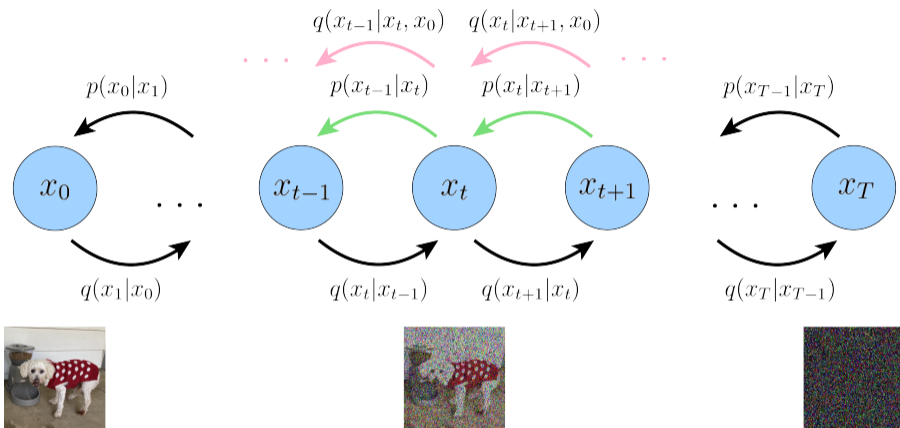


Diffusion Models

Bao Wang
Department of Mathematics
Scientific Computing and Imaging Institute
University of Utah

Diffusion Models



- **Forward diffusion process:** $q(x_t|x_{t-1}), \forall t = 1, 2, \dots, T$.
- **Reverse diffusion process:** $q(x_{t-1}|x_t, x_0), \forall t = T, T-1, \dots, 1$.
- **Diffusion model learns the reverse process $p(x_{t-1}|x_t)$ to match $q(x_{t-1}|x_t, x_0)$ for training data x_0 and $t = T, \dots, 1$. The learned reverse process can be applied to new sampled x_T to generate new data.**

Forward Diffusion Process

- Given a sample from a data distribution $\mathbf{x}_0 \sim q(\mathbf{x})$, we define a **forward diffusion process** by adding a small amount of Gaussian noise in T steps (with variance schedule $\{\beta_t \in (0, 1)\}_{t=1}^T$):

$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I}) \quad q(\mathbf{x}_{1:T} | \mathbf{x}_0) = \prod_{t=1}^T q(\mathbf{x}_t | \mathbf{x}_{t-1})$$

- Let $\alpha_t = 1 - \beta_t$ and $\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$, we have

$$q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I})$$

Usually, $\beta_1 < \beta_2 < \dots < \beta_T$. In DDPM, $\beta_1 = 10^{-4}$ and $\beta_T = 0.02$. When $T \rightarrow \infty$, \mathbf{x}_T is equivalent to an isotropic Gaussian distribution.

Reverse Diffusion Process

- If we can reverse the above process and sample from $q(\mathbf{x}_{t-1}|\mathbf{x}_t)$, we will be able to recreate the true sample from a Gaussian noise input, $\mathbf{x}_T \sim \mathcal{N}(0, \mathbf{I})$.
- We cannot estimate $q(\mathbf{x}_{t-1}|\mathbf{x}_t)$ because it needs to use the entire dataset, therefore, we need to learn a model p_θ to approximate these conditional probabilities to run the **reverse diffusion process**.

$$p_\theta(\mathbf{x}_{0:T}) = p(\mathbf{x}_T) \prod_{t=1}^T p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) \quad p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t))$$

- The reverse conditional probability is tractable when conditioned on \mathbf{x}_0 . Use Bayes' rule, we have

$$\begin{aligned} q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) &= q(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0) \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_0)}{q(\mathbf{x}_t|\mathbf{x}_0)} \\ &\propto \exp\left(-\frac{1}{2}\left(\left(\frac{\alpha_t}{\beta_t} + \frac{1}{1-\bar{\alpha}_{t-1}}\right)\mathbf{x}_{t-1}^2 - \left(\frac{2\sqrt{\alpha_t}}{\beta_t}\mathbf{x}_t + \frac{2\sqrt{\bar{\alpha}_{t-1}}}{1-\bar{\alpha}_{t-1}}\mathbf{x}_0\right)\mathbf{x}_{t-1} + C(\mathbf{x}_t, \mathbf{x}_0)\right)\right) \\ &:= \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_q(\mathbf{x}_t, \mathbf{x}_0), \tilde{\boldsymbol{\beta}}_t \mathbf{I}) \end{aligned}$$

where $C(\mathbf{x}_t, \mathbf{x}_0)$ is some function not involving \mathbf{x}_{t-1} and the details are omitted.

Reverse Diffusion Process

- The mean and variance can be parameterized as ($\alpha_t = 1 - \beta_t$ and $\bar{\alpha}_t = \prod_{i=1}^T \alpha_i$):

$$\begin{aligned}\tilde{\beta}_t &= 1 / \left(\frac{\alpha_t}{\beta_t} + \frac{1}{1 - \bar{\alpha}_{t-1}} \right) = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \cdot \beta_t \\ \tilde{\mu}_q(x_t, x_0) &= \left(\frac{\sqrt{\alpha_t}}{\beta_t} x_t + \frac{\sqrt{\bar{\alpha}_{t-1}}}{1 - \bar{\alpha}_{t-1}} x_0 \right) / \left(\frac{\alpha_t}{\beta_t} + \frac{1}{1 - \bar{\alpha}_{t-1}} \right) \\ &= \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} x_t + \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1 - \bar{\alpha}_t} x_0\end{aligned}$$

- Note that $q(x_t|x_0) = \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)I)$, we can represent then $x_0 = \frac{1}{\sqrt{\bar{\alpha}_t}}(x_t - \sqrt{1 - \bar{\alpha}_t}\epsilon_t)$; thus

$$\mu_q = \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} x_t + \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1 - \bar{\alpha}_t} \frac{1}{\sqrt{\bar{\alpha}_t}} (x_t - \sqrt{1 - \bar{\alpha}_t}\epsilon_t) = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_t \right)$$

Training: Maximize the Evidence Lower Bound (ELBO)

-

$$\log p_{\theta}(\mathbf{x}) = \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[\log \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q(\mathbf{z}|\mathbf{x})} \right] + D_{\text{KL}}(q(\mathbf{z}|\mathbf{x}) || p_{\theta}(\mathbf{z}|\mathbf{x})) \geq \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[\log \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q(\mathbf{z}|\mathbf{x})} \right]$$

- For the multi-level case, we have

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

We need to minimize L_{VLB} .

- Notice that we have

$$L_{\text{VLB}} = \mathbb{E}_q \left[\underbrace{D_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_0) || p_{\theta}(\mathbf{x}_T))}_{L_T} + \sum_{t=2}^T \underbrace{D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) || p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t))}_{L_{t-1}} - \underbrace{\log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)}_{L_0} \right]$$

- Every KL term in L_{VLB} (except for L_0) compares two Gaussian distributions and can be computed in closed form. L_T is constant and can be ignored during training because q has no learnable parameters and \mathbf{x}_T is a Gaussian noise. [May not be true for finite T]

Parameterization of L_t for Training Loss

- We need to learn a neural network to approximate the conditioned probability distributions in the reverse diffusion process, $p_\theta(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \Sigma_\theta(x_t, t))$.
- We train μ_θ to predict $\mu_q = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \epsilon_t \right)$. Because x_t is available as input at training time, we can reparameterize the Gaussian noise term instead to predict ϵ_t from the input x_t at time step t :

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

- The loss term L_t is parameterized to minimize the KL divergence (closed form of KL divergence):

$$\begin{aligned} L_t &= \mathbb{E}_{x_0, \epsilon} \left[\frac{1}{2 \|\Sigma_\theta(x_t, t)\|_2^2} \|\mu_q(x_t, x_0) - \mu_\theta(x_t, t)\|^2 \right] \\ &= \mathbb{E}_{x_0, \epsilon} \left[\frac{(1-\alpha_t)^2}{2\alpha_t(1-\bar{\alpha}_t) \|\Sigma_\theta\|_2^2} \|\epsilon_t - \epsilon_\theta(\sqrt{\bar{\alpha}_t}x_0 + \sqrt{1-\bar{\alpha}_t}\epsilon_t, t)\|^2 \right] \end{aligned}$$

- DDPM chooses to fix β_t as constants instead of making them learnable and set $\Sigma_\theta(x_t, t) = \sigma_t^2 I$, where σ_t is not learned but set to β_t or $\tilde{\beta}_t = \frac{1-\bar{\alpha}_t-1}{1-\bar{\alpha}_t} \cdot \beta_t$.

Training and Sampling DDPM

- Empirically, training DDPM works better with a simplified objective that ignores the weighting term:

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

Algorithm 1 Training

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

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: $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$
 - 5: **end for**
 - 6: **return** \mathbf{x}_0
-

Latent Diffusion Model

- The latent diffusion model (LDM) runs the diffusion process in the latent space instead of pixel space, making training cost lower and inference faster.
- An encoder \mathcal{E} is used to compress the input image $\mathbf{x} \in \mathbb{R}^{H \times W \times 3}$ to a smaller 2D latent vector $\mathbf{z} = \mathcal{E}(\mathbf{x}) \in \mathbb{R}^{h \times w \times c}$. Then an decoder \mathcal{D} reconstructs the images from the latent vector, $\tilde{\mathbf{x}} = \mathcal{D}(\mathbf{z})$.
- The diffusion and denoising processes happen on the latent vector \mathbf{z} . The denoising model is a time-conditioned U-Net, augmented with the cross-attention mechanism to handle flexible conditioning information for image generation (e.g. class labels).

L_t : A Score Viewpoint

- Tweedie's formula: $\mathbb{E}[\boldsymbol{\mu}_q|\mathbf{x}_t] = \mathbf{x}_t + (1 - \bar{\alpha}_t)\nabla \log p(\mathbf{x}_t)$, where $\nabla \log p(\mathbf{x}_t) := \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)$.
- We have shown that $q(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1 - \bar{\alpha}_t)I)$. Then, Tweedie's formula implies that

$$\sqrt{\bar{\alpha}_t}\mathbf{x}_0 = \mathbf{x}_t + (1 - \bar{\alpha}_t)\nabla \log p(\mathbf{x}_t) \Rightarrow \mathbf{x}_0 = \frac{\mathbf{x}_t + (1 - \bar{\alpha}_t)\nabla \log p(\mathbf{x}_t)}{\sqrt{\bar{\alpha}_t}}.$$

- The ground-truth denoising transition mean $\boldsymbol{\mu}_q(\mathbf{x}_t, \mathbf{x}_0)$ again and derive a new form:

$$\boldsymbol{\mu}_q(\mathbf{x}_t, \mathbf{x}_0) = \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\mathbf{x}_0}{1 - \bar{\alpha}_t} = \frac{1}{\sqrt{\alpha_t}}\mathbf{x}_t + \frac{1 - \alpha_t}{\sqrt{\alpha_t}}\nabla \log p(\mathbf{x}_t)$$

- Therefore, we can set our approximate denoising transition mean $\boldsymbol{\mu}_\theta(\mathbf{x}_t, t)$ as follows:

$$\boldsymbol{\mu}_\theta(\mathbf{x}_t, t) = \frac{1}{\sqrt{\alpha_t}}\mathbf{x}_t + \frac{1 - \alpha_t}{\sqrt{\alpha_t}}\mathbf{s}_\theta(\mathbf{x}_t, t).$$

- Then L_t can be written as follows ([score matching](#)):

$$L_t = D_{KL}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)||p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)) = \frac{1}{2\|\Sigma_\theta\|_2^2} \frac{(1 - \alpha_t)^2}{\alpha_t} \left[\|\mathbf{s}_\theta(\mathbf{x}_t, t) - \nabla \log p(\mathbf{x}_t)\|_2^2 \right]$$

Here, $\mathbf{s}_\theta(\mathbf{x}_t, t)$ is a neural network that learns to predict the score function $\nabla \log p(\mathbf{x}_t)$.

Some Score Matching Algorithms

- **Explicit score matching:** learns θ so that $\psi(\mathbf{x}; \theta)$ best matches the corresponding score of the true distribution, i.e. $\frac{\partial \log p(\mathbf{x})}{\partial \mathbf{x}}$. The corresponding objective function is

$$J_{ESM_p}(\theta) = \mathbb{E}_{p(\mathbf{x})} \left[\frac{1}{2} \left\| \psi(\mathbf{x}; \theta) - \frac{\partial \log p(\mathbf{x})}{\partial \mathbf{x}} \right\|^2 \right],$$

since p is unknown, we do not have explicit regression targets $\frac{\partial \log p(\mathbf{x})}{\partial \mathbf{x}}$.

- **Implicit score matching:** Notice that

$$\underbrace{\mathbb{E}_{p(\mathbf{x})} \left[\frac{1}{2} \left\| \psi(\mathbf{x}; \theta) - \frac{\partial \log p(\mathbf{x})}{\partial \mathbf{x}} \right\|^2 \right]}_{J_{ESM_p}(\theta)} = \underbrace{\mathbb{E}_{p(\mathbf{x})} \left[\frac{1}{2} \|\psi(\mathbf{x}; \theta)\|^2 + \sum_{i=1}^d \frac{\partial \psi_i(\mathbf{x}; \theta)}{\partial \mathbf{x}_i} \right]}_{J_{ISM_p}(\theta)} + C_1,$$

where $\psi_i(\mathbf{x}; \theta) = \psi(\mathbf{x}; \theta)_i$ and C_1 is a constant that does not depend on θ . This yields an implicit score matching objective J_{ISM_p} that no longer requires having an explicit score target for p .

- **Denoising score matching:**

> Given i.i.d. samples $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}\} \sim p(\mathbf{x})$, the Parzen window density estimator is

$$p_\sigma(\tilde{\mathbf{x}}) = \frac{1}{n} \sum_{t=1}^n p_\sigma(\tilde{\mathbf{x}}|\mathbf{x}^{(t)}) = \frac{1}{n} \sum_{t=1}^n \frac{1}{(2\pi)^{d/2}\sigma^d} e^{-\frac{1}{2\sigma^2}\|\tilde{\mathbf{x}}-\mathbf{x}^{(t)}\|^2}$$

> Explicitly matching $\psi(\mathbf{x}; \theta)$ with the score of Parzen windows density estimator gives

$$J_{ESM_{p_\sigma}}(\theta) = \mathbb{E}_{p_\sigma(\tilde{\mathbf{x}})} \left[\frac{1}{2} \left\| \psi(\tilde{\mathbf{x}}; \theta) - \frac{\partial \log p_\sigma(\tilde{\mathbf{x}})}{\partial \mathbf{x}} \right\|^2 \right].$$

> Consider pairs of clean and corrupted data $(\mathbf{x}, \tilde{\mathbf{x}}) \sim p_\sigma(\tilde{\mathbf{x}}, \mathbf{x}) = p_\sigma(\tilde{\mathbf{x}}|\mathbf{x})p_0(\mathbf{x})$, we define the following denoising score matching (DSM) objective:

$$J_{DSM_{p_\sigma}}(\theta) = \mathbb{E}_{p_\sigma(\mathbf{x}, \tilde{\mathbf{x}})} \left[\frac{1}{2} \left\| \psi(\tilde{\mathbf{x}}; \theta) - \frac{\partial \log p_\sigma(\tilde{\mathbf{x}}|\mathbf{x})}{\partial \tilde{\mathbf{x}}} \right\|^2 \right].$$

With Gaussian kernel, we have $\frac{\partial \log p_\sigma(\tilde{\mathbf{x}}|\mathbf{x})}{\partial \tilde{\mathbf{x}}} = \frac{1}{\sigma^2}(\mathbf{x} - \tilde{\mathbf{x}})$.

- ESM and ISM directly samples from the true distribution. In contrast, **DSM approximates the Parzen window density estimator from the samples – Gaussian mixture.**

Yuhao Huang, Qingsong Wang, Akwum Onwunta, Bao Wang, Efficient Score Matching via Deep Equilibrium Layers, ICLR, 2024.

Score-based Generative Models

Score-based Generative Models

- We choose a sequence of noise levels $\{\sigma_t\}_{t=1}^T$ and define a sequence of perturbed data distributions:

$$p_{\sigma_t}(\mathbf{x}_t) = \int p(\mathbf{x}) \mathcal{N}(\mathbf{x}_t; \mathbf{x}, \sigma_t^2 \mathbf{I}) d\mathbf{x}.$$

Notice that none of $p_{\sigma_t}(\mathbf{x}_t)$ is a Gaussian mixture, and DSM can be worse than SSM.

- Then, a neural network $\mathbf{s}_\theta(\mathbf{x}, t)$ is used to learn the score function for all noise levels simultaneously:

$$\arg \min_{\theta} \sum_{t=1}^T \lambda(t) \mathbb{E}_{p_{\sigma_t}(\mathbf{x}_t)} \left[\|\mathbf{s}_\theta(\mathbf{x}, t) - \nabla \log p_{\sigma_t}(\mathbf{x}_t)\|_2^2 \right],$$

where $\lambda(t)$ is a positive weighting function that conditions on noise level t .

- Furthermore, we use annealed Langevin dynamics sampling as a generative procedure, in which samples are produced by running Langevin dynamics for each $t = T, \dots, 2, 1$ in sequence.
- The initialization is chosen from a fixed prior, e.g. uniform, and each subsequent sampling step starts from the final samples of the previous simulation. The noise levels steadily decrease over timesteps t , and we reduce the step size over time. [\[Some related theoretical analysis.\]](#)

Score SDE Formulation

- **Forward SDE:** The diffusion process can be viewed as the following SDE:

$$dx = f(x, t)dt + g(t)d\mathbf{w}, \quad t \in [0, T],$$

where $p_0(\mathbf{x}) = p(\mathbf{x})$ is the data distribution. After perturbing $p(\mathbf{x})$ for a sufficiently long time T , $p_T(\mathbf{x})$ becomes close to a tractable noise distribution $\pi(\mathbf{x})$, called a **prior distribution**.

- **Reverse SDE:** The sampling is done via a corresponding reverse-time SDE of the forward process:

$$dx = [f(x, t) - g(t)^2 \nabla_x \log p_t(x)]d\bar{t} + g(t)d\bar{\mathbf{w}}, \quad t \in [0, T],$$

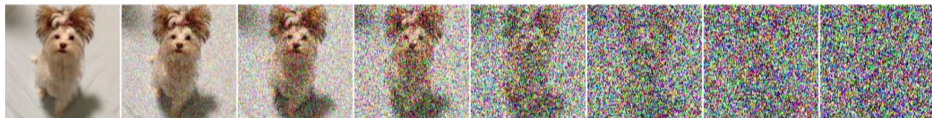
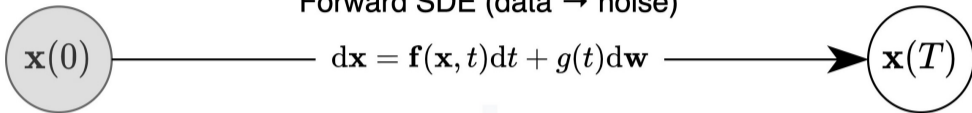
where \bar{t} denotes time traveling backward from T to 0 and $\nabla_x \log p_t(x)$ is the score function.

- The training objective is a weighted sum of score-matching objective:

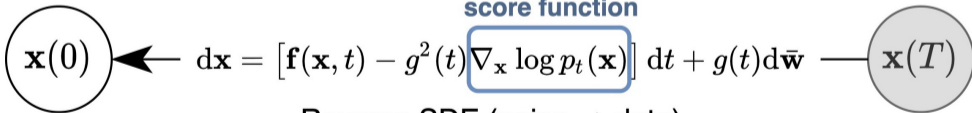
$$L := \mathbb{E}_t \left\{ \lambda(t) \mathbb{E}_{\mathbf{x}_0} \mathbb{E}_{q(\mathbf{x}_t | \mathbf{x}_0)} \left[\|\mathbf{s}_\theta(\mathbf{x}, t) - \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t | \mathbf{x}_0)\|_2^2 \right] \right\}$$

where \mathbf{x}_0 is sampled from the data distribution p_0 and $\lambda(t)$ is sampled from the data distribution p_0 and $\lambda(t)$ is the positive weighting function. $q(\mathbf{x}_t | \mathbf{x}_0)$ is the Gaussian transition kernel associated with the forward process. For example, $q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \mathbf{x}_0, \sigma^2(t)\mathbf{I})$.

Forward SDE (data \rightarrow noise)



score function



Reverse SDE (noise \rightarrow data)

Probability flow ODE

- **Probability flow ODE (PFO):** PFO is the continuous-time ODE that supports the deterministic process which shares the same marginal probability density with SDE. The corresponding PFO is

$$d\mathbf{x} = \left[f(\mathbf{x}, t) - \frac{1}{2}g(t)^2 \nabla_{\mathbf{x}} \log p_t(\mathbf{x}) \right] dt.$$

PFO can be solved with larger step sizes as they have no randomness.

- **Problem:** Enforcing guidance in the SDE design for particular applications.

Reference:

<https://arxiv.org/abs/2209.15408>

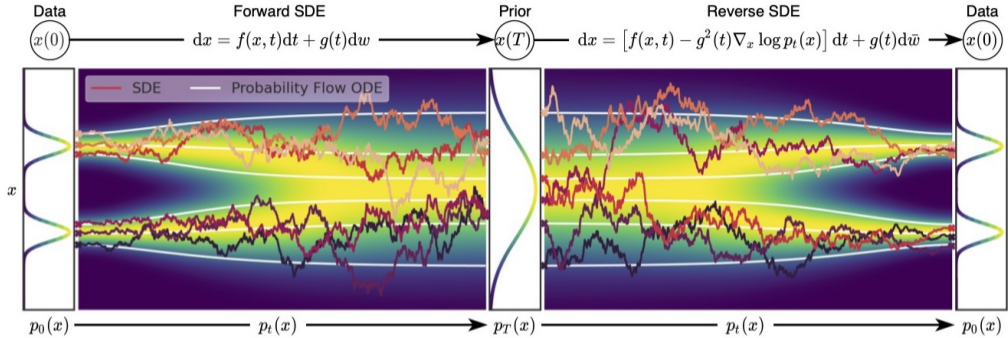


Figure: We map data to a noise distribution (prior) with an SDE, and reverse this SDE for generative modeling. We can also reverse the associated probability flow ODE, yielding a deterministic process that samples from the same distribution as the SDE.

Guided Diffusion Models

Guidance

- So far, we have focused on modeling just the data distribution $p(\mathbf{x})$. However, we are often also interested in learning conditional distribution $p(\mathbf{x}|y)$, enabling us to control the data we generated through conditioning information y . [This is the backbone of DALL-E2 and Imagen.](#)
- We can add arbitrary conditioning information y at each transition step as follows:

$$p(\mathbf{x}_{0:T}|y) = p(\mathbf{x}_T) \prod_{t=1}^T p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t, y).$$

For example, y could be a low-resolution image to perform super-resolution on or physical properties in material sciences and biophysics.

- The two most popular forms of guidance are [Classifier Guidance](#) and [Classifier-Free Guidance](#).

Classifier Guidance

- Consider learning $\nabla \log p(\mathbf{x}_t|y)$, the score of the conditional model, at arbitrary noise levels t .
- By Bayes rule, we have

$$\nabla \log p(\mathbf{x}_t|y) = \nabla \log \left(\frac{p(\mathbf{x}_t)p(y|\mathbf{x}_t)}{p(y)} \right) = \underbrace{\nabla \log p(\mathbf{x}_t)}_{\text{unconditional score}} + \underbrace{\nabla \log p(y|\mathbf{x}_t)}_{\text{adversarial gradient}}. \quad (1)$$

- In Classifier Guidance, the score of an unconditional diffusion model is learned as previously derived, alongside a classifier that takes in arbitrary noisy \mathbf{x}_t and attempts to predict conditional information y . In sampling, the overall conditional score function is used for annealed Langevin dynamics.
- We can rescale the adversarial gradient of the noisy classifier to get

$$\nabla \log p(\mathbf{x}_t|y) = \nabla \log p(\mathbf{x}_t) + \gamma \nabla \log p(y|\mathbf{x}_t). \quad (2)$$

Classifier-Free Guidance

- Train a separate classifier in favor of an unconditional diffusion model and a conditional diffusion model. To derive the score function under Classifier-Free Guidance, we can first rearrange (1) to

$$\nabla \log p(y|x_t) = \nabla \log p(x_t|y) - \nabla \log p(x_t). \quad (3)$$

- Substituting (3) into (2), we get

$$\begin{aligned} \nabla \log p(x_t|y) &= \nabla \log p(x_t) + \gamma(\nabla \log p(x_t|y) - \nabla \log p(x_t)) \\ &= \underbrace{\gamma \nabla \log p(x_t|y)}_{\text{conditional score}} + \underbrace{(1 - \gamma) \nabla \log p(x_t)}_{\text{unconditional score}} \end{aligned} \quad (4)$$

- Learning two separate diffusion models is expensive, we can learn both conditional and unconditional diffusion models together as a singular conditional model: the unconditional diffusion model can be queried by replacing the conditioning information with fixed constant values, such as zeros.

AI for Science

Controllable generation for inverse problem solving

- Accelerate MRI and CT: Achieving even better performance than supervised or unrolled deep learning approaches, while being more robust to different measurement processes at test time.
- Inverse problems are the same as Bayesian inference. Let the forward process of generating \mathbf{y} from \mathbf{x} is given by the transition probability $p(\mathbf{y}|\mathbf{x})$, the inverse problem is to compute $p(\mathbf{x}|\mathbf{y})$.
- From Bayes' rule, we have $p(\mathbf{x}|\mathbf{y}) = p(\mathbf{x})p(\mathbf{y}|\mathbf{x}) / \int p(\mathbf{x})p(\mathbf{y}|\mathbf{x})d\mathbf{x}$. Taking gradient w.r.t. \mathbf{x} on both sides, leading to the following Bayes' rule for score functions

$$\nabla_{\mathbf{x}} \log p(\mathbf{x}|\mathbf{y}) = \nabla_{\mathbf{x}} \log p(\mathbf{x}) + \nabla_{\mathbf{x}} \log p(\mathbf{y}|\mathbf{x}). \quad (5)$$

- Through score matching, we can train a model to estimate the score function of the unconditional data distribution, i.e. $\mathbf{s}_{\theta}(\mathbf{x}) \approx \nabla_{\mathbf{x}} \log p(\mathbf{x})$.
- This will allow us to compute the posterior score function $\nabla_{\mathbf{x}} \log p(\mathbf{x}|\mathbf{y})$ from the known forward process $p(\mathbf{y}|\mathbf{x})$ via (5), and sample from it with Langevin-type sampling.

ConfGF

References:

<http://proceedings.mlr.press/v139/shi21b/shi21b.pdf>

https://proceedings.neurips.cc/paper_files/paper/2021/file/a45a1d12ee0fb7f1f872ab91da18f899-Paper.pdf

- **Molecular Conformation Prediction:** Given an undirected graph $G = (V, E)$ with each node v_i associated with a nuclear charge Z_i and a 3D vector \mathbf{r}_i , indicating its atomic type and atomic coordinate, respectively. All distances between connected nodes $\mathbf{d} = \{d_{ij}\} \in \mathbb{R}^{|E|}$. We extend the original molecular graph by adding auxiliary edges, i.e., virtual bonds, between atoms that are 2 or 3 hops away. Given an extended molecular graph $G = (V, E)$, our goal is to learn a generative model to generate molecular conformations $\mathbf{R} = (\mathbf{r}_1, \mathbf{r}_2, \dots, \mathbf{r}_{|V|})$ based on the molecular graph.
- We aim to model the score function of the p.d.f. w.r.t. conformation $s(\mathbf{R}) = \nabla_{\mathbf{R}} \log p(\mathbf{R}|G)$. As $\log p(\mathbf{R}|G)$ is roto-translation invariant w.r.t. conformations, so is the score function $s(\mathbf{R})$.
- To ensure $s(\mathbf{R})$ is roto-translational equivariant, ConfGF parameterize the log-density of coordinates \mathbf{R} given the molecular graph G as $\log p_{\theta}(\mathbf{R}|G) := f_G \circ g_G(\mathbf{R}) = f_G(\mathbf{d})$, where $g_G : \mathbb{R}^{|V| \times 3} \rightarrow \mathbb{R}^{|E|}$ maps a set of atomic coordinates to a set of interatomic distances, and $f_G : \mathbb{R}^{|E|} \rightarrow \mathbb{R}$ is a GNN that estimates the negative energy of a molecule based on the interatomic distances \mathbf{d} and the graph G .
- Estimating the gradient fields of atomic coordinates can be translated to estimating the gradient fields of interatomic distances, which can be estimated using score-based generative models.

- Note that the interatomic distances \mathbf{d} are continuously differentiable w.r.t. atomic coordinates \mathbf{R} . Therefore, the two gradient fields, $\nabla_{\mathbf{R}} \log p_{\theta}(\mathbf{R}|G)$ and $\nabla_{\mathbf{d}} \log p_{\theta}(\mathbf{d}|G)$, are connected by chain rule:

$$\mathbf{s}_{\theta}(\mathbf{R})_i = \frac{\partial f_G(\mathbf{d})}{\partial \mathbf{r}_i} = \sum_{(i,j), \mathbf{e}_{ij} \in E} \frac{\partial f_G(\mathbf{d})}{\partial d_{ij}} \cdot \frac{\partial d_{ij}}{\partial \mathbf{r}_i} = \sum_{j \in N(i)} \frac{1}{d_{ij}} \cdot \frac{\partial f_G(\mathbf{d})}{\partial d_{ij}} \cdot (\mathbf{r}_i - \mathbf{r}_j) = \sum_{j \in N(i)} \frac{1}{d_{ij}} \cdot \mathbf{s}_{\theta}(\mathbf{d})_{ij} \cdot (\mathbf{r}_i - \mathbf{r}_j),$$

where $\mathbf{s}_{\theta}(\mathbf{R})_i = \nabla_{\mathbf{r}_i} \log p_{\theta}(\mathbf{R}|G)$ and $\mathbf{s}_{\theta}(\mathbf{d})_{ij} = \nabla_{d_{ij}} \log p_{\theta}(\mathbf{d}|G)$.

- By the above observation, we first train a noise conditional score network $\mathbf{s}_{\theta}(\mathbf{d}, \sigma)$ to approximate $\nabla_{\tilde{\mathbf{d}}} \log q_{\sigma}(\tilde{\mathbf{d}}|G)$. We then calculate the gradient fields of atomic coordinates $\mathbf{s}_{\theta}(\mathbf{R}, \sigma)$ from $\mathbf{s}_{\theta}(\mathbf{d}, \sigma)$. Such designed score network $\mathbf{s}_{\theta}(\mathbf{R}, \sigma)$ enjoys roto-translation equivariance.

- Let $\{\sigma_i\}_{i=1}^L$ be a geometric progression with common ratio γ , i.e. $\sigma_i/\sigma_{i-1} = \gamma$. We define the perturbed conditional distance distribution $q_{\sigma}(\tilde{\mathbf{d}}|G)$ to be $\int p(\mathbf{d}|G) \mathcal{N}(\tilde{\mathbf{d}}|\mathbf{d}, \sigma_i^2 \mathbf{I}) d\mathbf{d}$. We aim to learn a conditional score network to jointly estimate the scores of all perturbed distance distributions, i.e. $\forall \sigma \in \{\sigma_i\}_{i=1}^L : \mathbf{s}_{\theta}(\tilde{\mathbf{d}}, \sigma) \approx \nabla_{\tilde{\mathbf{d}}} \log q_{\sigma}(\tilde{\mathbf{d}}|G)$. Note that $\mathbf{s}_{\theta}(\tilde{\mathbf{d}}, \sigma) \in \mathbb{R}^{|E|}$ – **Much higher-dimension**.

- **ConfGF circumvents directly modeling atomic coordinates using intermediate geometric variables such as atomic distances, bond, and torsion angles, which are roto-translational invariant.**

- We embed node and edge attributes into low-dimensional feature spaces using feedforward networks

$$\mathbf{h}_i^0 = MLP(Z_i) \quad \forall v_i \in V; \quad \mathbf{h}_{e_{ij}} = MLP(e_{ij}, d_{ij}), \quad \forall e_{ij} \in E.$$

Then use graph isomorphism network to compute node embeddings based on graph structures:

$$\mathbf{h}_i^l = MLP(\mathbf{h}_i^{l-1} + \sum_{j \in N(i)} ReLU(\mathbf{h}_j^{l+1} + \mathbf{h}_{e_{ij}})).$$

We compute the final edge embeddings by concatenating the corresponding node embeddings:

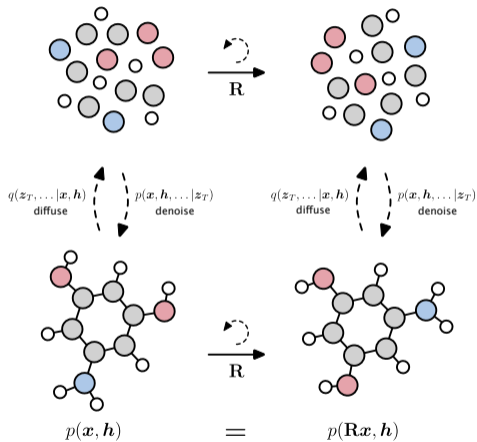
$$\mathbf{h}_{e_{ij}}^o = \mathbf{h}_i^N || \mathbf{h}_j^N || \mathbf{h}_{e_{ij}}.$$

- We parameterize the noise conditional network with $\mathbf{s}_\theta(\tilde{\mathbf{d}}, \sigma) = \mathbf{s}_\theta(\tilde{\mathbf{d}})/\sigma$ as follows

$$\mathbf{s}_\theta(\tilde{\mathbf{d}})_{ij} = MLP(\mathbf{h}_{e_{ij}}^o), \quad \forall e_{ij} \in E,$$

where $\mathbf{s}_\theta(\tilde{\mathbf{d}})$ is an unconditional score network, and the MLP maps edge embeddings to a scalar for each $e_{ij} \in E$. Note that $\mathbf{s}_\theta(\tilde{\mathbf{d}}, \sigma)$ is roto-translation invariant as we only use interatomic distances.

EDM



Reference: <https://arxiv.org/abs/2203.17003>

- EDM learns to denoise a diffusion process with an equivariant network that jointly operates on both continuous (atom coordinates) and categorical features (atom types).

- A conditional distribution $p(\mathbf{y}|\mathbf{x})$ is equivariant to the action of rotations or reflections when

$$p(\mathbf{y}|\mathbf{x}) = p(\mathbf{R}\mathbf{y}|\mathbf{R}\mathbf{x}), \text{ for all orthogonal } \mathbf{R}.$$

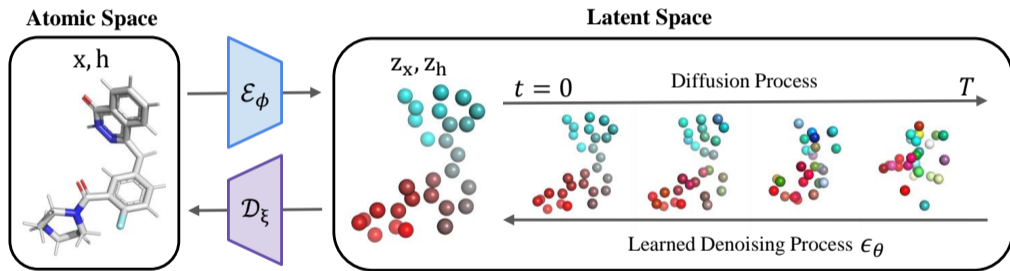
A distribution is invariant to \mathbf{R} transformation if $p(\mathbf{y}) = p(\mathbf{R}\mathbf{y})$, for all orthogonal \mathbf{R} .

- If $\mathbf{x} \sim p(\mathbf{x})$ is invariant to a group and the transition probabilities of a Markov chain $\mathbf{y} \sim p(\mathbf{y}|\mathbf{x})$ are equivariant, then the marginal distribution of \mathbf{y} at any time step is invariant to group transformations. It means that if $p(\mathbf{z}_T)$ is invariant and the neural network $p(\mathbf{x}_{t-1}|\mathbf{x}_t)$ is equivariant, then the marginal distribution $p(\mathbf{x})$ of the denoising model will be invariant.

- Parameterize diffusion model (the noise term) using EGNN.

- **Problem:** How to extend EDM to score-based diffusion models?

GeoLDM: Geometric Latent Diffusion Models for 3D Molecule Generation



Reference: <https://proceedings.mlr.press/v202/xu23n/xu23n.pdf>

- Each molecule is represented as point cloud $\mathcal{G} = (\mathbf{x}, \mathbf{h})$, where $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_N) \in \mathbb{R}^{N \times 3}$ are the atom coordinates and $\mathbf{h} = (\mathbf{h}_1, \dots, \mathbf{h}_N) \in \mathbb{R}^{N \times d}$ are the node features, such as atomic type and charges.

Unconditional generation. With a collection of molecules \mathcal{G} , learn parameterized generative models $p_\theta(\mathcal{G})$ which can generate diverse and realistic molecules $\hat{\mathcal{G}}$ in 3D.

Controllable generation. With molecules \mathcal{G} that has certain properties s , learn conditional generation models $p_\theta(\mathcal{G}|s)$, which can conduct controllable molecule generation given desired property value s .

- The learned likelihood should be invariant to roto-translation.
- **Geometric autoencoder:** The encoder \mathcal{E}_ϕ encodes \mathcal{G} into latent domain $\mathbf{z} = \mathcal{E}_\phi(\mathbf{x}, \mathbf{h})$ and the decoder \mathcal{D}_ξ learns to decode \mathbf{z} back to data domain $\tilde{\mathbf{x}}, \tilde{\mathbf{h}} = \mathcal{D}_\xi(\mathbf{z})$.
- **Parameterization of latent space as invariant scalar-valued features is challenging.** Geometric AE requires an additional function ψ to represent appropriate group actions for encoding, and aligning output and input positions for decoding, to solve the reconstruction task.
- GeoLDM incorporates equivariance into \mathcal{E} and \mathcal{D} by constructing latent features as point-structured variables $\mathbf{z} = (\mathbf{z}_x, \mathbf{z}_h) \in \mathbb{R}^{N \times (3+k)}$, which holds 3D equivariant and k D invariant latent features \mathbf{z}_x and \mathbf{z}_h for each node. GeoLDM implements it by parameterizing \mathcal{E} and \mathcal{D} with EGNN, which extracts both invariant and equivariant embeddings.

- The latent space \mathbf{z} contains not only scalars (invariant features) \mathbf{z}_h but also tensors (equivariant features) \mathbf{z}_x . This requires the distribution of latent DMs to satisfy $p_\theta(\mathbf{z}_x, \mathbf{z}_h) = p_\theta(\mathbf{R}\mathbf{z}_x, \mathbf{z}_h)$, $\forall \mathbf{R}$, and GeoLDM implements latent DMs using EGNN.
- For the number of nodes N , we need to sample different numbers N for generating molecules of different sizes. We use the distribution $p(N)$ of molecular sizes on the training set. Then for generation, we can first sample $N \sim p(N)$ and then generate latent variables and node features in size N .
- DMs are capable of controllable generation with given conditions s , by modeling conditional distributions $p(\mathbf{z}|s)$. This in DMs can be implemented with conditional denoising networks $\epsilon_\theta(\mathbf{z}, t, s)$, with the critical difference that it takes additional input s .

SyMat: Towards Symmetry-Aware Generation of Periodic Materials

Reference: <https://arxiv.org/abs/2307.02707>

- SyMat generates atom types and lattices of materials by generating atom type sets, lattice lengths, and lattice angles with VAE. SyMat further employs a symmetry-aware score-based diffusion model to generate atom coordinates of materials.
- Assume there are n atoms in a unit cell of the periodic material $\mathbf{M} = (\mathbf{A}, \mathbf{P}, \mathbf{L})$, where $\mathbf{A}, \mathbf{P} \in \mathbb{R}^{3 \times n}$, and $\mathbf{L} \in \mathbb{R}^{3 \times 3}$ are the atom type vector, coordinate matrix, and lattice matrix, resp.
- We are given a dataset $\mathcal{M} = \{\mathbf{M}_j\}_{j=1}^m$ where each \mathbf{M}_j is assumed to be sampled from the distribution $p(\cdot)$, which is invariant to the following symmetry transformation:
 - > **Permutation.** If we permute \mathbf{A} to get \mathbf{A}' and use the same order to permute the column vectors of \mathbf{P} to obtain \mathbf{P}' . Let $\mathbf{M}' = (\mathbf{A}', \mathbf{P}', \mathbf{L})$, then $p(\mathbf{M}) = p(\mathbf{M}')$.
 - > **Rotation and translation.** For any $\mathbf{M} = (\mathbf{A}, \mathbf{P}, \mathbf{L})$, we can rotate and translate its atom positions to obtain $\mathbf{M}' = (\mathbf{A}, \mathbf{P}', \mathbf{L}')$. Then $p(\mathbf{M}) = p(\mathbf{M}')$.
 - > **Periodic transformation.** A periodically equivalent coordinate matrix \mathbf{P}^K of \mathbf{P} can be obtained by periodic transformation as $\mathbf{P}^K = \mathbf{P} + \mathbf{L}\mathbf{K}$, where $\mathbf{K} \in \mathbb{Z}^{3 \times n}$. Let $\mathbf{M} = (\mathbf{A}, \mathbf{P}, \mathbf{L})$ and $\mathbf{M}^K = (\mathbf{A}, \mathbf{P}^K, \mathbf{L})$ then $p(\mathbf{M}) = p(\mathbf{M}^K)$.
- We learn a generative model $p_\theta(\cdot)$ from the given dataset \mathcal{M} to capture the real data distribution $p(\cdot)$, and the learned generation model can generate a valid material structure \mathbf{M} using $p_\theta(\mathbf{M})$.

SyMat: Atom Type and Lattice Generation

- Using VAE to directly generate and capture the distribution of \mathbf{A} and \mathbf{L} needs to include permutation invariance of \mathbf{A} and rotation invariance of \mathbf{L} . SyMat transforms \mathbf{A} and \mathbf{L} to items that are invariant to symmetry transformations and makes these items the direct generation targets of the VAE model.

- (1) We count the number of atoms for every type and represent \mathbf{A} with an unordered k -element atom type set $\mathbf{c} = \{(e_1, n_1), \dots, (e_k, n_k)\}$. We see that \mathbf{c} is invariant to any permutation on \mathbf{A} .

- (2) We use six rotation-invariant items for lattice \mathbf{L} generation, including three lattice lengths $\mathbf{l} = [l_1, l_2, l_3]$ and three lattice angles $\phi = [\phi_{12}, \phi_{13}, \phi_{23}]$, where ϕ_{ij} is the angle between l_i and l_j .

- The VAE model used for generating \mathbf{c} , \mathbf{l} , and ϕ consists of an encoder and a decoder. The encoder uses symmetry-aware SphereNet that takes a material $\mathbf{M} = (\mathbf{A}, \mathbf{P}, \mathbf{L})$ as inputs and outputs a latent variable $\mathbf{z}^{\mathbf{A}} \in \mathbb{R}^d$ and a latent variable $\mathbf{z}^{\mathbf{L}} \in \mathbb{R}^f$.

- For any material $\mathbf{M} = (\mathbf{A}, \mathbf{P}, \mathbf{L})$ in \mathcal{M} , the encoder first maps it to latent variable $\mathbf{z}^{\mathbf{A}}$ and $\mathbf{z}^{\mathbf{L}}$. Afterwards, the decoder uses $\mathbf{z}^{\mathbf{A}}$ and $\mathbf{z}^{\mathbf{L}}$ to reconstruct the exact \mathbf{c} , \mathbf{l} , and ϕ obtained from \mathbf{A} and \mathbf{L} .

SyMat: Coordinate Generation

- SyMat generates the coordinate matrix \mathbf{P} conditioned on \mathbf{A} and \mathbf{L} using a score-based diffusion model with score matrix $\nabla_{\mathbf{P}} \log p(\mathbf{P}|\mathbf{A}, \mathbf{L})$.
- $\nabla_{\mathbf{P}} \log p(\mathbf{P}|\mathbf{A}, \mathbf{L})$ satisfies:
 - (1) $\nabla_{\mathbf{P}} \log p(\mathbf{P}|\mathbf{A}, \mathbf{L})$ is equivariant to permutations on \mathbf{A} and \mathbf{P}
 - (2) $\nabla_{\mathbf{QP} + \mathbf{b1}^\top} \log p(\mathbf{QP} + \mathbf{b1}^\top|\mathbf{A}, \mathbf{QL}) = \mathbf{Q} \nabla_{\mathbf{P}} \log p(\mathbf{P}|\mathbf{A}, \mathbf{L})$ for any orthogonal matrix \mathbf{Q} and \mathbf{b}
 - (3) $\nabla_{\mathbf{P} + \mathbf{LK}} \log p(\mathbf{P} + \mathbf{LK}|\mathbf{A}, \mathbf{L}) = \nabla_{\mathbf{P}} \log p(\mathbf{P}|\mathbf{A}, \mathbf{L})$ holds for any $\mathbf{K} \in \mathbb{Z}^{3 \times n}$.
- We formulate $\nabla_{\mathbf{P}} \log p(\mathbf{P}|\mathbf{A}, \mathbf{L})$ as a score function of edge distances in a graph created by a [multi graph method](#). For \mathbf{M} with coordinates $\mathbf{P} = [\mathbf{p}_1, \dots, \mathbf{p}_n]$ and lattice matrix \mathbf{L} , multi graph method produces an n -node undirected graph on \mathbf{M} – the node i corresponds to the i -th atom in the unit cell, and an edge is added between any two nodes i, j if one of the interatomic distances between them is within cutoff r .
- The interatomic distances between i, j are the distances between their corresponding atoms in the complete infinite structure of \mathbf{M} , including both the distance within the unit cell and the distance crossing the unit cell boundary. Formally, the set of all edges in the graph constructed on \mathbf{M} can be written as $\mathbf{E}(\mathbf{M}) = \{(i, j, \mathbf{k}) : \|\mathbf{p}_i + \mathbf{Lk} - \mathbf{p}_j\|_2 \leq r, \mathbf{k} \in \mathbb{Z}^3, 1 \leq i, j \leq n\}$.

SyMat: Coordinate Generation

- Let $d_{i,j,k} = \|\mathbf{p}_i + \mathbf{L}\mathbf{k} - \mathbf{p}_j\|_2$ be the distance of the edge (i, j, \mathbf{k}) , we consider $\log p(\mathbf{P}|\mathbf{A}, \mathbf{L})$ as a function of the distances of all edges in $E(\mathbf{M})$. Denoting s_i as the score function of \mathbf{p}_i , then $\nabla_{\mathbf{P}} \log p(\mathbf{P}|\mathbf{A}, \mathbf{L}) = [\mathbf{s}_1, \dots, \mathbf{s}_n]$. From the chain rule, we can calculate s_i as

$$s_i = \sum_{(j,k) \in \mathcal{N}(i)} \nabla_{d_{i,j,k}} \log p(\mathbf{P}|\mathbf{A}, \mathbf{L}) \cdot \nabla_{\mathbf{p}_i} d_{i,j,k} = \sum_{(j,k) \in \mathcal{N}(i)} s_{i,j,k} \cdot \frac{\mathbf{p}_i + \mathbf{L}\mathbf{k} - \mathbf{p}_j}{d_{i,j,k}},$$

where $s_{i,j,k} = \nabla_{d_{i,j,k}} \log p(\mathbf{P}|\mathbf{A}, \mathbf{L})$ is the score function of the distance $d_{i,j,k}$. $\nabla_{\mathbf{P}} \log p(\mathbf{P}|\mathbf{A}, \mathbf{L})$ can be approximated by first approximating $s_{i,j,k}$ for every edge in the multi-graph representation of \mathbf{M} , then calculating each column vector of $\nabla_{\mathbf{P}} \log p(\mathbf{P}|\mathbf{A}, \mathbf{L})$.

- If $s_{i,j,k}$ is invariant to all symmetry transformations described before for every edge (i, j, \mathbf{k}) , the score matrix $\nabla_{\mathbf{P}} \log p(\mathbf{P}|\mathbf{A}, \mathbf{L})$ satisfies the desired symmetry.
- SyMat uses SphereNet as the score model $s_{\theta}(\cdot)$ to approximate $\nabla_{\mathbf{P}} \log p(\mathbf{P}|\mathbf{A}, \mathbf{L})$. $s_{\theta}(\cdot)$ takes the multi-graph representation of \mathbf{M} as input and outputs $o_{i,j,k}$ as the approximated $s_{i,j,k}$ at every edge (i, j, \mathbf{k}) in the input graph.

DiffCSP: Crystal Structure Prediction by Joint Equivariant Diffusion

Reference: <https://arxiv.org/abs/2309.04475>

- **Representation of crystal structures.** A 3D crystal can be represented as the infinite periodic arrangement of atoms in 3D space, and the smallest repeating unit is called a unit cell. A unit cell can be defined by a triplet $\mathcal{M} = (\mathbf{A}, \mathbf{X}, \mathbf{L})$, where $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_N] \in \mathbb{R}^{h \times N}$ denotes the list of the one-hot representation of atom types, $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N] \in \mathbb{R}^{3 \times N}$ consists of Cartesian coordinates of the atoms, and $\mathbf{L} = [l_1, l_2, l_3] \in \mathbb{R}^{3 \times 3}$ represents the lattice matrix consisting three basic vectors to describe the periodicity of the crystal. The infinite periodic crystal structure is represented by

$$\{(\mathbf{a}'_i, \mathbf{x}'_i) | \mathbf{a}'_i = \mathbf{a}_i, \mathbf{x}'_i = \mathbf{x}_i + \mathbf{L}\mathbf{k}, \forall \mathbf{k} \in \mathbb{Z}^{3 \times 1}\},$$

where the j -th element of the integral vector \mathbf{k} denotes the integral 3D translation in units of l_j .

Fractional coordinate system. The Cartesian coordinate system \mathbf{X} leverages three standard orthogonal bases as the coordinate axes. The fractional coordinate system uses the lattices (l_1, l_2, l_3) as the bases. A point represented by fractional coordinate vector $\mathbf{f} = [f_1, f_2, f_3]^T \in [0, 1)^3$ corresponds to Cartesian vector $\mathbf{x} = \sum_{i=1}^3 f_i l_i$. Using a fractional coordinate system, we denote the crystal by $\mathcal{M} = (\mathbf{A}, \mathbf{F}, \mathbf{L})$, where the fractional coordinates of all atoms in a cell compose the matrix $\mathbf{F} = [0, 1)^{3 \times N}$.

Task. CSP predicts for each unit cell the lattice matrix \mathbf{L} and the fractional matrix \mathbf{F} given its chemical composition \mathbf{A} , namely, learning the conditional distribution $p(\mathbf{L}, \mathbf{F} | \mathbf{A})$.

- Generating crystal structures requires not only modeling the distribution of the atom coordinates within every cell but also inferring how their bases (a.k.a. lattice vectors) are placed in 3D space.
- **Periodic E(3) invariance:** Any E(3) transformation of the crystal coordinates does not change the physical law and thus keeps the crystal distribution invariant. Moreover, crystal structure prediction (CSP) exhibits unique challenges due to the periodicity of the atom arrangement in crystals.
- DiffCSP jointly generates the lattice and atom coordinates for each crystal using a periodic E(3) equivariant diffusion model, to model the crystal geometry.
- DiffCSP leverages the fractional coordinate systems other than the Cartesian system, which encodes periodicity intrinsically. The fractional representation not only allows us to consider Wrapped Normal (WN) distribution to better model the periodicity, but also facilitates the design of the denoising model via the Fourier transformation.
- The denoising process in diffusion models acts like a force field that drives the atom coordinates towards the energy local minimum and thus is able to increase stability.

- We consider three types of symmetries in $p(\mathbf{L}, \mathbf{F}|\mathbf{A})$:
 - > (Permutation invariance) For any permutation $\mathbf{P} \in S_N$, $p(\mathbf{L}, \mathbf{F}|\mathbf{A}) = p(\mathbf{L}, \mathbf{FP}|\mathbf{AP})$, i.e. changing the order of atoms will not change the distribution.
 - > (O(3) invariance) For any orthogonal transformation $\mathbf{Q} \in \mathbb{R}^{3 \times 3}$, $p(\mathbf{QL}, \mathbf{F}|\mathbf{A}) = p(\mathbf{L}, \mathbf{F}|\mathbf{A})$.
 - > (Periodic translation invariance) For $\forall \mathbf{t} \in \mathbb{R}^{3 \times 1}$, $p(\mathbf{L}, w(\mathbf{F} + \mathbf{t}\mathbf{1}^\top)|\mathbf{A}) = p(\mathbf{L}, \mathbf{F}|\mathbf{A})$, where the function $w(\mathbf{F}) = \mathbf{F} - \lceil \mathbf{F} \rceil \in [0, 1)^{3 \times N}$ returns the fractional part of each element in \mathbf{F} .
- The permutation invariance is encapsulated by using GNNs. We focus on periodic E(3) invariance.
- With the help of the fractional system, the periodic E(3) invariance is made tractable by fulfilling O(3) invariance on \mathbf{L} and periodic translation invariance on \mathbf{F} .
- \mathbf{L} is a continuous variable; we use equivariant DDPM to accomplish the generation.

- The domain of fractional coordinates $[0, 1)^{3 \times N}$ induced by periodicity. DDPM is not suited to generate \mathbf{F} – the normal distribution used in DDPM cannot model the cyclical and bounded domain.
- We use a score matching-based framework along with Wrapped Normal (WN) distribution to fit the specificity here. During the forward process, we first sample each column of $\epsilon_{\mathbf{F}} \in \mathbb{R}^{3 \times N}$ from $\mathcal{N}(0, 1)$, and then acquire $\mathbf{F}_t = w(\mathbf{F}_0 + \sigma_t \epsilon_{\mathbf{F}})$. This truncated sampling implies the WN transition:

$$q(\mathbf{F}_t | \mathbf{F}_0) \propto \sum_{\mathbf{z} \in \mathbb{Z}^{3 \times N}} \exp\left(-\frac{\|\mathbf{F}_t - \mathbf{F}_0 + \mathbf{z}\|_{\mathbf{F}}^2}{2\sigma_t^2}\right).$$

This process ensures the probability distribution over $[z, z + 1)^{3 \times N}$ for any integer z to be the same to keep the crystal periodicity. $q(\mathbf{F}_t | \mathbf{F}_0)$ is periodic translation equivariant, and approaches a uniform distribution $U(0, 1)$ if σ_T is sufficiently large.

- The marginal distribution $p(\mathbf{F}_0)$ is periodic translation invariant if $\hat{\epsilon}_{\mathbf{F}}(\mathcal{M}_t, t)$ is periodic translation invariant, i.e. $\hat{\epsilon}_{\mathbf{F}}(\mathbf{L}_t, \mathbf{F}_t, \mathbf{A}, t) = \hat{\epsilon}_{\mathbf{F}}(\mathbf{L}_t, w(\mathbf{F}_t + \mathbf{t} \mathbf{1}^T), \mathbf{A}, t), \forall \mathbf{t} \in \mathbb{R}^3$.
- Training objective: $\mathcal{L}_{\mathbf{F}} = \mathbb{E}_{\mathbf{F}_t \sim q(\mathbf{F}_t | \mathbf{F}_0), t \sim U(1, T)} \left[\lambda_t \|\nabla_{\mathbf{F}_t} \log q(\mathbf{F}_t | \mathbf{F}_0) - \hat{\epsilon}_{\mathbf{F}}(\mathcal{M}_t, t)\|_2^2 \right],$