

# **Generative Models for Inverse Imaging Problems**

From mathematical foundations to physics-driven  
applications

# CONTENTS

- Introduction
- Foundations for inverse imaging problems with generative models
  - VAEs
  - GANs
  - Score-based generative models
- Physics-driven applications
  - Generative modeling for cryo-EM analysis
  - Score-based generative models for sparse-view CT and accelerated MRI
- Summary and future outlook

# Introduction

In computational imaging, an image sensor measurement  $\mathbf{y} \in \mathcal{Y}$  from an underlying unknown image  $\mathbf{x} \in \mathcal{X}$  is usually described by

$$\mathbf{y} = \mathbf{H}(\mathbf{x}) + \varepsilon \quad (1)$$

where  $\varepsilon$  is measurement noise and  $\mathbf{H} : \mathcal{X} \mapsto \mathcal{Y}$  is a forward mapping arising from the imaging physics.

- Penalized least squares  
computationally expensive iterative solvers  
highly dependent upon the regularization
- Deep learning (DL)  
usually based on supervised training
- Generative neural network  
potential to be used for unsupervised learning  
learn the distribution of the real-world data  
can incorporate the relevant physical laws and constraints of the measurement process

# Foundations for inverse imaging problems with generative models

- $P(\mathbf{x}|\mathbf{y}) = \frac{P(\mathbf{y}|\mathbf{x})P(\mathbf{x})}{P(\mathbf{y})}$  *Bayesian approach*
- The measurement likelihood  $P(\mathbf{y}|\mathbf{x})$  is sometimes available from the physics of the measurement process.
- The prior distribution  $P(\mathbf{x})$  is usually hard to model for computational imaging problems.
- To model the unknown prior distributions, we need to use data-driven methods. **Generative models** are especially useful in learning the prior distribution.

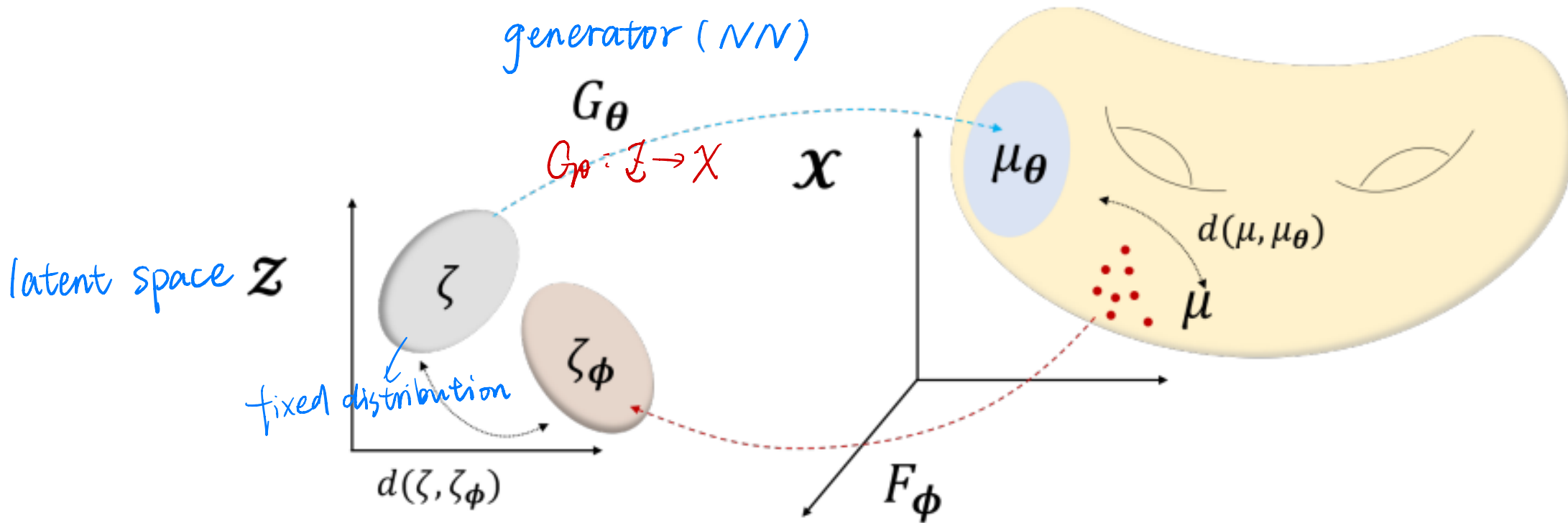


Fig. 5. Geometric view of deep generative models. Fixed distribution  $\zeta$  in  $\mathcal{Z}$  is pushed to  $\mu_\theta$  in  $\mathcal{X}$  by the network  $G_\theta$ , so that **the mapped distribution  $\mu_\theta$  approaches the real distribution  $\mu$** . In VAE,  $G_\theta$  works as a decoder to generate samples, while  $F_\phi$  acts as an encoder, additionally constraining  $\zeta_\phi$  to be as close to  $\zeta$ . With such geometric view, auto-encoding generative models (e.g. VAE), and GAN-based generative models can be seen as variants of this single illustration.

# VAEs (Variational AutoEncoder)

The standard VAE [6] models learn the data distribution  $p_{\theta}(\mathbf{x})$  by assuming that the data are generated by some random process involving an unobserved continuous random variable  $\mathbf{z}$  in the latent space  $\mathcal{Z}$ , i.e.,

$$p_{\theta}(\mathbf{x}) = \int_{\mathcal{Z}} p_{\theta}(\mathbf{x}, \mathbf{z}) d\mathbf{z} = \int_{\mathcal{Z}} p_{\theta}(\mathbf{x} | \mathbf{z}) \underbrace{p(\mathbf{z})}_{\text{prior}} d\mathbf{z}. \quad (3)$$

To match the distribution  $p_{\theta}(\mathbf{x})$  of generated samples to the true unknown data distribution  $p(\mathbf{x})$ , we wish to learn (identify) the model parameters  $\theta$  from training data, which can be done, in principle, by maximum likelihood. However, it is hard to learn  $\theta$  by directly maximizing the likelihood in (3) since it is intractable to evaluate the integral in (3) in high-dimensional space. Instead of directly evaluating the likelihood, we can compute a lower bound for the log-likelihood of a data point, called the *evidence lower bound (ELBO)*, or in the context of variational inference, the variational lower bound. This is done using an approximate posterior distribution  $q_{\phi}(\mathbf{z} | \mathbf{x})$  for the latent variable. As we will see shortly, this is where the probabilistic encoder comes in; its role is to model  $q_{\phi}(\mathbf{z} | \mathbf{x})$ .

Encoder:  $\mathbf{x} \rightarrow \mathbf{z}$ , is used only in learning  $\theta$  from data

Decoder:  $\mathbf{z} \rightarrow \mathbf{x}$ ,  $\mathbf{z} \sim p(\mathbf{z})$  prior  $\xrightarrow{\text{Decoder}}$   $p_{\theta}(\mathbf{x} | \mathbf{z})$   
 $\downarrow$  (3)  
 $p_{\theta}(\mathbf{x})$

$$\begin{aligned} \log p_{\theta}(\mathbf{x}) &= \log \left( \int p_{\theta}(\mathbf{x} | \mathbf{z}) p(\mathbf{z}) \frac{q_{\phi}(\mathbf{z} | \mathbf{x})}{q_{\phi}(\mathbf{z} | \mathbf{x})} d\mathbf{z} \right) \\ &\stackrel{\text{Jensen}}{\geq} \int \log \left( p_{\theta}(\mathbf{x} | \mathbf{z}) \frac{p(\mathbf{z})}{q_{\phi}(\mathbf{z} | \mathbf{x})} \right) q_{\phi}(\mathbf{z} | \mathbf{x}) d\mathbf{z} \\ &= -D_{\text{KL}}(q_{\phi}(\mathbf{z} | \mathbf{x}) \| p(\mathbf{z})) + \mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z} | \mathbf{x})} [\log p_{\theta}(\mathbf{x} | \mathbf{z})] \\ &\equiv L_{\text{ELBO}}(\phi, \theta; \mathbf{x}) \end{aligned} \quad (4)$$

$$\log p_{\theta}(\mathbf{x}) = \int_{\mathcal{Z}} q_{\phi}(\mathbf{z} | \mathbf{x}) \log p_{\theta}(\mathbf{x}) d\mathbf{z} \rightarrow$$

$$L_{\text{ELBO}}(\phi, \theta; \mathbf{x}) = \log p_{\theta}(\mathbf{x}) - D_{\text{KL}}(q_{\phi}(\mathbf{z} | \mathbf{x}) \| p_{\theta}(\mathbf{z} | \mathbf{x})). \quad (5)$$

It follows that maximizing the ELBO simultaneously attempts to maximize the data log-likelihood with respect to  $\theta$  and minimize the divergence between the true posterior  $p_{\theta}(\mathbf{z} | \mathbf{x})$  and its approximant  $q_{\phi}(\mathbf{z} | \mathbf{x})$ .

The approximate posterior  $q_\phi(\mathbf{z}|\mathbf{x})$  is assumed to be a multivariate normal distribution with diagonal covariance, and the latent variable  $\mathbf{z}$  is reparameterized as

$$\mathbf{z} = \boldsymbol{\mu}_\phi(\mathbf{x}) + \boldsymbol{\sigma}_\phi(\mathbf{x}) \odot \boldsymbol{\varepsilon} \quad (6)$$

where  $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ , and  $\odot$  denotes the Hadamard product or elementwise product. As shown in Figure 1, the encoder (realized by an NN parameterized by  $\phi$ ) outputs the mean  $\boldsymbol{\mu}_\phi(\mathbf{x})$  and the standard deviation  $\boldsymbol{\sigma}_\phi(\mathbf{x})$  given the input  $\mathbf{x}$ .

Assuming that the prior  $p(\mathbf{z})$  is a standard normal distribution, the loss function for training the VAE model end to end is to minimize the negative ELBO with respect to  $\boldsymbol{\theta}, \phi$

$$\begin{aligned} \mathcal{L}_{\text{VAE}}(\boldsymbol{\theta}, \phi) &= \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} [-L_{\text{ELBO}}(\phi, \boldsymbol{\theta}; \mathbf{x})] \\ &\approx \sum_{i=1}^N \left[ \frac{1}{2} (\|\boldsymbol{\mu}_\phi(\mathbf{x}^{(i)})\|_2^2 + \|\boldsymbol{\sigma}_\phi(\mathbf{x}^{(i)})\|_2^2 - 2\|\log(\boldsymbol{\sigma}_\phi(\mathbf{x}^{(i)}))\|_1 \right. \\ &\quad \left. - 1) - \frac{1}{L} \sum_{l=1}^L \log p_\theta(\mathbf{x}^{(i)} | \mathbf{z}^{(i,l)}) \right] \end{aligned} \quad (7)$$

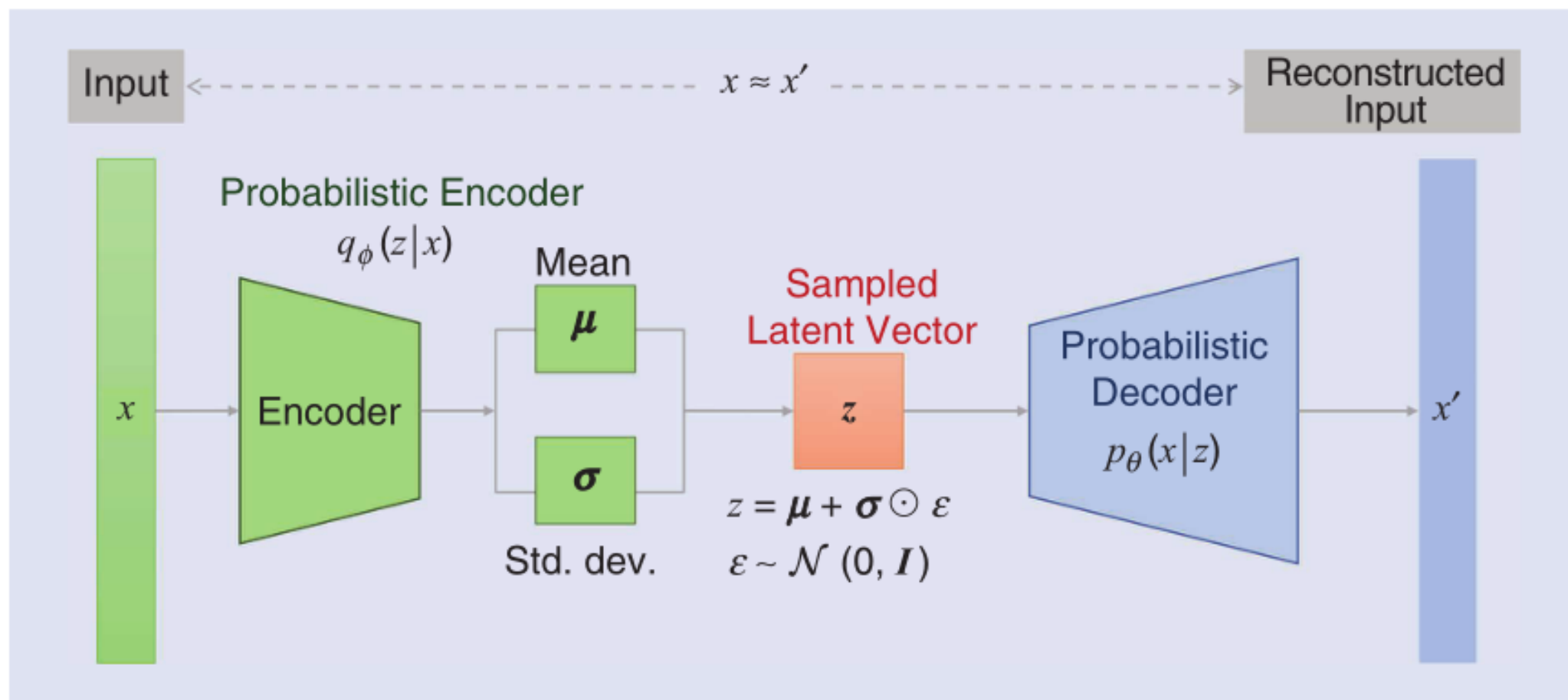
- Advantage  
easy and fast to train

- Limitation

$q_\phi(\mathbf{z}|\mathbf{x})$  is often restricted to simple distributions

inherent discrepancy between  $q_\phi(\mathbf{z}|\mathbf{x})$  and  $p_\theta(\mathbf{z}|\mathbf{x})$  remains

it cannot provide the exact likelihood



**FIGURE 1.** The standard VAE structure with probabilistic encoder and decoder. The sampled latent vector is constructed through the reparameterization trick in (6). Std. dev.: standard deviation.

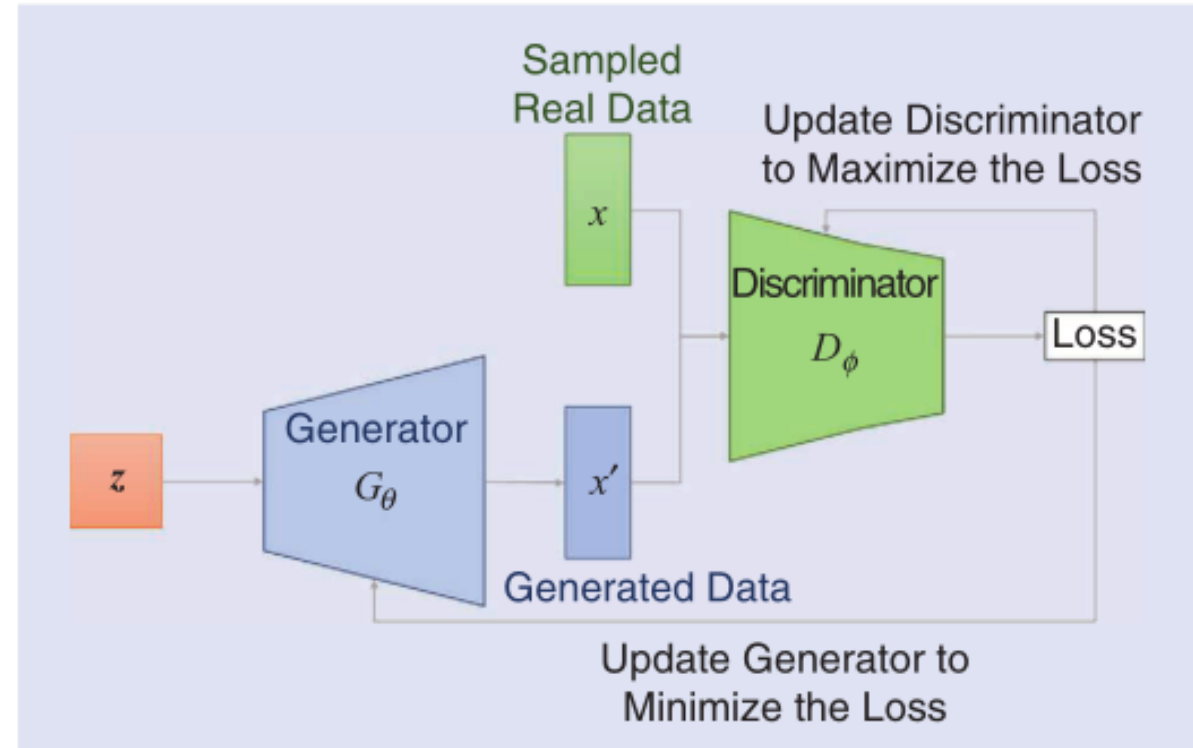


# GANs (Generative Adversarial Network)

A GAN [7], [8], [23] (Figure 2) is composed of a generator  $G_\theta$ , which, when driven by Gaussian random vectors  $z \sim N(\mathbf{0}, \mathbf{I})$ , generates samples from the generator distribution  $p_{\text{gen}}$ , and a discriminator  $D_\phi$  that compares the data distribution  $p_{\text{data}}$  with  $p_{\text{gen}}$ . Both are implemented by NNs, with parameters  $\theta, \phi$ . GANs pose generative modeling as a problem of minimizing a statistical distance or divergence between probability distributions. The  $f$ -GAN [23] uses the  $f$ -divergence as the measure of the divergence between two distributions  $P$  and  $Q$  with density function  $p$  and  $q$

$$D_f(P \parallel Q) = \int_{\mathcal{X}} q(x) f\left(\frac{p(x)}{q(x)}\right) dx \quad (8)$$

where the function  $f: \mathbb{R}_+ \rightarrow \mathbb{R}$  is a convex lower-semicontinuous function satisfying  $f(1) = 0$ . For example, taking  $f(s) = s \log s$ , the  $f$ -divergence becomes the KL divergence.



**FIGURE 2.** The standard GAN architecture.

For  $D_\phi$ ,  $\max_{D_\phi} \mathcal{L}_{GAN}$   $\begin{cases} D_\phi(x) \rightarrow 1 & \text{real} \\ D_\phi(G_\theta(z)) \rightarrow 0 & \text{fake} \end{cases}$

For  $G_\theta$ ,  $\min_{G_\theta} \mathcal{L}_{GAN}$ ,  $D_\phi(G_\theta(z)) \rightarrow 1$  real

$$\min_{G_\theta} \max_{D_\phi} \mathcal{L}_{GAN}(D_\phi, G_\theta) \quad (9)$$

where

$$\begin{aligned} \mathcal{L}_{GAN}(D_\phi, G_\theta) := & \mathbb{E}_{x \sim p_{\text{data}}} [\log D_\phi(x)] \\ & + \mathbb{E}_z [\log(1 - D_\phi(G_\theta(z)))]. \end{aligned} \quad (10)$$

The discriminator output  $D_\phi(x) \in [0, 1]$ , with a value close to one, indicating a sample likely drawn from the data distribution, whereas a value close to zero indicates a “fake” (unlikely) sample. The min-max problem (9) is therefore interpreted as a game between two adversaries—the discriminator, trying to improve its detection of fake samples, and the generator, trying to improve its ability to fool the discriminator.

It was found that GANs using statistical discrepancy measures such as the Jensen–Shannon divergence are especially hard to train, in part because the divergence measure between two distributions “saturates” (to infinity) when their support does not overlap. This scenario is common in the case of image data, which typically lie on (or close to) a low-dimensional manifold in the ambient space, which with high probability does not intersect with the corresponding manifold of the untrained generator. This leads to the vanishing of the gradients of the training loss function, preventing convergence.

- Wasserstein-GAN

$p$ -Wasserstein distance between distributions  $p_{\text{data}}$  and  $p_{\text{gen}}$  is defined as

$$W_p(p_{\text{data}}, p_{\text{gen}}) = \inf_{\gamma \in \Pi(p_{\text{data}}, p_{\text{gen}})} (\mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \gamma} [\|\mathbf{x} - \mathbf{y}\|^p])^{1/p} \quad (11)$$

Estimating the Wasserstein distance in high-dimensional space is, however, not straightforward. Arjovsky et al. [8] applied the Kantorovich–Rubinstein duality to the 1-Wasserstein distance, which states

$$W_1(p_{\text{data}}, p_{\text{gen}}) = \sup_{\|f\|_{\text{Lip}} \leq 1} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [f(\mathbf{x})] - \mathbb{E}_{\mathbf{x} \sim p_{\text{gen}}} [f(\mathbf{x})] \quad (12)$$

For  $D_\phi$ ,  $\max_{\|D_\phi\|_{\text{Lip}} \leq 1} L \cdot \begin{cases} D_\phi(x) \rightarrow 1 & \text{real} \\ D_\phi(G_\theta(z)) \rightarrow 0 & \text{fake} \end{cases}$

For  $G_\theta$ ,  $\min_{G_\theta} L, D_\phi(G_\theta(z)) \rightarrow 1$  real

where the supremum is taken over all 1-Lipschitz functions,  $\|f\|_{\text{Lip}} \leq 1$ , which means that  $|f(\mathbf{x}_1) - f(\mathbf{x}_2)| \leq \|\mathbf{x}_1 - \mathbf{x}_2\|_1$  for all  $\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{X}$ . The function  $f$  is approximated by a deep NN  $D_\phi$  with parameters  $\phi$ . Therefore, the Wasserstein GAN is formulated as

$$\min_{G_\theta} \max_{\|D_\phi\|_{\text{Lip}} \leq 1} \{ \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [D_\phi(\mathbf{x})] - \mathbb{E}_z [D_\phi(G_\theta(z))] \} \quad (13)$$

and various ways have been suggested to enforce the Lipschitz constraint [24].

However, they're hard to train because of the min-max nature of optimization objectives.

# Score-based generative models

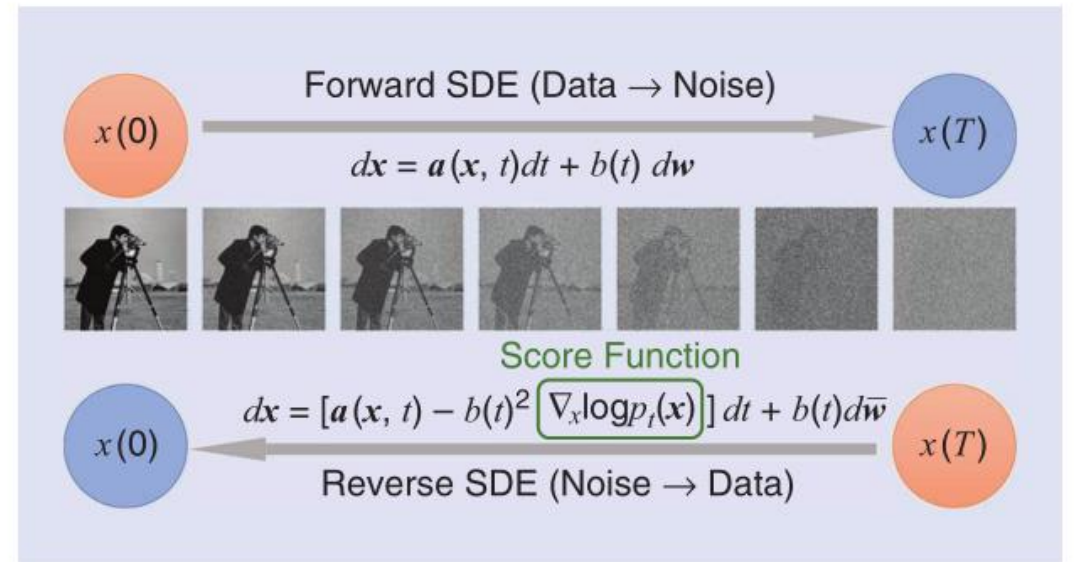
Score-based generative models [11], [12], [14] generate new data from noise through learning the gradient of the log probability of the data, also known as the *score function* [29]. To see why score matching is useful in learning a data distribution, let's first consider a probability density function defined in terms of a parameterized function  $f_\theta$  as  $p_\theta(\mathbf{x}) = (e^{-f_\theta(\mathbf{x})}/Z_\theta)$ , where  $Z_\theta$  (a so-called partition function) is a normalizing constant such that  $\int_{\mathcal{X}} p_\theta(\mathbf{x}) d\mathbf{x} = 1$ . The parameterized function  $f_\theta(\mathbf{x})$  is often called an *energy-based model*.

$$s_\theta(\mathbf{x}) = \nabla_{\mathbf{x}} \log p_\theta(\mathbf{x}) = -\nabla_{\mathbf{x}} f_\theta(\mathbf{x}) - \underbrace{\nabla_{\mathbf{x}} \log Z_\theta}_{=0} = -\nabla_{\mathbf{x}} f_\theta(\mathbf{x}). \quad (14)$$

Therefore, the score-based models are trained by minimizing the Fisher divergence between the model and the data distributions,  $\mathbb{E}_{p_\theta} \|s_\theta(\mathbf{x}) - \nabla_{\mathbf{x}} \log p_{\text{data}}(\mathbf{x})\|_2^2$ .

↓ model  
↓ data distribution  
approximate the score of  $p_\theta(\mathbf{x})$

Another important ingredient in the score-based generative models is modeling the process that transforms the data to noise and its reverse process. *diffusion model*



**FIGURE 3.** An illustration of the score-based generative model. The forward continuous time SDE transforms data to an image sampled from a simple noise distribution. This process can be reversed if we know  $\nabla_{\mathbf{x}} \log p_t(\mathbf{x})$  (the score function of the distribution) at each intermediate time step. We can sample a data point through evolving a noise image through the reverse-time SDE.



$$p_{0t}(\mathbf{x}(t) | \mathbf{x}(0)) = \mathcal{N}(\mathbf{x}(t) | \alpha(t)\mathbf{x}(0), \beta^2(t)\mathbf{I})$$

add noise

image

$$\mathbf{x}(t) = \alpha(t)\mathbf{x}(0) + \beta(t)\boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}). \quad (16)$$

The forward process in (16) is used to generate the training data  $\mathbf{x}(t)$  for  $t \in (0, T]$  to learn the gradient of the log probability of the data (score function).

The model is trained by minimizing the weighted Fisher divergence between the model and data distributions and over the time interval  $[0, T]$

$$\min_{\boldsymbol{\theta}} \mathbb{E}_{t \sim \text{Uniform}(0, T)} [\lambda(t) \mathbb{E}_{\mathbf{x}(0) \sim p_{\text{data}}} \mathbb{E}_{\mathbf{x}(t) | \mathbf{x}(0)} \| s_{\boldsymbol{\theta}}(\mathbf{x}(t), t) - \nabla_{\mathbf{x}(t)} \log p_{0t}(\mathbf{x}(t) | \mathbf{x}(0)) \|^2] \quad (18)$$

where the weight  $\lambda(t) > 0$  is typically chosen as  $1/\lambda(t) \propto \mathbb{E}_{\mathbf{x}(0) \sim p_{\text{data}}} \mathbb{E}_{\mathbf{x}(t) | \mathbf{x}(0)} \| \nabla_{\mathbf{x}} \log p_{0t}(\mathbf{x}(t) | \mathbf{x}(0)) \|^2$ . In the training, the expectations in the loss function (18) are replaced by empirical means by sampling  $\mathbf{x}(0)$  from the training data, and for each  $\mathbf{x}(0)$ , sampling  $\mathbf{x}(t)$  according to (16).

After training, there are a few different ways to sample the distribution and generate new data  $\mathbf{x}(0) \sim p_{\text{data}}$  from  $\mathbf{x}(T) \sim p_T$ .

### Reverse SDE

One is to use a numerical solver for the SDE in (17), such as the Euler–Maruyama approach where  $dt$  is approximated by  $\Delta t = (T/N)$ , with discrete time steps  $0 = t_0 < t_1 \cdots < t_N = T$ . Another way is to use the predictor–corrector (PC) method. PC sampling uses one step of a numerical solver, which is called the *predictor step*, to generate sample  $\mathbf{x}(t_{i-1})$  from  $\mathbf{x}(t_i)$ .

A fundamentally different approach to transform noise to data is based on the theoretical result that the reverse-time SDE can be converted into an ordinary differential equation (ODE) without changing its time-evolving distribution  $\{p_t(\mathbf{x})\}_{t \in [0, T]}$

# Physics-driven applications

(Cryo-electron microscopy)

- Generative modeling for cryo-EM analysis

Cryo-EM aims to recover the 3D structure of a particle of interest  $V: \mathbb{R}^3 \rightarrow \mathbb{R}$  (called *volume*) from a collection of noisy and blurred line-integral projection images  $y_1, \dots, y_N$  of the volume taken at unknown projection angles with unknown shifts. These images are obtained by examining a frozen sample containing multiple randomly oriented and shifted copies of the volume. The generation of image  $y$  can be modeled as

$$y = C_d * S_t P_R(V) + \epsilon \quad (21)$$

where the  $P_R(V)$  is the tomographic projection of  $V$  rotated by a rotation matrix  $R \in SO(3)$ , the group of 3D rotations

$$P_R(V)(r_x, r_y) = \int_{\mathbb{R}} V(R^T \mathbf{r}) dr_z, \quad \mathbf{r} = (r_x, r_y, r_z)^T. \quad (22)$$

The Fourier transform of the resulting image  $\hat{y}$  is modulated pointwise by the contrast transfer function (CTF)  $\hat{C}_d$ , which is determined by the defocus value and other parameters of the electron microscope  $d$  and is assumed to be known. In the spatial domain, this effect corresponds to the convolution of the inverse Fourier transform of the CTF with the projection image. We assume that  $\epsilon$  is additive WGN with noise variance  $\sigma^2$ . The image size is  $D \times D$ . According to the Fourier slice theorem, the 2D Fourier transform of the image is given by

$$\hat{y}(k_x, k_y) = \hat{C}_d(k_x, k_y) \tilde{S}_t(k_x, k_y) \hat{V}(R^T [k_x, k_y, 0]^T) + \hat{\epsilon}(k_x, k_y) \quad (23)$$

2D Fourier transform of a projection of a 3D object can be computed as the set of Fourier transforms of the object along lines

St shift the centered 2D projection image by  $t = (t_x, t_y)$  perpendicular to the direction of the projection.

### Single conformation model

The traditional cryo-EM reconstruction [43] uses the **maximum a posteriori (MAP) estimate** of  $\mathbf{V}$  from  $N$  projection images, marginalizing over the posterior distribution of the  $\boldsymbol{\varphi} = (\mathbf{R}, t)$

$$\mathbf{V}_{\text{rec}} = \underset{\mathbf{V}}{\operatorname{argmax}} \sum_{i=1}^N \log \left( \int_{SO(3) \times \mathbb{R}^2} p(\mathbf{y}_i | \boldsymbol{\varphi}, \mathbf{V}) p(\boldsymbol{\varphi}) d\boldsymbol{\varphi} \right) + \log p(\mathbf{V}). \quad (24)$$

### Multiple conformation model *Structural heterogeneity*

In multiclass refinement [43], the image formation model is extended to assume that images are generated from  $K$  independent volumes,  $\mathbf{V}_1, \dots, \mathbf{V}_K$ , and the inference requires marginalization over both the pose parameters  $\boldsymbol{\varphi}_i$  and the class assignment.

$$\underset{\mathbf{V}_1, \dots, \mathbf{V}_K}{\operatorname{argmax}} \sum_{i=1}^N \log \sum_{j=1}^K \left( \pi_j \int_{SO(3) \times \mathbb{R}^2} p(\mathbf{y}_i | \boldsymbol{\varphi}, \mathbf{V}_j) p(\boldsymbol{\varphi}) d\boldsymbol{\varphi} \right) + \sum_{j=1}^K \log p(\mathbf{V}_j). \quad (25)$$

In CryoSPARC [44], the class assignment probabilities are assumed to be uniform, i.e.,  $\pi_j = 1/K$ . This approach requires **prior knowledge of the number of classes** and is computationally feasible **only for a small number of classes**. However, **the protein conformation changes are continuous** and may be poorly approximated with a small number of discrete volumes.

# VAE-based generative models

CryoDRGN (Deep Reconstructing Generative Network) [47] employs VAE to learn a continuous low-dimensional manifold over a protein's conformational states from 2D cryo-EM images in an unsupervised way and perform ab initio (i.e., not as a refinement of an existing volume) reconstruction of the volumes.

The model contains a standard VAE probabilistic encoder  $q_\phi(\mathbf{z}|\hat{\mathbf{y}}_i)$ .

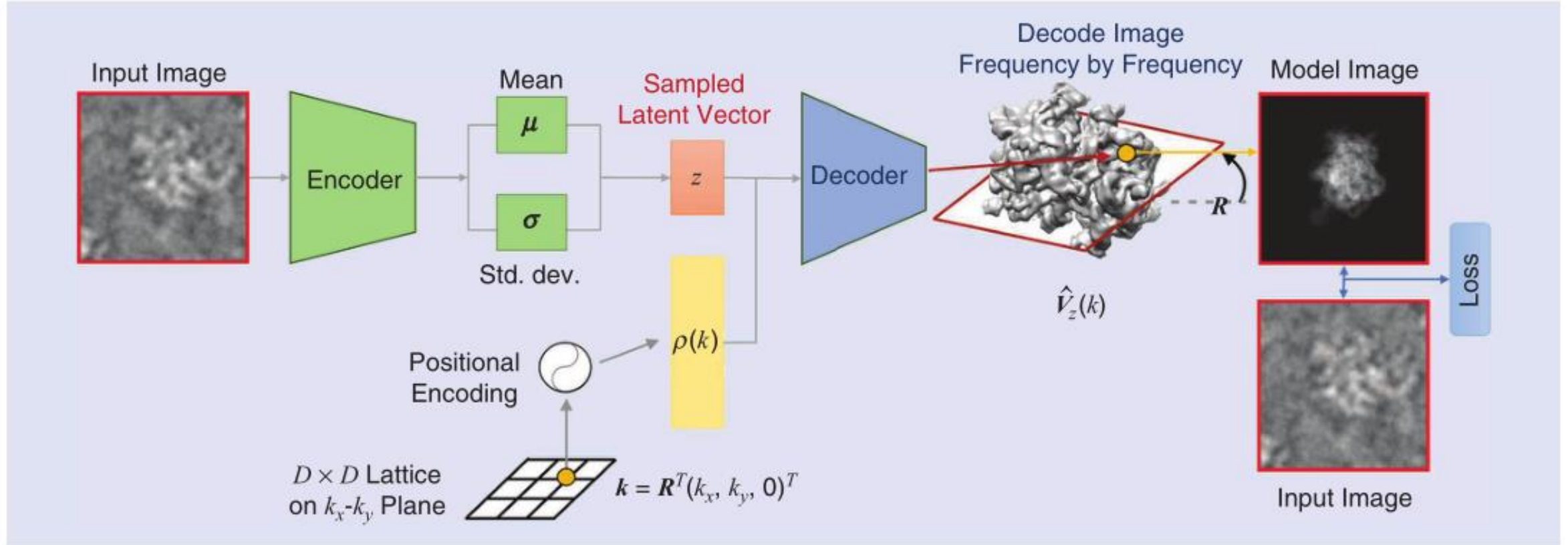
The Fourier transform  $\hat{V}$  of the volume in CryoDRGN is modeled by a probabilistic decoder  $p_\theta(\hat{V}|\mathbf{z}, \mathbf{k})$  with two separate inputs: 1) the latent variable  $\mathbf{z}$  from the encoder and 2) the frequency  $\mathbf{k}$ .

3D coordinate  $\mathbf{k} \in [-0.5, 0.5]^3$ , CryoDRGN use positional encoding  $\boldsymbol{\rho} : [-0.5, 0.5]^3 \rightarrow [-1, 1]^{3D}$  to map the 3D coordinate to a higher dimensional vector. For  $\mathbf{k} = (k_1, k_2, k_3)^T$ , the positional encoding consists of sine and cosine waves, and the mapping is defined as

$$\begin{aligned}\boldsymbol{\rho}^{(2i)}(k_j) &= \sin\left(k_j D \pi \left(\frac{2}{D}\right)^{\frac{2i}{D}}\right), \\ \boldsymbol{\rho}^{(2i+1)}(k_j) &= \cos\left(k_j D \pi \left(\frac{2}{D}\right)^{\frac{2i}{D}}\right), \\ i &= 1, \dots, \frac{D}{2}; \quad j = 1, 2, 3.\end{aligned}\tag{26}$$

It is empirically observed that using this encoding works well for clean data. For noisy data, it is required to exclude 10% of the high-frequency components. The decoder outputs the volume at frequency  $\mathbf{k}$  as  $\hat{V}_z(\mathbf{k}) = G_\theta(\mathbf{z}, \boldsymbol{\rho}(\mathbf{k}))$ .





**FIGURE 4.** The CryoDRGN model architecture. The VAE is used to perform approximate inference for latent variable  $z$  denoting structural heterogeneity. The decoder reconstructs an image frequency by frequency in  $k$ -space, given  $z$  and  $\rho(k)$ , the positional encoding of 3D Cartesian  $k$ -space coordinates. The 3D frequency coordinates corresponding to each input image Fourier coefficient are obtained by rotating a  $D \times D$  lattice on the  $k_x - k_y$  plane by  $R$ , the orientation of the particle. The latent orientation  $R$  for each image is inferred by a branch and bound global optimization procedure. (Source: Figure courtesy of the authors of [37].)

Using the Fourier slice theorem in (23) and the decoder  $G_\theta$ , we can evaluate the negative log-likelihood of the image  $\hat{y}$  given the latent variable  $z \sim q_\phi(z|\hat{y})$  and the pose parameters  $(\mathbf{R}, t)$

$$\begin{aligned} & -\log p_\theta(\hat{y}|z, \mathbf{R}, t) \\ &= \frac{1}{2\sigma^2} \sum_{k_x} \sum_{k_y} |\hat{y}(k_x, k_y) - \hat{\mathbf{C}}_d(k_x, k_y) \tilde{\mathbf{S}}_t(k_x, k_y) \\ & \quad \times G_\theta(z, \rho(\mathbf{R}^T [k_x, k_y, 0]^T))|^2 + D^2 \log(\sqrt{2\pi} \sigma). \quad (27) \end{aligned}$$

CryoDRGN jointly estimates the pose parameters and the VAE network parameters by alternating between updating those two sets of parameters. When the pose parameters  $(\mathbf{R}, t)$  are fixed, CryoDRGN minimizes the negative ELBO using the Adam optimizer [48] with respect to  $\phi$  and  $\theta$  following the standard VAE framework

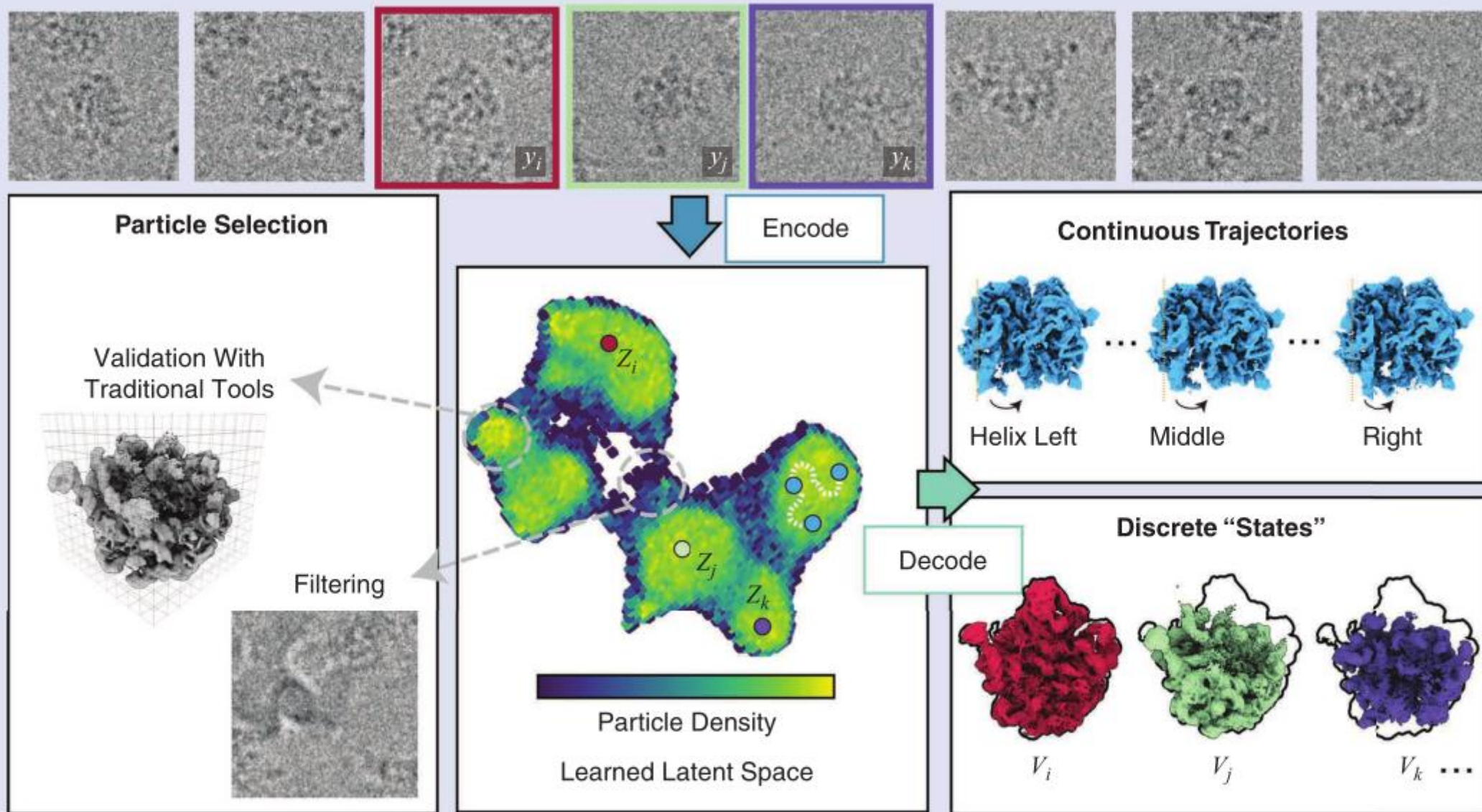
$$\begin{aligned} \mathcal{L}(\phi, \theta) = & \sum_{m=1}^B D_{KL}(q_\phi(z|\hat{y}_{n_m}) \| p(z)) \\ & - \mathbb{E}_{z \sim q_\phi(z|\hat{y}_{n_m})} [\log p_\theta(\hat{y}_{n_m}|z, \mathbf{R}, t)] \quad (28) \end{aligned}$$

where  $B$  is the batch size, and  $\{n_m\}_{m=1}^B$  is a set of random indices for projection images chosen at each training iteration.

The physics-based measurement model is incorporated in the negative log-likelihood in (27), which appears in the negative ELBO (28) for training the VAE model. With fixed VAE network parameters  $\phi$  and  $\theta$ , a global search over  $SO(3) \times \mathbb{R}^2$  is performed for the maximum-likelihood estimation of the pose  $(\mathbf{R}, t)$  for each image given the decoded volume  $\hat{V}_z(\mathbf{k}) = G_\theta(z, \rho(\mathbf{k}))$ . An efficient joint maximum-likelihood estimator pose estimation and a VAE training scheme for cryo-EM are detailed in [49].

After training, CryoDRGN provides per-particle image latent encoding. The encoder network outputs  $z_i = \mu_\phi(\hat{y}_i)$ , which corresponds to the approximate MAP estimate of the latent variable  $z_i$  given  $\hat{y}_i$ , i.e.,  $z_i = \operatorname{argmax}_z q_\phi(z|\hat{y}_i) = \mu_\phi(\hat{y}_i)$ . The trained decoder network can then generate 3D volumes given arbitrary values of the latent variable  $z$  via  $\hat{V}_z(\mathbf{k}) = G_\theta(z, \rho(\mathbf{k}))$ .





**FIGURE 5.** A structural heterogeneity analysis. With the trained CryoDRGN encoder, all particle images of the ribosome dataset [51] are encoded into the latent space, as shown in the center panel. The trained decoder can map the latent variable to the 3D volumes, as indicated in the right panel. In addition, the latent space representation can be used to filter particle images and remove impurities. (Source: Figure courtesy of the authors of [47].)

# GAN-based generative models

CryoGAN [4] and Multi-CryoGAN [18] address the 3D reconstruction problem in cryo-EM from the perspective of distribution matching. These methods adopt the **Wasserstein GAN structure** and use a **physics-based cryo-EM data simulator** to generate the projection images. Unlike the maximum likelihood-based approaches, the adversarial learning-based approaches **do not require the estimation of the individual pose parameters**. Instead, in CryoGAN [4] and Multi-CryoGAN [18], the imaging parameters (rotations, in-plane translations, and CTF parameters)  $\boldsymbol{\varphi} = (\mathbf{R}, t, d)$  are drawn from a **known prior distribution  $p_{\boldsymbol{\varphi}}$** .

CryoGAN is used to reconstruct a single 3D volume from the image data, assuming a single conformation. **Given the 3D volume  $V$** , the distribution of the generated images is denoted by  $p_{\text{gen}}(\mathbf{y}; V)$ . The reconstruction is achieved by minimizing the 1-Wasserstein distance between  $p_{\text{gen}}(\mathbf{y}; V)$  and  $p_{\text{data}}(\mathbf{y})$  and becomes the following min-max optimization problem:

$$V_{\text{rec}} = \underset{V}{\operatorname{argmin}} \max_{\|D_{\phi}\|_{\text{Lip}} \leq 1} (\mathbb{E}_{\mathbf{y} \sim p_{\text{data}}} [D_{\phi}(\mathbf{y})] - \mathbb{E}_{\mathbf{y} \sim p_{\text{gen}}(\mathbf{y}; V)} [D_{\phi}(\mathbf{y})]). \quad (29)$$

The Lipschitz constraint  $\|D_{\phi}\|_{\text{Lip}}$  is enforced by penalizing the norm of the gradient of  $D_{\phi}$  with respect to its input. Based on empirical samples, the loss in (29) is reformulated as

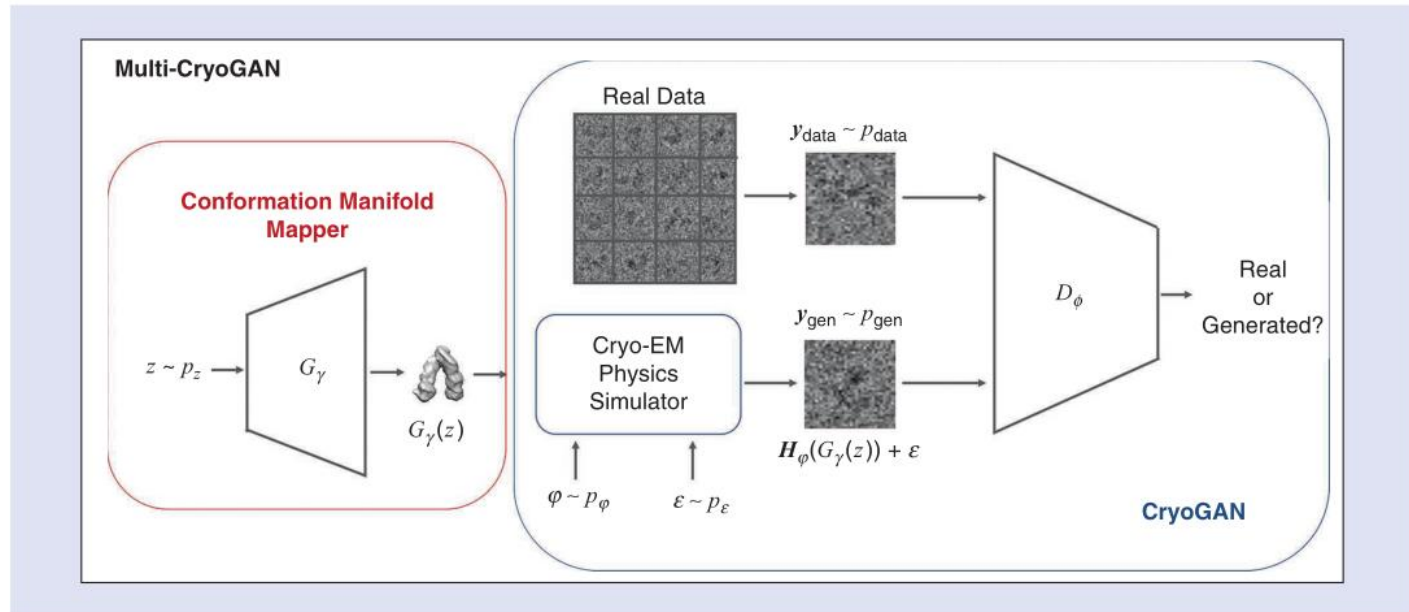
$$\mathcal{L}(V, \boldsymbol{\phi}) = \sum_{m=1}^B (D_{\phi}(\mathbf{y}_{\text{data}}^{n_m}) - D_{\phi}(\mathbf{y}_{\text{gen}}^m) - \lambda[(\|\nabla_{\mathbf{y}} D_{\phi}(\mathbf{y}_{\text{int}}^m)\| - 1)^2]) \quad (30)$$

where  $B$  is the batch size;  $\{n_m\}_{m=1}^B$  is a set of random indices for projection images selected at each iteration;  $\mathbf{y}_{\text{gen}}^m$  denotes the projections from the current estimate  $V$  generated by the cryo-EM physics simulator according to (21); and  $\lambda \in \mathbb{R}_+$  is a gradient penalty coefficient. The interpolated sample used in the gradient penalty is  $\mathbf{y}_{\text{int}}^m = a_m \mathbf{y}_{\text{data}}^{n_m} + (1 - a_m) \mathbf{y}_{\text{gen}}^m$ , where  $a_m$  is sampled from a uniform distribution on the interval  $[0, 1]$ . Using the interpolated samples ensures that the learned discriminator function is 1-Lipschitz on the domain spanned by both the generated and real data.



To address the multiconformation setting, Multi-CryoGAN [18] (Figure 6) adds a convolutional NN (CNN)  $G_\gamma$  to the CryoGAN architecture to learn a mapping from a latent space to the 3D conformation distribution. It samples a latent vector  $z$  from a prior distribution  $p_z$ . Then, a CNN  $G_\gamma$  maps the input latent variable  $z$  to the conformation manifold, i.e.,  $V_z = G_\gamma(z)$ . Based on the generated volume  $V_z$ , the cryo-EM physics-based simulator generates noisy projection images. The distribution of the generated projection images is denoted by  $p_{\text{gen}}(\mathbf{y}; G_\gamma)$ . To find the network parameters in the generator  $G_\gamma$ , Multi-CryoGAN minimizes the 1-Wasserstein distance between  $p_{\text{data}}(\mathbf{y})$  and  $p_{\text{gen}}(\mathbf{y}; G_\gamma)$ , which results in the following min-max optimization problem.

$$\begin{aligned} \gamma^* &= \underset{\gamma}{\operatorname{argmin}} W_1(p_{\text{data}}, p_{\text{gen}}) \\ &= \underset{\gamma}{\operatorname{argmin}} \max_{\|D_\phi\|_{\text{Lip}} \leq 1} (\mathbb{E}_{\mathbf{y} \sim p_{\text{data}}} [D_\phi(\mathbf{y})] - \mathbb{E}_{\mathbf{y} \sim p_{\text{gen}}(\mathbf{y}; G_\gamma)} [D_\phi(\mathbf{y})]). \end{aligned} \quad (31)$$



**FIGURE 6.** The model architecture for Multi-CryoGAN, which contains the architecture for (a) conformation manifold mapper (in a red box) and (b) CryoGAN (in a blue box). (Source: Figure courtesy of the authors of [18].)

- Score-based generative models for sparse-view CT and accelerated MRI

In sparse-view CT and accelerated MRI, we consider the general linear measurement model

$$\mathbf{y} = \mathbf{A}\mathbf{x} + \boldsymbol{\varepsilon} = \mathcal{P}(\boldsymbol{\Lambda})\mathbf{M}\mathbf{x} + \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (42)$$

where  $\mathbf{M}$  corresponds to the Radon transform in sparse-view CT and Fourier transform in accelerated (undersampled) MRI, respectively, and  $\boldsymbol{\Lambda}$  is an  $n \times n$  diagonal matrix where  $\Lambda_{ii} \in \{0, 1\}$  with  $\text{tr}(\boldsymbol{\Lambda}) = m$ , and  $\mathcal{P}(\boldsymbol{\Lambda}) \in \{0, 1\}^{m \times n}$  is an operator that keeps only the rows of  $\boldsymbol{\Lambda}$  that  $\Lambda_{i,i} \neq 0$ . Matrix  $\mathbf{A}$  is assumed to have full row rank.

In the “Score-Based Generative Models” section, we discussed the unconditional sampling of  $\mathbf{x}(t)$  to generate data  $\mathbf{x}(0)$  from noise  $\mathbf{x}(T)$  by following the reverse-time SDE. Given the measurement data  $\mathbf{y}$ , in principle, the reconstruction should be achieved by generating approximate samples from the conditional stochastic process  $\{\mathbf{x}(t) | \mathbf{y}\}$  from  $t = T$  to  $t = 0$ . However, it is difficult to directly solve the conditional reverse-time SDE  $\{\mathbf{x}(t) | \mathbf{y}\}_{t \in [0, T]}$  without using paired training data to learn the conditional score function. For simplicity, we set  $T = 1$  for the discussions later.

To overcome the difficulty in sampling directly from  $\{\mathbf{x}(t) \mid \mathbf{y}\}$ , [1] introduced the following stochastic process  $\{\mathbf{y}(t)\}_{t \in [0,1]}$ , given the unconditional stochastic process  $\{\mathbf{x}(t)\}_{t \in [0,1]}$

$$\begin{aligned} \mathbf{y}(t) &\triangleq \mathbf{A}\mathbf{x}(t) + \alpha(t)\boldsymbol{\varepsilon} \stackrel{(a)}{=} \mathbf{A}(\alpha(t)\mathbf{x}(0) + \beta(t)\mathbf{z}) + \alpha(t)\boldsymbol{\varepsilon} \\ &\stackrel{(b)}{=} \alpha(t)\mathbf{y} + \beta(t)\mathbf{A}\mathbf{z} \end{aligned} \quad (43)$$

where  $\mathbf{z} \in \mathbb{R}^n \sim \mathcal{N}(0, \mathbf{I})$ . Equality (a) follows from the Gaussian transition kernel  $p_{0t}(\mathbf{x}(t) \mid \mathbf{x}(0)) = \mathcal{N}(\mathbf{x}(t) \mid \alpha(t)\mathbf{x}(0), \beta^2(t)\mathbf{I})$

Therefore, it is tractable to generate samples  $\tilde{\mathbf{y}}(t) \sim p_t(\mathbf{y}(t) \mid \mathbf{y})$  according to (43). In addition, we can modify the iterative sampling algorithm designed for the unconditional stochastic process  $\{\mathbf{x}(t)\}_{t \in [0,1]}$  by encouraging data consistency between the samples  $\tilde{\mathbf{x}}(t) \sim p_t(\mathbf{x}(t))$  and the samples  $\tilde{\mathbf{y}}(t)$  at each time  $t$  by constructing the intermediate samples

$$\begin{aligned} \tilde{\mathbf{x}}'(t) &= \underset{\mathbf{v} \in \mathbb{R}^n}{\operatorname{argmin}} \{ (1 - \lambda) \|\mathbf{v} - \tilde{\mathbf{x}}(t)\|_M^2 + \lambda \min_{\mathbf{u} \in \mathbb{R}^n} \|\mathbf{v} - \mathbf{u}\|_M^2 \}, \\ &\text{subject to } \mathbf{A}\mathbf{u} = \tilde{\mathbf{y}}(t) \end{aligned} \quad (44)$$

where  $\lambda \in [0, 1]$  is a parameter that balances between data consistency with the unconditional generation  $\tilde{\mathbf{x}}(t)$  and the conditional generation  $\tilde{\mathbf{y}}(t)$ . The weighted norm  $\|\cdot\|_M^2$  is defined as  $\|\mathbf{x}\|_M^2 = \|\mathbf{M}\mathbf{x}\|_2^2$ .

When  $\lambda = 0$ , we have  $\tilde{\mathbf{x}}'(t) = \mathbf{x}(t)$ , and we don't incorporate the stochastic process of  $\mathbf{x}(t)$  associated with the observation  $\mathbf{y}$ . This parameter can be tuned on a validation dataset using Bayesian optimization. In summary, the physics of the measurement process is incorporated into the conditional generation through two steps: 1) define the process  $\{\mathbf{y}(t)\}$  according to (43) and 2) encourage data consistency according to (44).

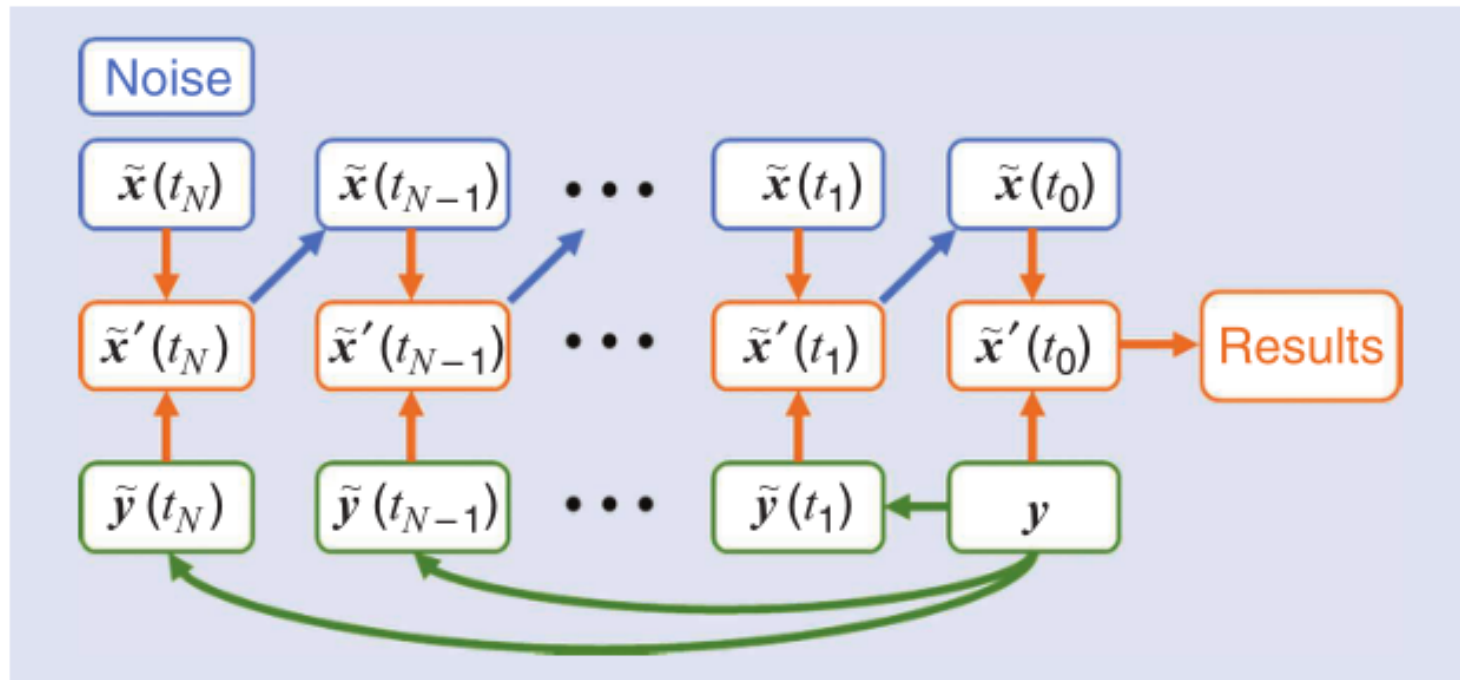
The sampling process selects a sequence of times  $0 = t_0 < t_1 < \dots < t_N = 1$  and iterates according to

$$\tilde{\mathbf{y}}(t_i) = \alpha(t_i)\mathbf{y} + \beta(t_i)\mathbf{A}\mathbf{z}_i, \quad \mathbf{z}_i \in \mathbb{R}^n \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (45a)$$

$$\begin{aligned} \tilde{\mathbf{x}}'(t_i) &= \mathbf{M}^{-1} [\lambda \boldsymbol{\mathcal{P}}^\dagger(\boldsymbol{\Lambda})\tilde{\mathbf{y}}(t_i) + (1 - \lambda)\boldsymbol{\Lambda}\mathbf{M}\tilde{\mathbf{x}}(t_i) \\ &\quad + (\mathbf{I} - \boldsymbol{\Lambda})\mathbf{M}\tilde{\mathbf{x}}(t_i)] \end{aligned} \quad (45b)$$

$$\tilde{\mathbf{x}}(t_{i-1}) = \tilde{\mathbf{x}}'(t_i) - \frac{a(t_i)\tilde{\mathbf{x}}'(t_i)}{N} + \frac{b(t_i)^2 s_{\theta'}(\tilde{\mathbf{x}}'(t_i), t)}{N} + \frac{b(t_i)\mathbf{z}_i}{\sqrt{N}} \quad (45c)$$

where  $\boldsymbol{\mathcal{P}}^\dagger(\boldsymbol{\Lambda})$  denotes the right inverse of  $\boldsymbol{\mathcal{P}}(\boldsymbol{\Lambda})$ . The process runs from  $t = 1$  backward to  $t = 0$  and draws approximate samples  $\tilde{\mathbf{x}}'(0)$  from  $p(\mathbf{x} \mid \mathbf{y})$  as the reconstructed images. This is especially useful in quantifying the uncertainty of the reconstructions by directly evaluating the mean and variance of the reconstructed images for the same measurement  $\mathbf{y}$ .



**FIGURE 9.** An illustration of the iterative sampling method in [1] for solving the inverse problems. The green arrows indicate the process of generating samples  $\tilde{y}(t_i)$  from the observation  $y$  according to (45a). The orange arrows show the process of combining the unconditional generative samples  $\tilde{x}(t_i)$  with  $\tilde{y}(t_i)$  according to (45b). The blue arrows indicate the reverse-time stochastic process in (45c).



# Summary and future outlook

- cryo-EM
  - more systematic quantitative comparisons to choices and hyperparameters
  - provide physical or biological meaning to the distance between the latent space of conformation variables
  - still lack a standardized measure or gold-standard task to evaluate how well a method is able to capture the heterogeneity
- score-based approaches
  - slow convergence