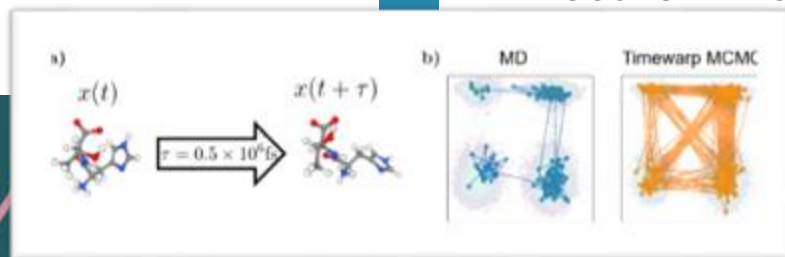# Diffusion Models

Chin-Wei Huang & Victor Garcia Satorras
Senior Researchers @ AI for Science – Microsoft Research Amsterdam
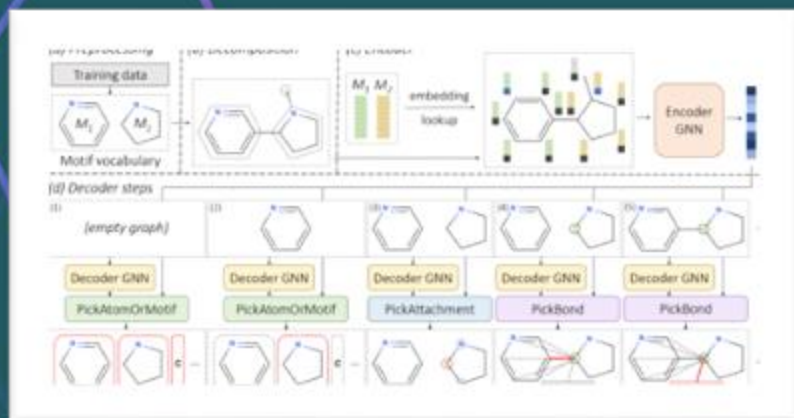
# Microsoft Research AI for Science
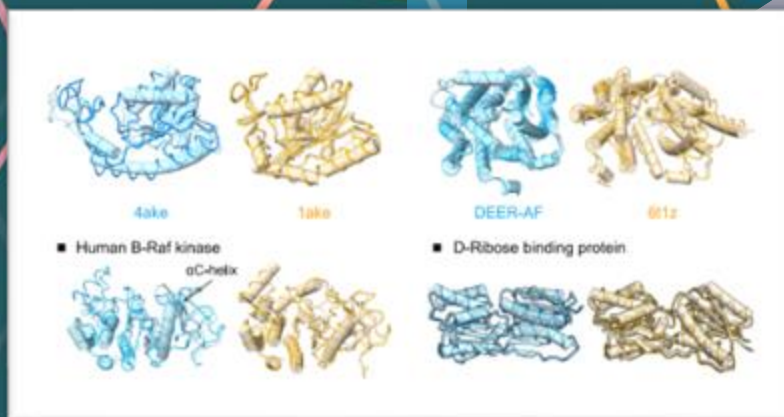
**Timewarp** for **long time-scale** MD simulation
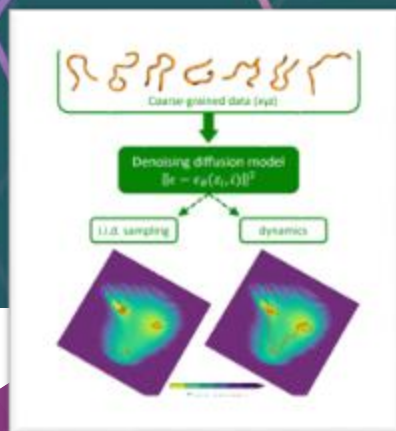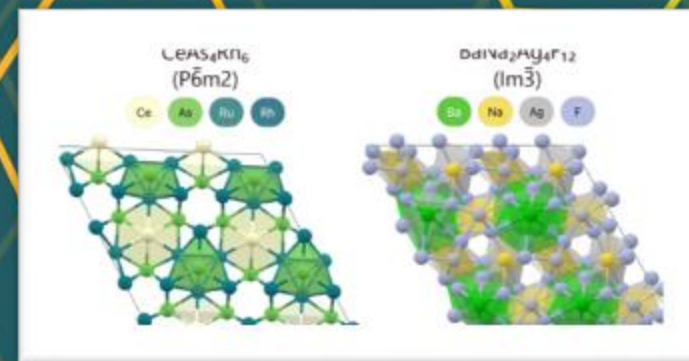
**MoLeR** for **drug-like molecular** graph generation

**Graphormer** for **protein structure** generation

**Two-for-one**: **coarse-grained** molecular simulation

**MatterGen** for **inorganic materials** design

LENGTH SCALE

TIME SCALE

# Outline

1. Introduction to Denoising Diffusion Models

2. Equivariant Diffusion Models for Molecule Generation in 3D

3. Workshop: Code & Practice

*Slides borrowed from*

Rianne van den Berg
Principle Research Manager
AI4Science, MSR
Amsterdam

Daniel Zuegner
Senior Researcher
AI4Science, MSR
Berlin

Sarah Lewis
Senior RSDE
AI4Science, MSR
Cambridge, UK

# Outline

1. **Introduction to Denoising Diffusion Models**
2. Equivariant Diffusion Models for Molecule Generation in 3D
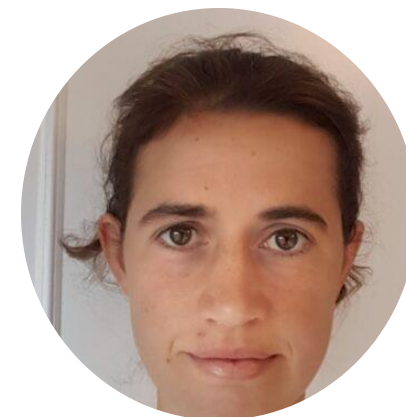3. Workshop: Code & Practice

# Very deep VAE / normalizing flows



Two networks $p_\theta$ and $q_\phi$ to optimize

Focus on data -> prior transformation

# Major drawbacks of deep VAE / flow

- Space-time complexity grows with depth
- Arbitrary dynamics



Image from Papamakarios et al JMLR 2021

# What is a diffusion model?

- Forward diffusion process gradually destroys information in the data.



Forward diffusion process (fixed)

Data → Noise

# What is a diffusion model?

- Forward diffusion process gradually destroys information in the data.

- A generative model that learns to **revert** a **diffusion process**.

# What is a diffusion model?

Sohl-Dickstein et al., ICML 2015, Ho et al., NeurIPS 2020, Song et al., ICLR 2021

- Forward diffusion process gradually destroys information in the data.

- A generative model that learns to **revert** a **diffusion process**.

- **Alternative perspectives** on diffusion models:
  - Deep generative models with latent variables & ELBO maximization
  - (Denoising) score matching



Forward diffusion process (fixed)

Data — Noise

Reverse denoising process (generative)

# Recap: Latent Variable Generative Model

A latent (unobserved) random variable

$$p_\theta(x) = \int p_\theta(x, z)\, \mathrm{d}z$$

Goal: maximize $\mathbb{E}_{q(x)}[\log p_\theta(x)]$ & sample $x \sim p_\theta(x)$

Defining joint distribution as $p_\theta(x, z) = p_\theta(x|z)p(z)$ allows ancestral sampling of $(x, z) \sim p_\theta(x, z)$.

# Evidence lower bound (ELBO)

$$\log p_\theta(x) = \log \int \ldots \int p_\theta(x, z)\, dz$$

# Evidence lower bound (ELBO)

$$\log p_\theta(x) = \log \int \dots \int p_\theta(x, \textcolor{red}{z}) \, d\textcolor{red}{z}$$

$$= \log \int \dots \int \frac{q(z|x)}{q(z|x)} p_\theta(x, z) \, dz$$

$$= \log \mathbb{E}_{q(z_1, \dots, z_T | x)} \left[ \frac{p_\theta(x, z)}{q(z|x)} \right]$$

$$\geq \mathbb{E}_{q(z_1, \dots, z_T | x)} \log \left[ \frac{p_\theta(x, z)}{q(z|x)} \right]$$

# Introducing (a lot of) latent variables

$$\log p_\theta(x) = \log \int \dots \int p_\theta(x, z)\, dz$$

$$= \log \int \dots \int \frac{q(z|x)}{q(z|x)} p_\theta(x, z)\, dz$$

$$= \log \mathbb{E}_{q(z_1,\dots,z_T|x)}\left[\frac{p_\theta(x, z)}{q(z|x)}\right]$$

$$\geq \mathbb{E}_{q(z_1,\dots,z_T|x)} \log\left[\frac{p_\theta(x, z)}{q(z|x)}\right]$$

$$= \log \int \dots \int p_\theta(x, z_1, \dots, z_T)\, dz_1 \dots dz_T$$

$$= \log \int \dots \int \frac{q(z_1, \dots, z_T|x)}{q(z_1, \dots, z_T|x)} p_\theta(x, z_1, \dots, z_T)\, dz_1 \dots dz_T$$

$$= \log \mathbb{E}_{q(z_1,\dots,z_T|x)}\left[\frac{p_\theta(x, z_1, \dots, z_T)}{q(z_1, \dots, z_T|x)}\right]$$

$$\geq \mathbb{E}_{q(z_1,\dots,z_T|x)} \log\left[\frac{p_\theta(x, z_1, \dots, z_T)}{q(z_1, \dots, z_T|x)}\right]$$

How to factorize $p_\theta(x, z_1, \dots, z_T)$ and $q(z_1, \dots, z_T|x)$?

# How to factorize $q(z_1, \ldots, z_T | x)$ & $p_\theta(z_T, \ldots z_1, x)$?

Sohl-Dickstein et al., ICML 2015



("Inference" of latent variables)

$x$    $z_1$    $z_2$          $z_T$

Data                                                 Noise

Latent
variables

Forward: Define a *Markov* Chain to "infer" the latents given the data:

$$q(z_1, \ldots z_T | x) = q(z_1 | x) q(z_2 | z_1) \ldots q(z_T | z_{T-1})$$

Backward:

$$p_\theta(z_T, \ldots z_1, x) = p(z_T) p_\theta(z_{T-1} | z_T) \ldots p_\theta(z_1 | z_2) p_\theta(x | z_1)$$

# Organizing the negative ELBO

Sohl-Dickstein et al., ICML 2015, Ho et al., NeurIPS 2020, Song et al., ICLR 2021

Definitions of p and q:

$$p_\theta(z_T, \ldots z_1, x) = p(z_T)p_\theta(z_{T-1}|z_T) \ldots p_\theta(z_1|z_2)p_\theta(x|z_1)$$
$$q(z_1, \ldots, z_T|x) = q(z_1|x)q(z_2|z_1) \ldots q(z_T|z_{T-1})$$

Plugging it all in:

$$-\log p_\theta(x) \leq -\mathbb{E}_{q(z_1, \ldots, z_T|x)} \log \left[ \frac{p_\theta(x, z_1, \ldots, z_T)}{q(z_1, \ldots, z_T|x)} \right]$$

$$= -\mathbb{E}_{q(z_1, \ldots, z_T|x)} \log \left[ \frac{p(z_T)p_\theta(z_{T-1}|z_T) \ldots p(z_{t-1}|z_t) \ldots p_\theta(z_1|z_2)p_\theta(x|z_1)}{q(z_T|x)q(z_{T-1}|z_T, x) \ldots q(z_{t-1}|z_t, \ldots, z_T, x) \ldots q(z_1|z_2, \ldots, z_T, x)} \right]$$

$$= -\mathbb{E}_{q(z_1, \ldots, z_T|x)} \log \left[ \frac{p(z_T)p_\theta(z_{T-1}|z_T) \ldots p(z_{t-1}|z_t) \ldots p_\theta(z_1|z_2)p_\theta(x|z_1)}{q(z_T|x)q(z_{T-1}|z_T, x) \ldots q(z_{t-1}|z_t, x) \ldots q(z_1|z_2, x)} \right]$$

$$= KL(q(z_T|x)\|p(z_T)) + \sum_{t=2}^{T} \mathbb{E}_{q(z_t|x)}[KL(q(z_{t-1}|z_t, x)\|p(z_{t-1}|z_t))] - \mathbb{E}_{q(z_1|x)}[\log p_\theta(x|z_1)]$$

# Training diffusion models



$$q(z_1, \dots z_T | x) = q(z_1 | x) q(z_2 | z_1) \dots q(z_T | z_{T-1})$$

Data

$x$   $z_1$   $z_2$   $z_T$   Noise

$$p_\theta(z_T, \dots z_1, x) = p(z_T) p_\theta(z_{T-1} | z_T) \dots p_\theta(z_1 | z_2) p_\theta(x | z_1)$$

$$L_{vb} = \mathbb{E}_{q(x)} \left[ -\mathbb{E}_{q(z_1 | x)} \log p_\theta(x | z_1) \right] + \sum_{t=2}^{T} \mathbb{E}_{q(z_t | x)} [KL[q(z_{t-1} | z_t, x) || p_\theta(z_{t-1} | z_t)]]$$
$$+ \ KL[q(z_T | x) || p_\theta(z_T)]$$

Practical requirements for $q$ and $p_\theta$ to allow for efficient "simulation-free" training of $p_\theta$:

1. Efficient sampling of $z_t$ from $q(z_t | x)$ for arbitrary time $t$

2. Tractable expression for $q(z_{t-1} | z_t, x)$ (and the KL divergence).

For Gaussian or binomial $q(z_t | z_{t-1})$ (and $p_\theta(z_{t-1} | z_t)$): ☑ ☑

# Diffusion models with Gaussian distributions

Gaussian forward distributions: $q(z_t|z_{t-1}) = \mathcal{N}(z_t|\sqrt{\beta_t}z_{t-1}, (1-\beta_t)I)$     $(\beta_t < 1)$

1. Sampling arbitrary timesteps in one shot:

$$q(z_t|x) = \mathcal{N}(z_t|\sqrt{\bar{\alpha}_t}x, (1-\bar{\alpha}_t)I)$$

$$\alpha_t = 1 - \beta_t$$
$$\bar{\alpha}_t = \prod_{s=1}^{t} \alpha_s$$

2. Tractable expression for posterior:

$$q(z_{t-1}|z_t, x) = \mathcal{N}(z_{t-1}|\tilde{\mu}_t(z_t, x), \tilde{\beta}_t I)$$

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

# Additional guarantees

$$-\log p_\theta(x) \leq -\mathbb{E}_{q(z_1|x)} \log p_\theta(x|z_1) + \sum_{t=2}^{T} \mathbb{E}_{q(z_t|x)}[KL[q(z_{t-1}|z_t,x)||p_\theta(z_{t-1}|z_t)]]$$

$$+ KL[q(z_T|x)||p(z_T)]$$

When $T \to \infty$:

1. For many choices of the noise schedule, we get a stationary distribution
   $\lim_{T' \to \infty} q(z_{T'}|x) = N(0, I)$.

   $\to$ pick $p(z_T) = N(0, I)$

2. As $T \to \infty$, the optimal reverse (generative) becomes Gaussian-like.

   $\to p_\theta(z_{t-1}|z_t) = \mathcal{N}(z_{t-1}|\mu_\theta(z_t, t), \sigma_t I)$

# Parameterizing $\mu_\theta(z_t, t)$

Ho et al., NeurIPS 2020

$$\mathbb{E}_{q(z_t|x)}[KL[q(z_{t-1}|z_t, x)||p_\theta(z_{t-1}|z_t)]] = \mathbb{E}_{q(z_t|x)}\left[\frac{1}{2\sigma_t^2}||\tilde{\mu}_t(z_t, x) - \mu_\theta(z_t, t)||^2\right] + C$$

**Simplest option**: $\mu_\theta(z_t, t) = \text{nn}_\theta(z_t, t)$

Recall:

$$z_t = \sqrt{\bar{\alpha}_t}x + \sqrt{1 - \bar{\alpha}_t}\epsilon$$

**Alternative**:

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

Predict the noise:

$$\mu_\theta(z_t, t) = \frac{1}{\sqrt{\alpha_t}}\left(z_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}}\epsilon_\theta(z_t, t)\right)$$

# Loss, noise-scheduling, training and sampling

Ho et al., NeurIPS 2020

$$L_{vlb} = \mathbb{E}_{q(x)}\mathbb{E}_{\epsilon\sim\mathcal{N}(0,I)}\mathbb{E}_{t\sim U(2,T)} \left[ \boxed{\frac{\beta_t^2}{2\sigma_t^2\alpha_t\,(1-\bar{\alpha}_t)}} \left\| \epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t}x + \sqrt{1-\bar{\alpha}_t}\epsilon, t) \right\|_2^2 \right] + \text{reconstruct}$$

$$\lambda_t$$

Ho et al. 2020 found that setting $\lambda_t = 1$ improves sample quality, i.e., the training loss is:

$$L_{simple} = \mathbb{E}_{q(x)}\mathbb{E}_{\epsilon\sim\mathcal{N}(0,I)}\mathbb{E}_{t\sim U(2,T)} \left[ \left\| \epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t}x + \sqrt{1-\bar{\alpha}_t}\epsilon, t) \right\|_2^2 \right] + \text{reconstruct}$$

**Algorithm 1** Training

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

**Algorithm 2** Sampling

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

# Results from Ho et al. 2020

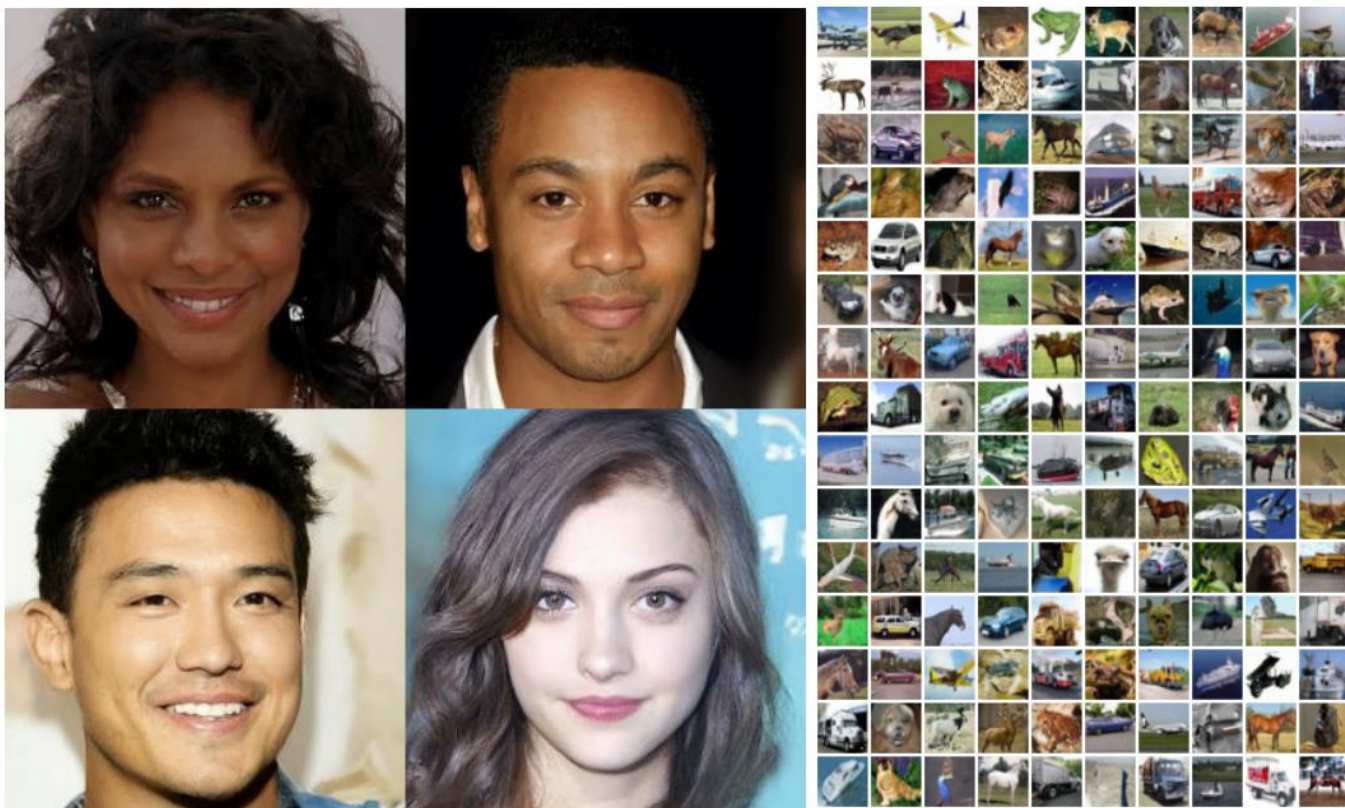Sohl-Dickstein et al., ICML 2015, Ho et al., NeurIPS 2020



Figure 1: Generated samples on CelebA-HQ 256 × 256 (left) and unconditional CIFAR10 (right)

Table 1: CIFAR10 results. NLL measured in bits/dim.

| Model | IS | FID | NLL Test (Train) |
|---|---|---|---|
| **Conditional** | | | |
| EBM [11] | 8.30 | 37.9 | |
| JEM [17] | 8.76 | 38.4 | |
| BigGAN [3] | 9.22 | 14.73 | |
| StyleGAN2 + ADA (v1) [29] | **10.06** | **2.67** | |
| **Unconditional** | | | |
| Diffusion (original) [53] | | | $\leq 5.40$ |
| Gated PixelCNN [59] | 4.60 | 65.93 | 3.03 (2.90) |
| Sparse Transformer [7] | | | **2.80** |
| PixelIQN [43] | 5.29 | 49.46 | |
| EBM [11] | 6.78 | 38.2 | |
| NCSNv2 [56] | | 31.75 | |
| NCSN [55] | $8.87 \pm 0.12$ | 25.32 | |
| SNGAN [39] | $8.22 \pm 0.05$ | 21.7 | |
| SNGAN-DDLS [4] | $9.09 \pm 0.10$ | 15.42 | |
| StyleGAN2 + ADA (v1) [29] | $\mathbf{9.74} \pm 0.05$ | 3.26 | |
| Ours ($L$, fixed isotropic $\Sigma$) | $7.67 \pm 0.13$ | 13.51 | $\leq 3.70$ (3.69) |
| **Ours** ($L_{\text{simple}}$) | $9.46 \pm 0.11$ | **3.17** | $\leq 3.75$ (3.72) |

# Major drawbacks of deep VAE / flow

- Space-time complexity grows with depth

  → Simple choice of $q(z_{t-1}|z_t)$ which allows for simulation-free training since we can directly sample from $q(z_{t-1}|x)$

- Arbitrary dynamics

  → Fix $q(z_{t-1}|z_t)$ during training, and only train the decoder to revert the noise-injecting dynamics

# Backup slides

# Connection to denoising score matching

$$\nabla_\theta \mathbb{E}_{z_t} \left[ \left\| s_\theta(z_t, t) - \nabla_{z_t} \log q(z_t) \right\|^2 \right]$$

$$= \nabla_\theta \mathbb{E}_{z_t, x} \left[ \left\| s_\theta(z_t, t) - \nabla_{z_t} \log q(z_t | x) \right\|^2 \right]$$

$$\propto \nabla_\theta \mathbb{E}_{\epsilon, x} \left[ \| \epsilon_\theta(z_t, t) - \epsilon \|^2 \right]$$

set $s_\theta(\cdot, t) = -\dfrac{\epsilon_\theta(\cdot, t)}{\sqrt{1 - \overline{\alpha}_t}}$

$$\boxed{z_t = \sqrt{\overline{\alpha}_t} x + \sqrt{1 - \overline{\alpha}_t} \epsilon}$$

---

**Algorithm 2** Sampling

1: $z_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
2: **for** $t = T, \ldots, 1$ **do**
3: $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ if $t > 1$, else $\epsilon = \mathbf{0}$
4: $z_{t-1} = \frac{1}{\sqrt{1-\beta_t}} \left( z_t - \frac{\beta_t}{\sqrt{1-\alpha_t}} \epsilon_\theta(z_t, t) \right) + \sigma_t \epsilon$
5: **end for**
6: **return** $z_0$

---

$$z_t' \leftarrow \frac{1}{\sqrt{1 - \beta_t}} \left( z_t + \sqrt{\beta_t} \sqrt{1 - \overline{\alpha}_t} s_\theta(z_t, t) \right) + \sigma_t \epsilon$$
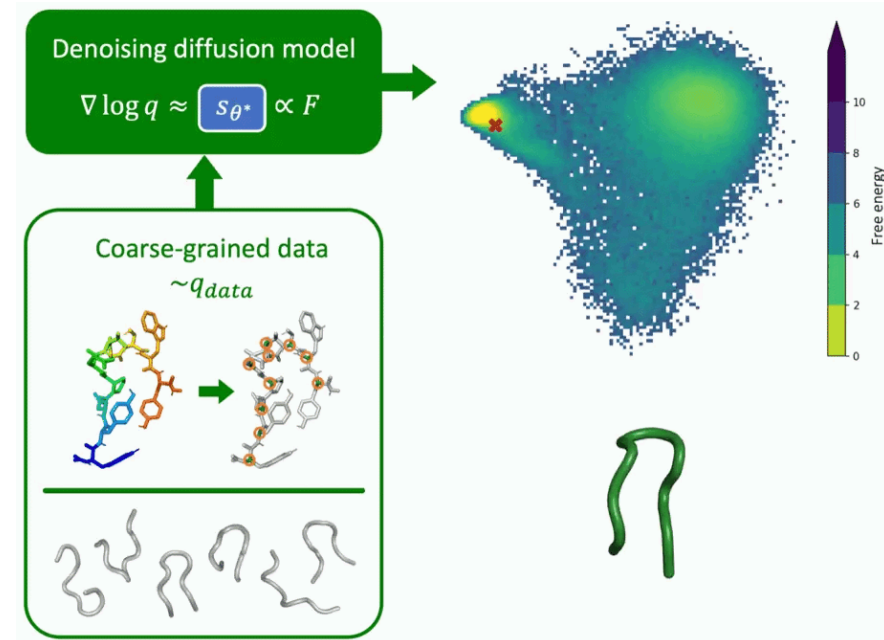
$$z_t' \leftarrow z + \gamma \nabla_z \log q(z_t) + \sqrt{2\gamma} \epsilon \qquad \text{(underdamped Langevin dynamics)}$$

# What can we use the score for

- Energy is composable -> classifier-based / -free guidance
- Accelerated sampling or corrector
- Extracted forces can be used for molecular dynamics



**Two for One: Diffusion Models and Force Fields for Coarse-Grained Molecular Dynamics**

Marloes Arts,[*,†,‡,@] Victor Garcia Satorras,[*,¶,@] Chin-Wei Huang,[¶] Daniel Zügner,[§] Marco Federici,[†,‖] Cecilia Clementi,[§,⊥] Frank Noé,[§] Robert Pinsler,[#] and Rianne van den Berg[¶]
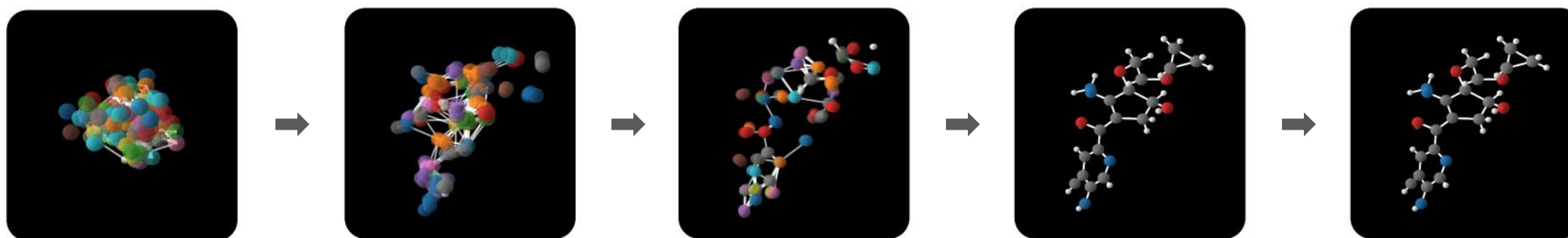
# Outline

1. Introduction to Denoising Diffusion Models
2. **Equivariant Diffusion Models for Molecule Generation in 3D**
3. Workshop: Code & Practice

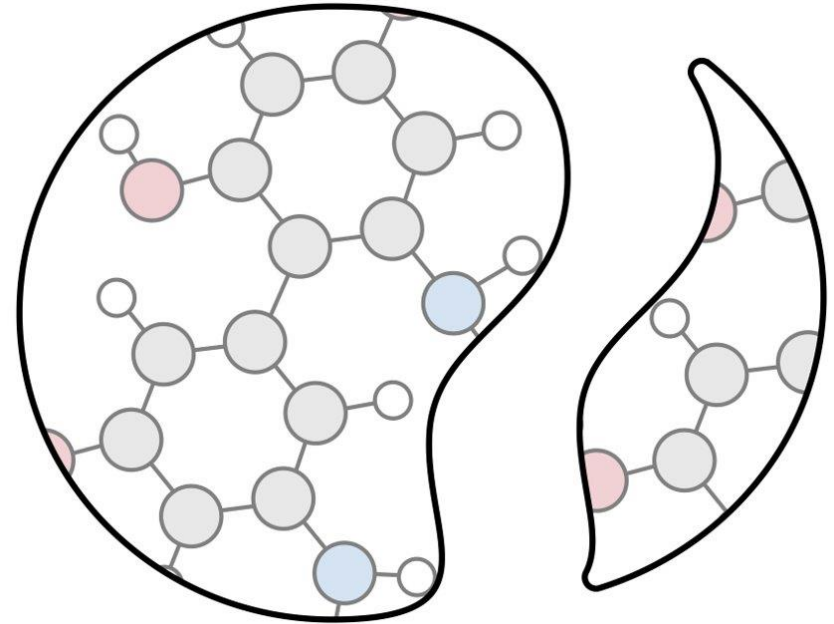# Equivariant Diffusion for Molecule Generation in 3D

Emiel Hoogeboom [*1]   Victor Garcia Satorras [*1]   Clément Vignac [*2]   Max Welling [1]

# Why Generate Molecules in 3D?

- Drug Discovery

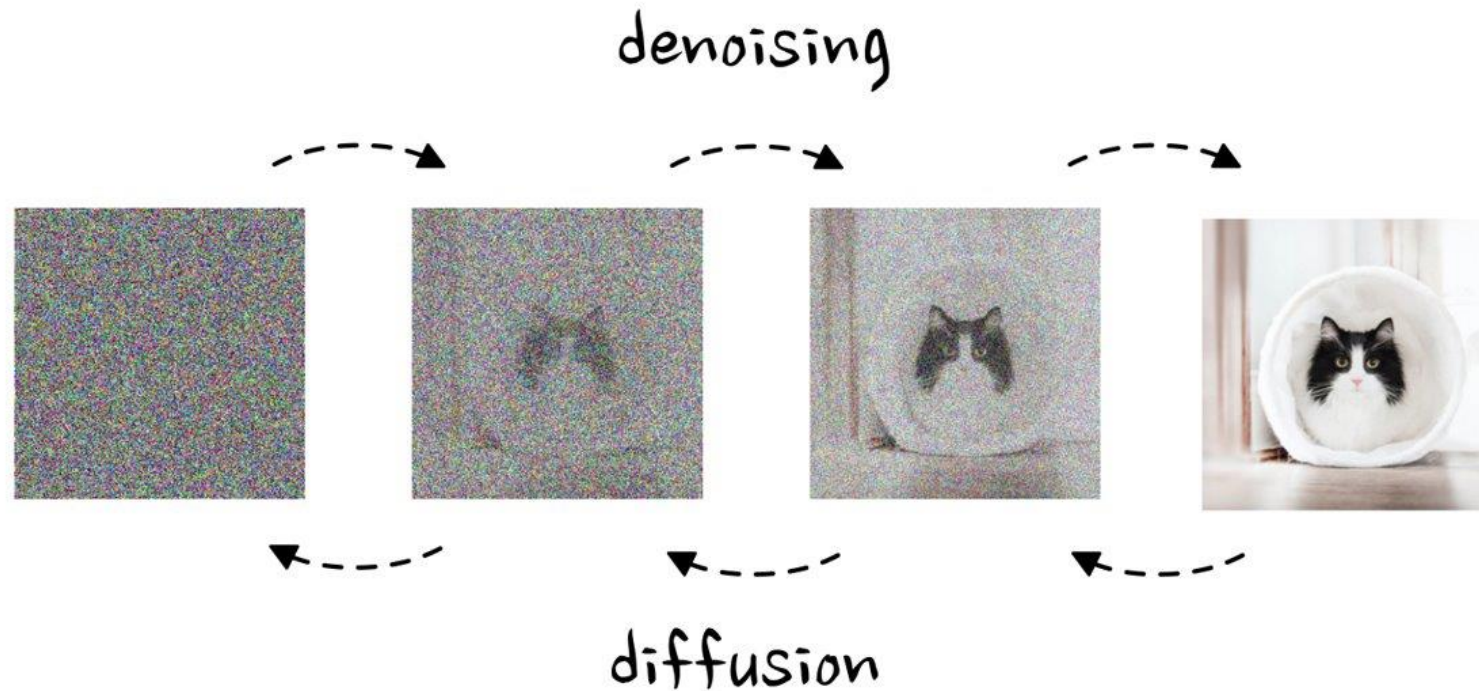- Catalyst Design

- Material Discovery

- Docking Problems

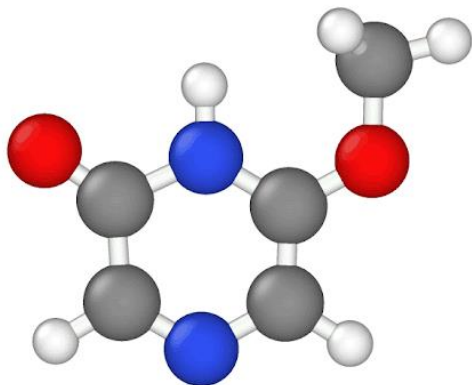# Background: Denoising Diffusion Process

Diffuse: Destroy signal towards Normal Gaussian.

Denoise: Learn a denoising process to generate.

# Background: Equivariance

Invariance / Equivariance



$$f(\boldsymbol{x}) = f(\mathbf{R}\boldsymbol{x})$$

# Background: Equivariance

Invariance / Equivariance



$$\mathbf{R}f(x) = f(\mathbf{R}x)$$

# Background: Equivariance



INPUT
32x32

C1: feature maps
6@28x28

C3: f. maps 16@10x10

S2: f. maps
6@14x14

S4: f. maps 16@5x5

C5: layer
120

F6: layer
84

OUTPUT
10

Convolutions        Subsampling        Convolutions        Subsampling        Full connection        Gaussian connections

Full connection

CNN        Label: 0

# Background: Equivariance

# Background: Equivariance

In molecular modelling we are interested in equivariance w.r.t.:

- Translations 3D

- Rotations 3D (possibly reflections)

- Permutations

# Background: Equivariance

In molecular modelling we are interested in equivariance w.r.t.:

- Translations 3D

- Rotations 3D (possibly reflections)

E(3) Equivariance

- Permutations

# Background: Equivariance

In molecular modelling we are interested in equivariance w.r.t.:

- Translations 3D

- Rotations 3D (possibly reflections)
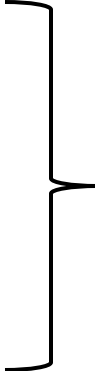
E(3) Equivariance

- Permutations -> Graph Neural Networks / Transformers

# Background: Equivariance



E(n) Equivariant Graph Neural Networks

Victor Garcia Satorras[1]   Emiel Hoogeboom[1]   Max Welling[1]

# EDM: Equivariant Diffusion Models

- Diffusion process to destroy info

- Learn denoising process to generate

- Handles continuous and discrete data



$q(z_T, \ldots | x, h)$ diffuse    $p(x, h, \ldots | z_T)$ denoise

$q(z_T, \ldots | x, h)$ diffuse    $p(x, h, \ldots | z_T)$ denoise

$p(x, h)$    $=$    $p(\mathbf{R}x, h)$

# Diffusion Process

- Adds Gaussian noise over time steps t = 0, ..., T

- Is equivariant to rotations and translations

# Diffusion Process

- Jump to time step "t"

$$q(z_t|x, h) = \mathcal{N}_{xh}(z_t|\sqrt{\bar{\alpha_t}}[x, h], \sigma_t^2 \mathbf{I})$$

To generate molecules
**Learn the *reverse* denoising process**

# Learnable Denoising Process

- Denoise with distributions.

$$p\left(\boldsymbol{z}_{t-1} \mid \boldsymbol{z}_t\right) \approx \mathcal{N}\left(\boldsymbol{z}_{t-1} \mid \hat{\boldsymbol{\mu}}_{t \to t-1}\left(\boldsymbol{z}_t\right), \sigma^2_{t \to t-1}\mathbf{I}\right)$$

- Chosen to have same form as the true denoising process.

# Equivariant Diffusion Models

# Optimization



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

$$q(\boldsymbol{z}_t | \boldsymbol{x}, \boldsymbol{h})$$

$$z_T \quad \bullet \bullet \bullet \quad z_t \quad \xrightarrow{} \quad z_{t-1} \quad \bullet \bullet \bullet \quad z_0 \quad \xrightarrow{} \quad \boldsymbol{x}, \boldsymbol{h}$$

$$p(\boldsymbol{z}_{t-1} | \boldsymbol{z}_t) \qquad\qquad p(\boldsymbol{x}, \boldsymbol{h} | \boldsymbol{z}_0)$$

# Optimization

# Optimization



$$q(\boldsymbol{z}_t | \boldsymbol{x}, \boldsymbol{h})$$

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

$$\boldsymbol{z}_T \quad \cdots \quad \boldsymbol{z}_t \quad \boldsymbol{z}_{t-1} \quad \cdots \quad \boldsymbol{z}_0 \quad \boldsymbol{x}, \boldsymbol{h}$$

$$p(\boldsymbol{z}_{t-1} | \boldsymbol{z}_t) \qquad p(\boldsymbol{x}, \boldsymbol{h} | \boldsymbol{z}_0)$$

$$\boldsymbol{\epsilon}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$
$$\boldsymbol{z}_t = \alpha_t [\boldsymbol{x}, \boldsymbol{h}] + \sigma_t \boldsymbol{\epsilon}_t$$

# Optimization



$$q(\boldsymbol{z}_t|\boldsymbol{x}, \boldsymbol{h})$$

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

$$\boldsymbol{z}_T \quad \cdots \quad \boldsymbol{z}_t \quad \xrightarrow{p(\boldsymbol{z}_{t-1}|\boldsymbol{z}_t)} \quad \boldsymbol{z}_{t-1} \quad \cdots \quad \boldsymbol{z}_0 \quad \xrightarrow{p(\boldsymbol{x}, \boldsymbol{h}|\boldsymbol{z}_0)} \quad \boldsymbol{x}, \boldsymbol{h}$$

$$\min ||\boldsymbol{\epsilon}_t - \phi(\boldsymbol{z}_t, t)||^2$$

diffusion noise    neural net

$$\boldsymbol{\epsilon}_t \sim \mathcal{N}(0, 1)$$

$$\boldsymbol{z}_t = \sqrt{\overline{\alpha_t}}[\boldsymbol{x}, \boldsymbol{h}] + \sqrt{1 - \overline{\alpha_t}}\boldsymbol{\epsilon}_t$$

# Optimization



$q(\boldsymbol{z}_t | \boldsymbol{x}, \boldsymbol{h})$

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

$\boldsymbol{z}_T$ $\cdots$ $\boldsymbol{z}_t$ $\boldsymbol{z}_{t-1}$ $\cdots$ $\boldsymbol{z}_0$ $\boldsymbol{x}, \boldsymbol{h}$

$p(\boldsymbol{z}_{t-1} | \boldsymbol{z}_t)$ $p(\boldsymbol{x}, \boldsymbol{h} | \boldsymbol{z}_0)$

$\min ||\boldsymbol{\epsilon}_t - \boxed{\phi(\boldsymbol{z}_t, t)}||^2$

diffusion noise    neural net

EGNN

$\boldsymbol{\epsilon}_t \sim \mathcal{N}(0,1)$

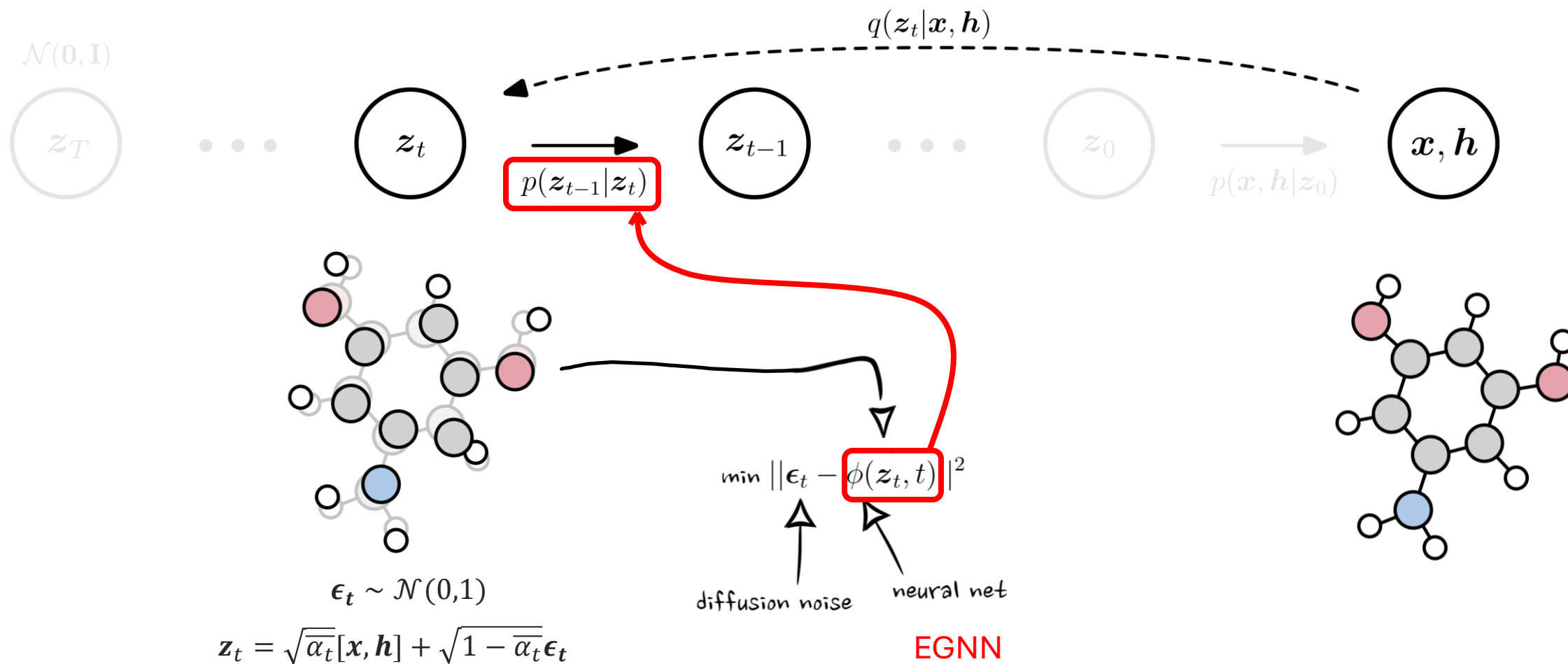$\boldsymbol{z}_t = \sqrt{\overline{\alpha_t}}[\boldsymbol{x}, \boldsymbol{h}] + \sqrt{1 - \overline{\alpha_t}}\boldsymbol{\epsilon_t}$
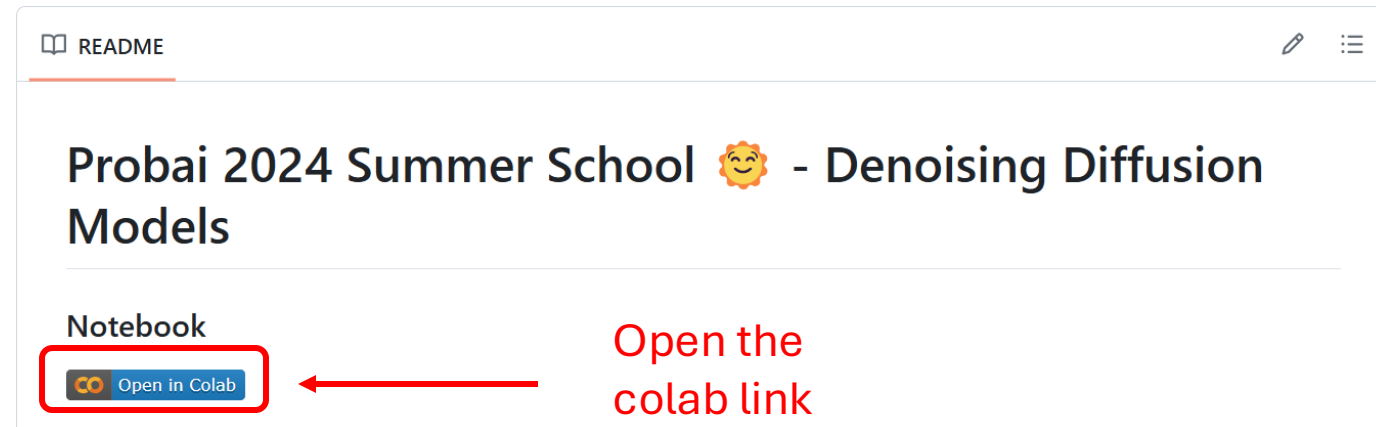
# Optimization

# Outline

1. Introduction to Denoising Diffusion Models
2. Equivariant Diffusion Models for Molecule Generation in 3D
3. **Workshop: Code & Practice**

# Workshop: Code & Practice

1. Open github repository: https://github.com/vgsatorras/probai24

2. Open the Collab:



**Open the colab link**

3. Copy the colab file to your personal drive:



**Press here**