➔ Super-resolution



➔ Outpainting



## Other task 1

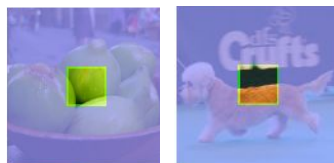We can solve many of these tasks through the usage of a mask

➜ Wide mask

➜ Thin mask

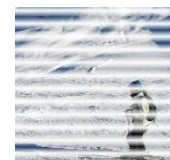➜ Outer mask for expanding the image

➜ 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.*

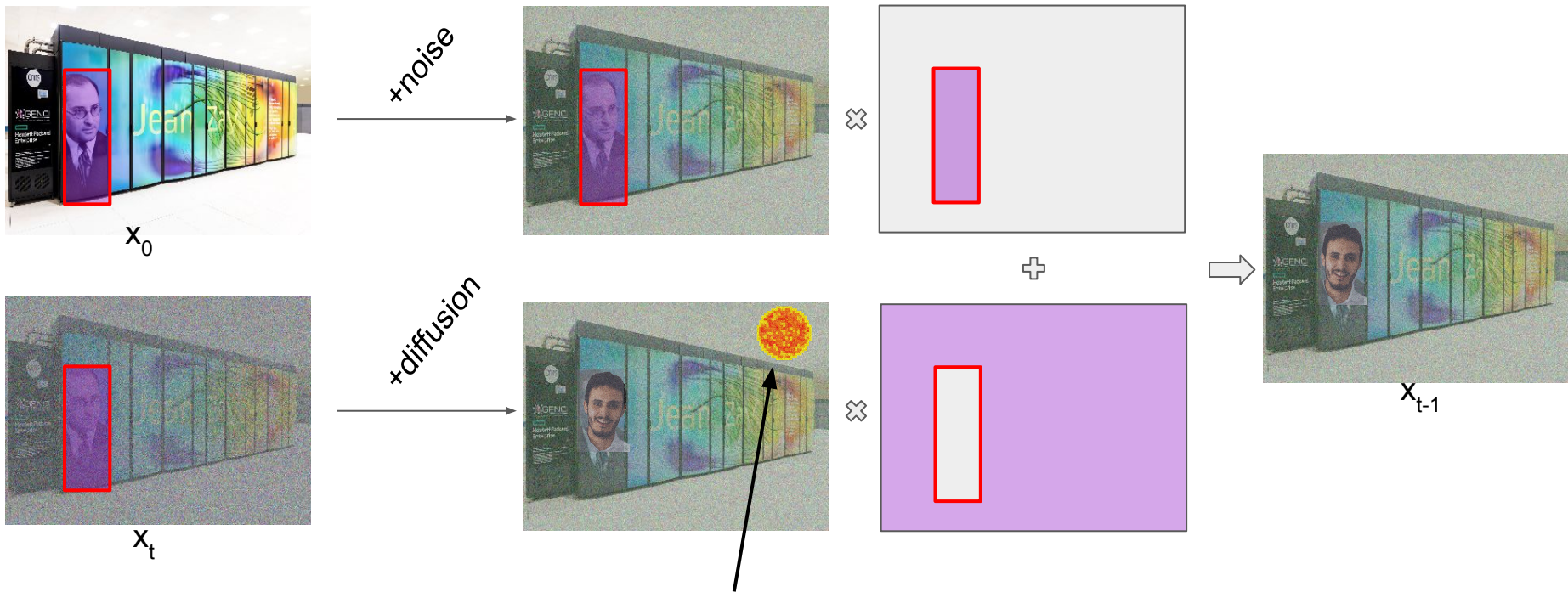## Other task 2

**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.*

## Other task 3

$x_0$

$x_t$

+noise

+diffusion

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

$x_{t-1}$

**Other task 4**

**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 (IDDPM) (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)
- Classifier-free diffusion guidance (https://arxiv.org/pdf/2207.12598)
- Score-Based Generative Modeling through Stochastic Differential Equation (https://openreview.net/pdf?id=PxTIG12RRHS)
- 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)

**Online Resources**