

Diffusion

Romain Loiseau

Slides mainly taken from a talk from **Yannis Siglidis** (IMAGINE/ENPC)

Diffusion

- 1. What is diffusion and how does it work?
- 2. Conditioning diffusion to match a description ?
- 3. Diffusion for 3D !

I didn't knew anything about diffusion 10 days ago so **please interrupt me** if:

- You don't understand anything I'm saying and want me to clarify
- You want to add more details you know
- You want me to stop explaining useless math things because it's too painful
- I'm saying stupid things

Key ideas for each slide in red here

History

Papers on Diffusion





First results mainly p.o.c. (cifar, texture synthesis, etc) Sohl-Dickstein et al., 2015 -> Nonequilibrium Thermodynamics

"DDPM": Realistic faces, churches, bedrooms, etc. Ho et al. 2020,

BigGAN



Real



"Diffusion models beat gans on image synthesis." (SoTA + conditional generation) Nichol & Dhariwal (2021)

State of the art:

DALL-E 2







Stable Diffusion

Stable Diffusion is a machine learning, text-to-image model to generate digital images from natural language descriptions. The underlying approach was developed at LMU Munich and then extended by a collaboration of StabilityAI, LMU, and Runway with support from EleutherAI and LAION. Wikipedia

Compositional Generalization!

A generative model.





The Markov chain diffusion process of generating a sample

Ho et al. 2020

Backward

$$\begin{array}{c} \mathbf{x}_{T} \xrightarrow{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_{t})} \underbrace{\mathbf{x}_{t-1}}_{q(\mathbf{x}_{t}|\mathbf{x}_{t-1})} \xrightarrow{q(\mathbf{x}_{t}|\mathbf{x}_{t-1})} \underbrace{\mathbf{x}_{t-1}}_{q(\mathbf{x}_{t}|\mathbf{x}_{t-1})} \xrightarrow{q(\mathbf{x}_{t}|\mathbf{x}_{t-1})} \underbrace{\mathbf{x}_{t-1}}_{q(\mathbf{x}_{t}|\mathbf{x}_{t-1})} \xrightarrow{q(\mathbf{x}_{t}|\mathbf{x}_{t-1})} \underbrace{\mathbf{x}_{t-1}}_{q(\mathbf{x}_{t}|\mathbf{x}_{t-1})} \xrightarrow{q(\mathbf{x}_{t}|\mathbf{x}_{t-1})} \underbrace{\mathbf{x}_{t-1}}_{q(\mathbf{x}_{t}|\mathbf{x}_{t-1})} \xrightarrow{q(\mathbf{x}_{t}|\mathbf{x}_{t-1})} \underbrace{\mathbf{x}_{t-1}}_{q(\mathbf{x}_{t}|\mathbf{x}_{t-1})} \underbrace{\mathbf{x}_{t-1}}_{q(\mathbf{x}_{t}|\mathbf{x}_{t-1})} \underbrace{\mathbf{x}_{t-1}}_{q(\mathbf{x}_{t}|\mathbf{x}_{t-1})} \xrightarrow{q(\mathbf{x}_{t}|\mathbf{x}_{t-1})} \underbrace{\mathbf{x}_{t-1}}_{q(\mathbf{x}_{t}|\mathbf{x}_{t-1})} \underbrace{\mathbf{x}_{t-1}}_{q(\mathbf{x}_{t}|\mathbf{x}_{t-1})}$$

$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I}) \xrightarrow{\alpha_t = 1 - \beta_t} q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I})$$
$$\bar{\alpha}_t = \prod_{i=1}^T \alpha_i$$

Sohl-Dickstein et al., 2015

Step Independent Computation is key for the training to be tractable

Forward
$$q(\mathbf{x}_{t}|\mathbf{x}_{0}) = \mathcal{N}(\mathbf{x}_{t}; \sqrt{\bar{\alpha}_{t}}\mathbf{x}_{0}, (1 - \bar{\alpha}_{t})\mathbf{I}) \xrightarrow{\mathbf{x}_{T}} \longrightarrow \cdots \longrightarrow \mathbf{x}_{t} \xrightarrow{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_{t})} \xrightarrow{\mathbf{x}_{t-1}} \longrightarrow \cdots \longrightarrow \mathbf{x}_{0}$$
$$\overbrace{q(\mathbf{x}_{t}|\mathbf{x}_{t-1})}^{p_{\theta}(\mathbf{x}_{t}|\mathbf{x}_{t})} \xrightarrow{\mathbf{x}_{t-1}} \overbrace{\mathbf{x}_{t-1}}^{p_{\theta}(\mathbf{x}_{t}|\mathbf{x}_{t})} \longrightarrow \mathbf{x}_{0}$$
$$\overbrace{q(\mathbf{x}_{t-1}|\mathbf{x}_{t}, \mathbf{x}_{0})}^{q(\mathbf{x}_{t}|\mathbf{x}_{0}, \mathbf{x}_{0})} = q(\mathbf{x}_{t}|\mathbf{x}_{t-1}, \mathbf{x}_{0}) \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_{0})}{q(\mathbf{x}_{t}|\mathbf{x}_{0})} \sim \mathcal{N}(\mathbf{x}_{t-1}; \tilde{\boldsymbol{\mu}}(\mathbf{x}_{t}, \mathbf{x}_{0}), \tilde{\boldsymbol{\beta}}_{t}\mathbf{I})$$
$$\widetilde{\boldsymbol{\mu}}_{t} \equiv \tilde{\boldsymbol{\mu}}_{t}(\mathbf{x}_{t}, \mathbf{x}_{0}) = \frac{1}{\sqrt{\alpha_{t}}} \left(\mathbf{x}_{t} - \frac{\beta_{t}}{\sqrt{1 - \bar{\alpha}_{t}}}\mathbf{z}_{t}\right) \qquad \tilde{\boldsymbol{\beta}}_{t} = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_{t}} \cdot \boldsymbol{\beta}_{t}$$
Close form transition kernel conditioned on \mathbf{x}_{0}

$$p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_{t}) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_{\theta}(\mathbf{x}_{t}, t), \boldsymbol{\Sigma}_{\theta}(\mathbf{x}_{t}, t))$$
$$\boldsymbol{\mu}_{\theta}(\mathbf{x}_{t}, t) = \frac{1}{\sqrt{\alpha_{t}}} \Big(\mathbf{x}_{t} - \frac{\beta_{t}}{\sqrt{1 - \bar{\alpha}_{t}}} \mathbf{z}_{\theta}(\mathbf{x}_{t}, t) \Big)$$
$$\text{Original: } \boldsymbol{\Sigma}_{\theta}(\mathbf{x}_{t}, t) = \sigma_{t}^{2} \mathbf{I}$$
Model

Sohl-Dickstein et al., 2015 Ask the model to denoise -> estimate the mean and variance of the applied noise



Sohl-Dickstein et al., 2015

Train it using a single loss at time t

Sohl-Dickstein et al., 2015

Ask the model to denoise

Training

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

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

: **until** converged

Ho et al. (2020)

6

Sampling



Algorithm 2 Sampling $p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_{\theta}(\mathbf{x}_t, t), \boldsymbol{\Sigma}_{\theta}(\mathbf{x}_t, t))$

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

6: return \mathbf{x}_0

Ho et al. (2020)





 $oldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0},\mathbf{I})$ must be kept the same across all steps.

Ho et al. (2020)

cvpr2022-diffusion

Why does it work?

Fourier Transform $\mathbf{x}_{t} = \sqrt{\bar{\alpha}_{t}}\mathbf{x}_{0} + \sqrt{(1 - \bar{\alpha}_{t})}\epsilon$ $\mathcal{F}(\mathbf{x}_{t}) = \sqrt{\bar{\alpha}_{t}}\mathcal{F}(\mathbf{x}_{0}) + \sqrt{(1 - \bar{\alpha}_{t})}\mathcal{F}(\epsilon)$



Low frequency dominance of image responses ——— high frequency information is perturbed faster.

You all know everything about how diffusion works know

Let's dive into some fun stuffs now

Conditioning the Model

An astronaut riding a horse in a photorealistic style



Conditional Model: $ilde{m{\epsilon}}_{ heta}\left(\mathbf{z}_{t},t,\mathbf{c}
ight)$

E.g. The network can be a UNET and **c** can be a BERT embedding.

$$\begin{array}{l} \text{Sample with: } \tilde{\boldsymbol{\epsilon}}_{\theta}\left(\mathbf{z}_{t},t,\mathbf{c}\right)=(1{+}w)\underbrace{\boldsymbol{\epsilon}_{\theta}\left(\mathbf{z}_{t},t,\mathbf{c}\right)}_{\text{Conditional}}{-}w\underbrace{\boldsymbol{\epsilon}_{\theta}\left(\mathbf{z}_{t},t,\mathbf{c}=\mathbf{0}\right)}_{\text{Unconditional}} \end{array}$$

good balance between FID (distinguish between synthetic and generated images) and IS (quality and diversity)

Ho & Salimans (2021)

Classifier Free Guidance

Proof

Bayes formula

$$\nabla_{\mathbf{x}_{t}} \log p(y|\mathbf{x}_{t}) = \nabla_{\mathbf{x}_{t}} \log p(\mathbf{x}_{t}|y) - \nabla_{\mathbf{x}_{t}} \log p(\mathbf{x}_{t})$$

$$= -\frac{1}{\sqrt{1 - \bar{\alpha}_{t}}} \left(\epsilon_{\theta}(\mathbf{x}_{t}, t, y) - \epsilon_{\theta}(\mathbf{x}_{t}, t) \right)$$
gradient of the log likelihood of an auxiliary classifier model p $\theta(y|\mathbf{x}_{t})$

$$\bar{\epsilon}_{\theta}(\mathbf{x}_{t}, t, y) = \epsilon_{\theta}(\mathbf{x}_{t}, t, y) - \sqrt{1 - \bar{\alpha}_{t}} w \nabla_{\mathbf{x}_{t}} \log p(y|\mathbf{x}_{t})$$

$$= \epsilon_{\theta}(\mathbf{x}_{t}, t, y) + w \left(\epsilon_{\theta}(\mathbf{x}_{t}, t, y) - \epsilon_{\theta}(\mathbf{x}_{t}, t) \right)$$

$$= (w + 1)\epsilon_{\theta}(\mathbf{x}_{t}, t, y) - w\epsilon_{\theta}(\mathbf{x}_{t}, t)$$

Latent Space Diffusion



Latent Space Diffusion



DREAMFUSION, Poole et al. (2022)



zoomed out view of Tower Bridge made out of gingerbread and candy[‡]

a robot and dinosaur playing chess, high resolution*

a squirrel gesturing in front of an easel showing colorful pie charts

DREAMFUSION, Poole et al. (2022)



$$\nabla_{\theta} \mathcal{L}_{\text{SDS}}(\phi, \mathbf{x} = g(\theta)) \triangleq \mathbb{E}_{t,\epsilon} \left[w(t) \left(\hat{\epsilon}_{\phi}(\mathbf{z}_t; y, t) - \epsilon \right) \frac{\partial \mathbf{x}}{\partial \theta} \right]$$

create 3D models that look like good images when rendered from random angles

3D-Diffusion, Watson et al. (2022)



3D-Diffusion, Watson et al. (2022)



3D-Diffusion, Watson et al. (2022)



🖉 craiyon

AI model drawing images from any prompt!

