



# Diffusion models for protein design

Elodie Laine

Laboratory of Computational and Quantitative Biology

e-mail: [elodie.laine@sorbonne-universite.fr](mailto:elodie.laine@sorbonne-universite.fr)



@LaineElodie



CHARLES  
UNIVERSITY



SORBONNE  
UNIVERSITÉ



UNIVERSITÄT  
HEIDELBERG  
ZUKUNFT  
SEIT 1386



UNIVERSITY  
OF WARSAW



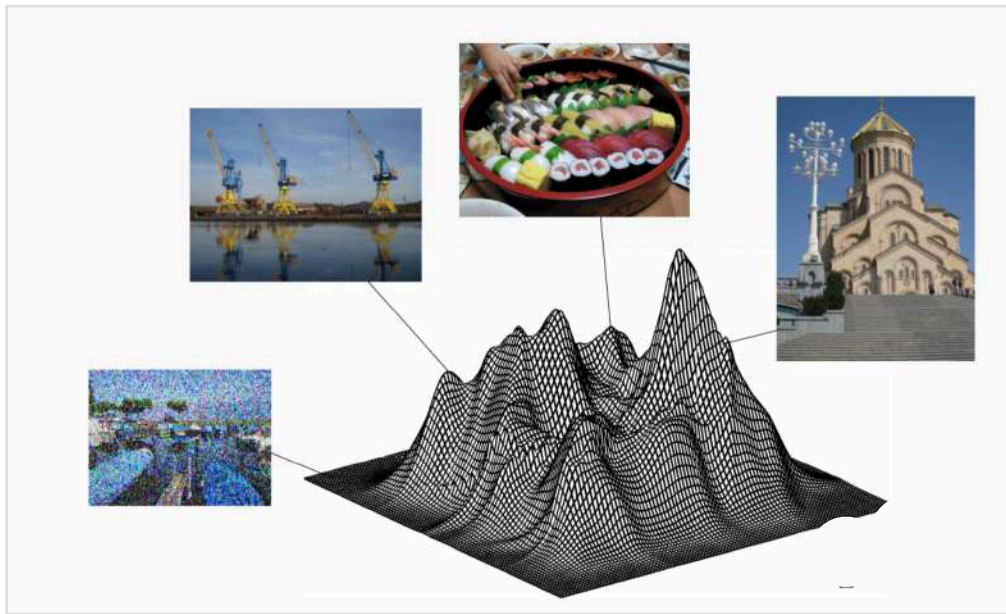
UNIVERSITÀ  
DEGLI STUDI  
DI MILANO



EUROPEAN  
UNIVERSITY  
ALLIANCE

# The problem of generative modelling

*Parametric density estimation: we want to learn the distribution underlying the training samples to generate new **plausible** samples that are similar but not identical.*



credits: Charles Deledalle

## What do we have?

A finite set of samples,  $\{x^{(n)}\}_{n=1}^N$ , from a true but unknown data distribution  $q_{data}(x)$ .

## What do we want?

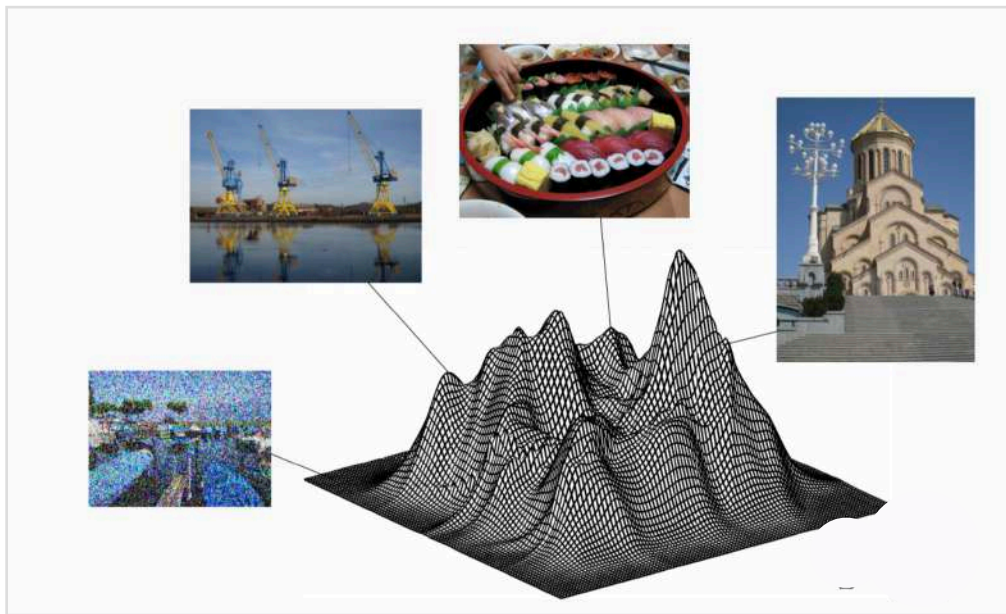
Estimate  $q_{data}(x)$  with a suitable model family  $p_{\theta}(x)$ , with unknown parameters  $\theta$ .

The problem boils down to maximising the average likelihood (w.r.t.  $\theta$ ) of all the samples under the model,

$$\theta^* = \arg \max_{\theta} \mathbb{E}_{x \sim q_{data}(x)} [\log p_{\theta}(x)] \approx \arg \max_{\theta} \frac{1}{N} \sum_{n=1}^N \log p_{\theta}(x^{(n)})$$

# The problem of generative modelling

*Parametric density estimation: we want to learn the distribution underlying the training samples to generate new **plausible** samples that are similar but not identical.*



credits: Charles Deledalle

## Is it simple?

No! The likelihood is intractable because one has to consider all the data points in the distribution to estimate any point's probability.

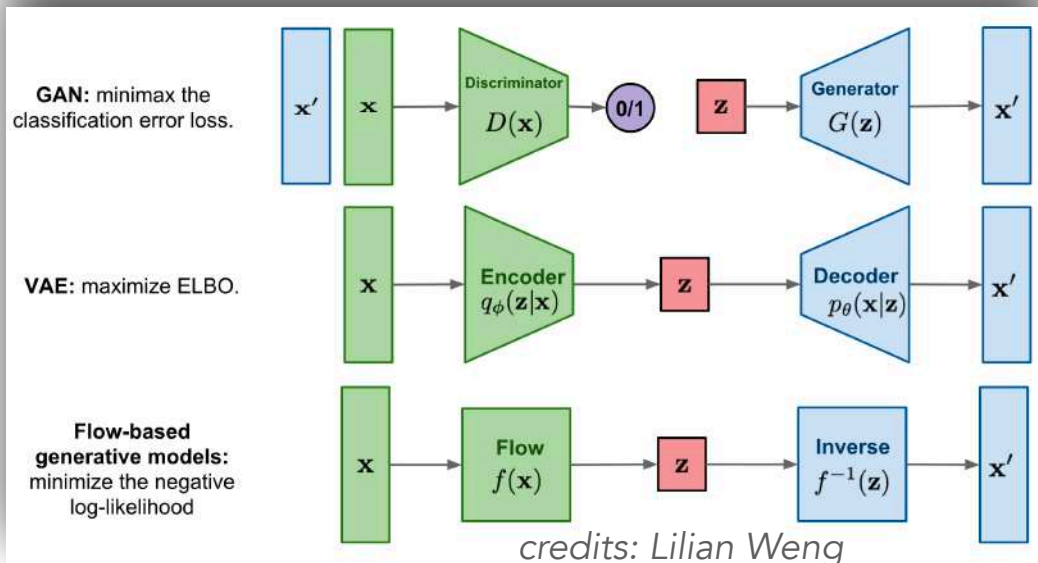
$$p_{\theta}(x) = \frac{\tilde{p}_{\theta}(x)}{\int_x \tilde{p}_{\theta}(x)}$$

Normalisation constant

## What can we do?

- Restrict the architecture to ensure tractability
- Rely on surrogate objectives
- Avoid evaluating probabilities (implicit models)

# Sampling is essential!



Sampling-centric generative models operate by transforming simple probabilistic densities.

$$x = f_\theta(z), \text{ where } z \sim \mathcal{N}(0, I)$$

Functional transformation (neural network)

Simple density (often standard normal)

In **diffusion models**, the transformation  $f_\theta$  is a sequence of invocations of a neural function, denoted as  $s_\theta$ , along with some additional computations, denoted as  $g(\cdot)$ .

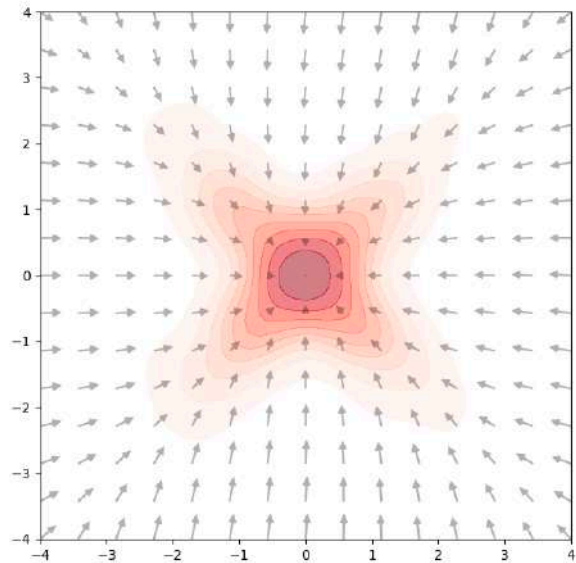
$$x = g_1(g_2(g_3(\cdots z \cdots, s_\theta), s_\theta), s_\theta), \text{ where } z \sim \mathcal{N}(0, I)$$

# Outline

- Why the score  $s_\theta$ ?
- Why diffusion? (backward)
- Why forward?
- Why so popular?

☐☐☐☐

# Why the score?



Vector field representing the score of the true distribution at any point

Diffusion models aim at estimating the (Stein) **score function**, defined as the **gradient of the log of the density of the data**,

$$\nabla_x \log q_{data}(x) \triangleq s(x)$$

It reflects the direction of steepest increase in log-likelihood at any given point in the data space.

At a point  $x$ , the **score** gives us the best direction to step into (with little step size  $\delta$ ) if we would like to see a point  $x'$  with slightly higher likelihood,

$$x' = x + \delta \cdot \nabla_x \log q_{data}(x) \Big|_{x=x}$$

**Why is it more convenient to estimate the score instead of the density?**

$\Rightarrow$  because the score **does not depend** on the normalising constant!

$$s_{\theta}(\mathbf{x}) = \nabla_{\mathbf{x}} \log p_{\theta}(\mathbf{x}) = \nabla_{\mathbf{x}} \tilde{p}_{\theta}(\mathbf{x}) - \underbrace{\nabla_{\mathbf{x}} \log \int_{\mathbf{x}} \tilde{p}_{\theta}(\mathbf{x})}_{=0} = \nabla_{\mathbf{x}} \tilde{p}_{\theta}.$$

# Outline

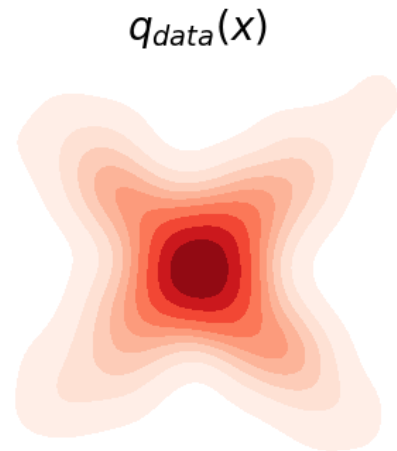
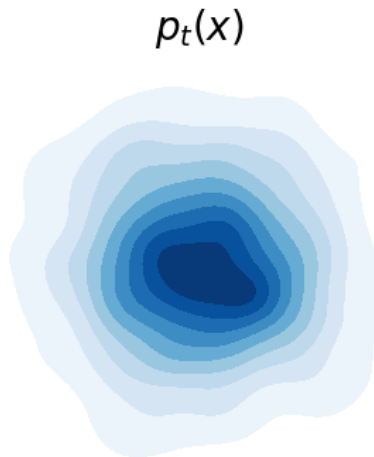
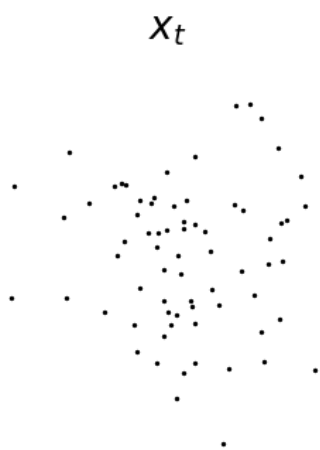
- Why the score  $s_\theta$ ?
- Why diffusion? (backward)
- Why forward?
- Why so popular?



# Sampling from the true score

$$x' = x + \delta \cdot \nabla_x \log q_{data}(x) \Big|_{x=x} \xrightarrow{\delta \rightarrow 0} dx = \nabla_x \log q_{data}(x) dt$$

Iteratively updating  $x$  from a normal distribution along the score **concentrate** the points in the most likely region.





# Langevin dynamics

French physicist *Paul Langevin* modelled the evolution of the positions of particles over time  $t$  under a potential energy field  $U(x)$  as,  $dx = -\nabla_x U(x)dt + \sqrt{2}dB_t$ .

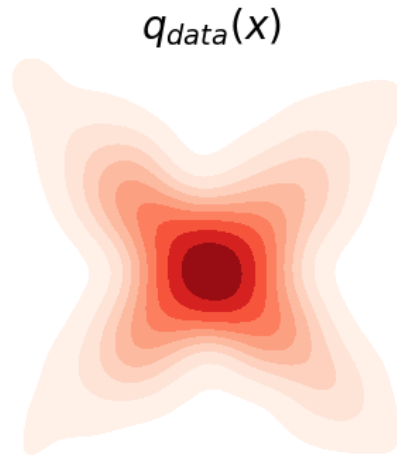
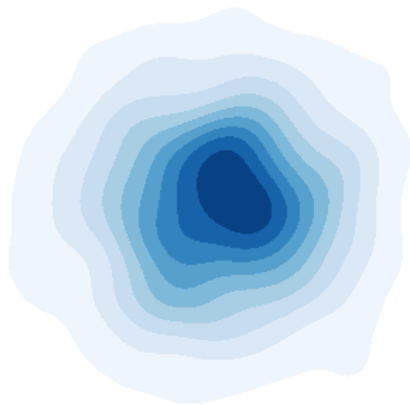
The energy is linked to the probability density,

$$q_{data}(x) = \frac{e^{-U(x)}}{Z} \quad \text{s.t.} \quad dx = \nabla_x \log q_{data}(x)dt + \sqrt{2}dB_t$$

$x_t$                        $p_t(x)$                        $q_{data}(x)$

Brownian motion (noise)

$$dB_t = \mathcal{N}(0, dt) \implies dB_t = \sqrt{dt} \cdot z, \text{ where } z \sim \mathcal{N}(0, I)$$



# Langevin dynamics

French physicist *Paul Langevin* modelled the evolution of the positions of particles over time  $t$  under a potential energy field  $U(x)$  as,  $dx = -\nabla_x U(x)dt + \sqrt{2}dB_t$ .

Brownian motion (noise)

$$dB_t = \mathcal{N}(0, dt) \implies dB_t = \sqrt{dt} \cdot z, \\ \text{where } z \sim \mathcal{N}(0, I)$$

The energy is linked to the probability density,

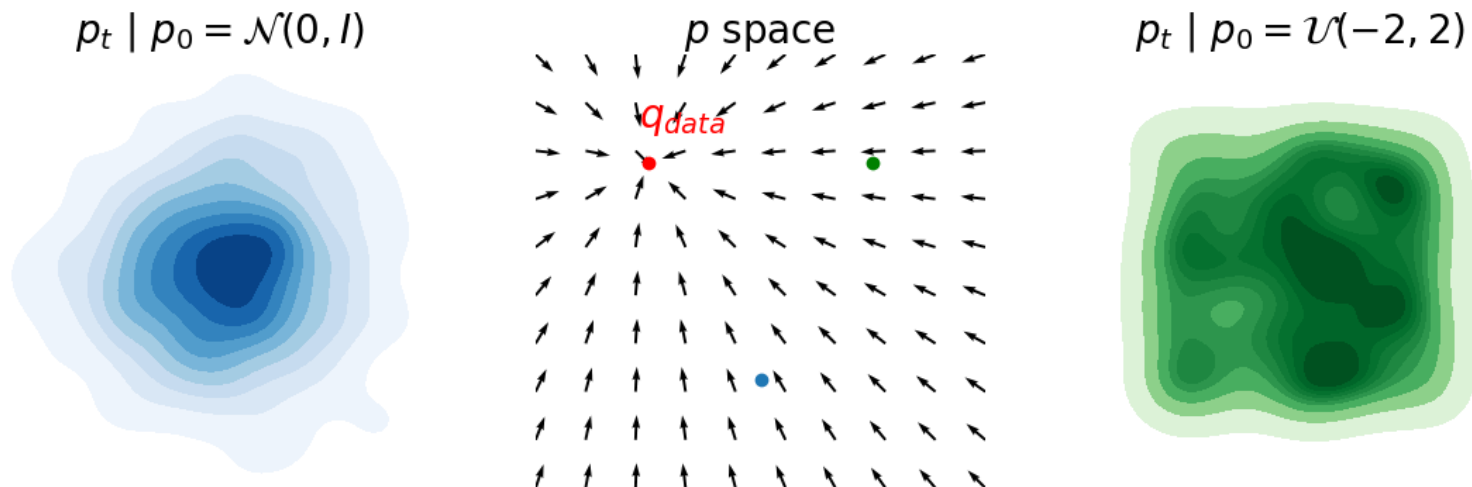
$$q_{data}(x) = \frac{e^{-U(x)}}{Z} \quad \text{s.t.} \quad dx = \nabla_x \log q_{data}(x)dt + \sqrt{2}dB_t$$

This stochastic process induces a time-varying distribution  $p_t(x)$  described by the *Fokker-Planck* equation,

$$\frac{\partial}{\partial t} p_t(x) = -\frac{\partial}{\partial x} \left[ p_t(x) \nabla_x \log q_{data}(x) \right] + \frac{1}{2} \frac{\partial^2}{\partial x^2} \left[ p_t(x) \sqrt{2} \right]$$

The probability distribution **provably converges** to the data distribution  $p_\infty(x) = q_{data}(x)$ .

# Illustration of the convergence



# Going backward

$q_{data}$

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



$x(0)$



$x(T)$

$$dx = \nabla_x \log q_{data}(x) dt + \sqrt{2} dB_t$$

**Backward process**

# Outline

- Why the score  $s_\theta$ ?
- Why diffusion? (backward)
- Why forward?
- Why so popular?



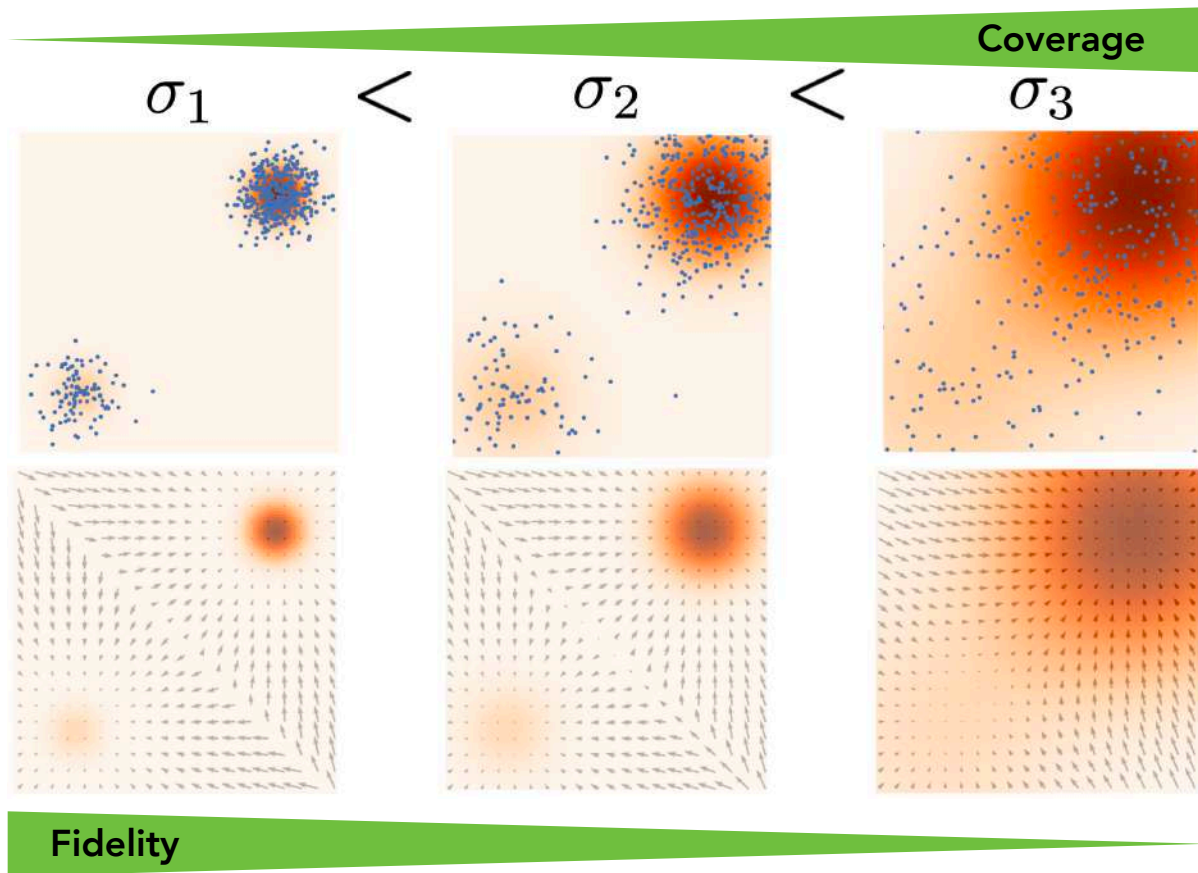
# Score matching

To estimate the score, we would ideally minimise a loss function that reflects the error to the true score.

$$J(\theta) = \frac{1}{2} \mathbb{E}_{x \sim q_{data}(x)} \left[ ||s_{\theta}(x) - \nabla_x \log q_{data}(x)||^2 \right]$$

**Problem:** we do not have access to the true score.

# Augmenting data with noise



Estimating empirical scores at different scales of noise, improves the accuracy of the score estimation!

credits: Yang Song

# Score matching

To estimate the score, we would ideally minimise a loss function that reflects the error to the true score.

$$J(\theta) = \frac{1}{2} \mathbb{E}_{x \sim q_{data}(x)} \left[ ||s_{\theta}(x) - \nabla_x \log q_{data}(x)||^2 \right]$$

**Problem:** we do not have access to the true score.

*Pascal Vincent* established in 2011 the connection between score matching and denoising autoencoders,

$$J_D(\theta) = \mathbb{E}_{x \sim q_{data}(x), \epsilon \sim \mathcal{N}(0, I)} \left[ \frac{1}{2} \left\| s_{\theta} \left( \underbrace{x + \sigma \epsilon}_{\tilde{x}} \right) - \left( -\frac{\epsilon}{\sigma} \right) \right\|^2 \right]$$

A connection between  
score matching and  
denoising  
autoencoders  
P. Vincent. Neural  
computation. MIT  
Press. 2011.

We estimate the score at  
a noisy version of  $x$

We compare with the noise added  
to  $x$  (which we know by design!)



# Going forward

## Forward process

$$dx = \nabla_x \log \mathcal{N}(0, I) dt + \sqrt{2} dB_t$$

$$dx = -x dt + \sqrt{2} dz$$

$q_{data}$    $\mathcal{N}(0, I)$

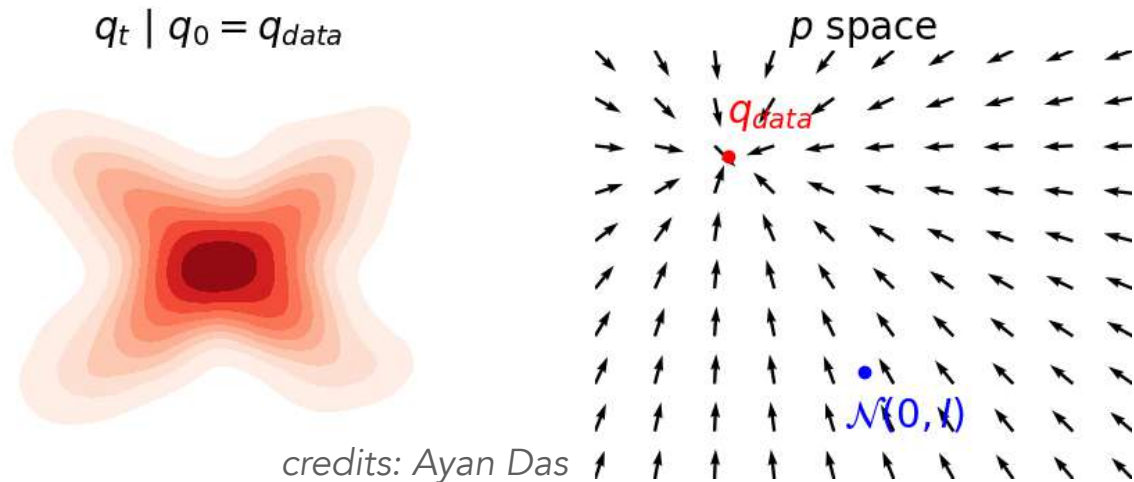


$x(0)$    $x(T)$

$$dx = \nabla_x \log q_{data}(x) dt + \sqrt{2} dB_t$$

## Backward process

# Illustration of the convergence



$$x_{t+dt} = (1 - dt)x_t + \sqrt{2dt} z$$



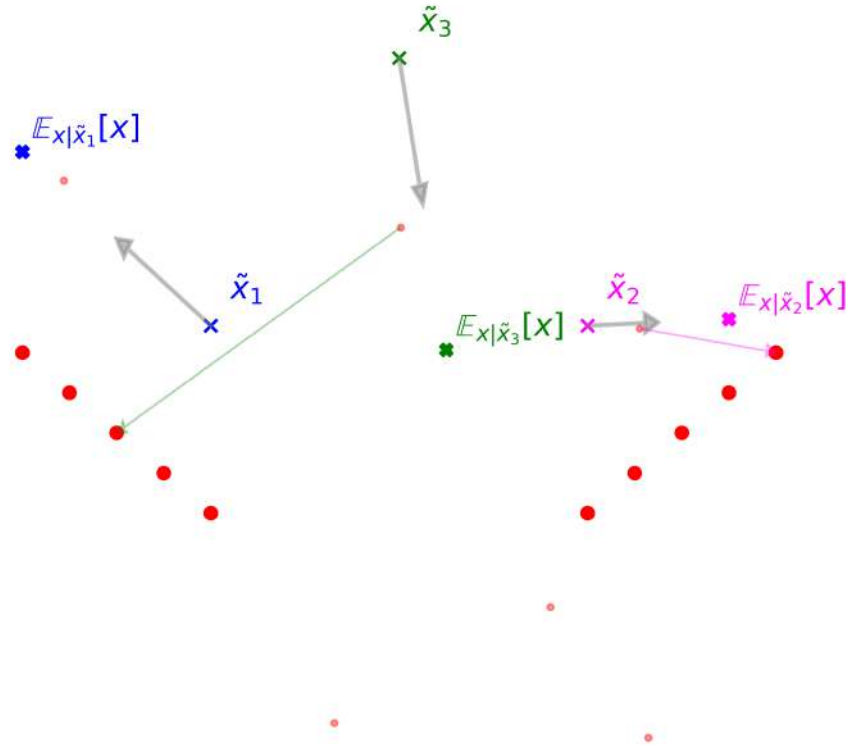
Convergence is guaranteed at infinite time. We may re-parametrize time to escalate the forward process (noise schedule).

Since the process is fully known, one can jump to any time without going through the sequential process.

$$x_{t+2dt} = (1 - dt)x_{t+dt} + \sqrt{2dt} z_2 = (1 - dt)[(1 - dt)x_t + \sqrt{2dt} z_1] + \sqrt{2dt} z_2$$

$$\Rightarrow x_{t+2dt} \sim \mathcal{N}((1 - 2 \cdot dt)x_t, 2 \cdot 2dtI)$$

# Illustration of the denoising



# Outline

- Why the score  $s_\theta$ ?
- Why diffusion? (backward)
- Why forward?
- Why so popular?



# Bayesian inversion

According to Bayes theorem,  $\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t | y) = \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) + \nabla_{\mathbf{x}_t} \log p(y | \mathbf{x}_t)$

Score of the unconditional denoising

Score of a classifier or regressor of a given property  $y$

Generalised to many properties,

$$\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t | y_1, y_2, \dots, y_M) = \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) + \sum_{i=1}^M \nabla_{\mathbf{x}_t} \log p(y_i | \mathbf{x}_t)$$

Provided classifiers able to predict the properties based on noisy samples, the diffusion model can be used for conditional generation **without re-training**.

Ingraham, John B., et al. "Illuminating protein space with a programmable generative model." *Nature* 623.7989 (2023): 1070-1078.



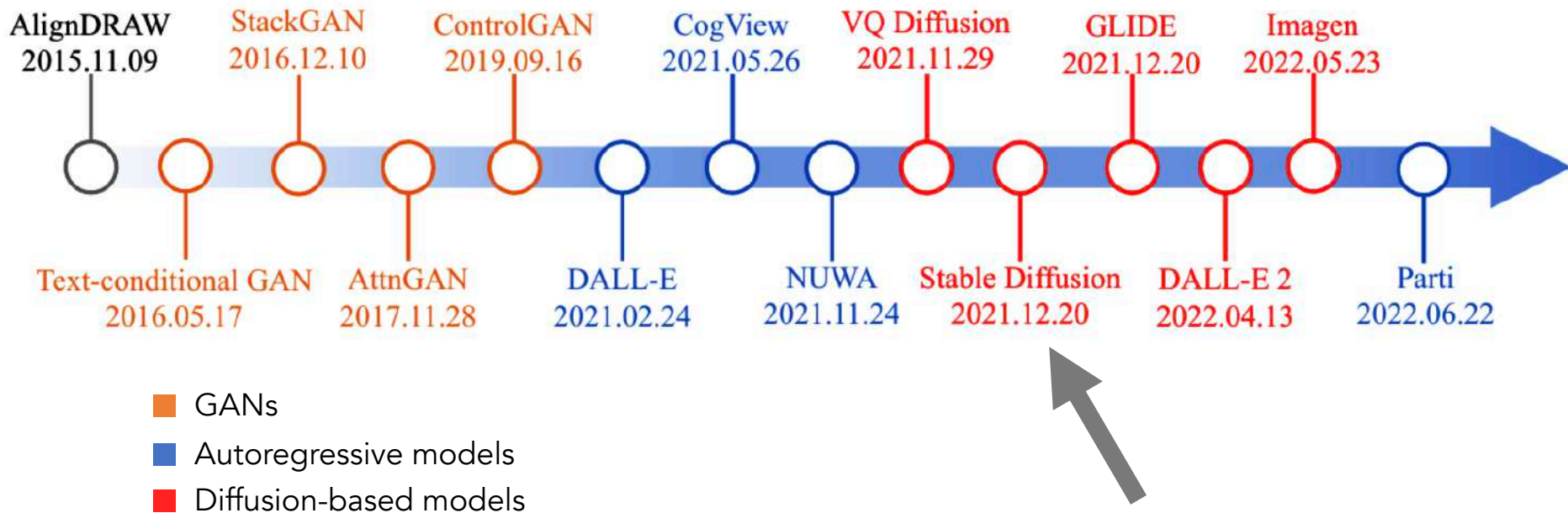
# Outline

- Why the score  $s_\theta$ ?
- Why diffusion? (backward)
- Why forward?
- Why so popular?



- **From text to images (Stable Diffusion)**
- Generating proteins
  - ▶ RFDiffusion
  - ▶ CHROMA
  - ▶ AlphaFold3...

# Text-to-Image architectures





# Stable diffusion

## Text-to-Image Synthesis on LAION. 1.45B Model.

'A street sign that reads  
"Latent Diffusion" '

'A zombie in the  
style of Picasso'

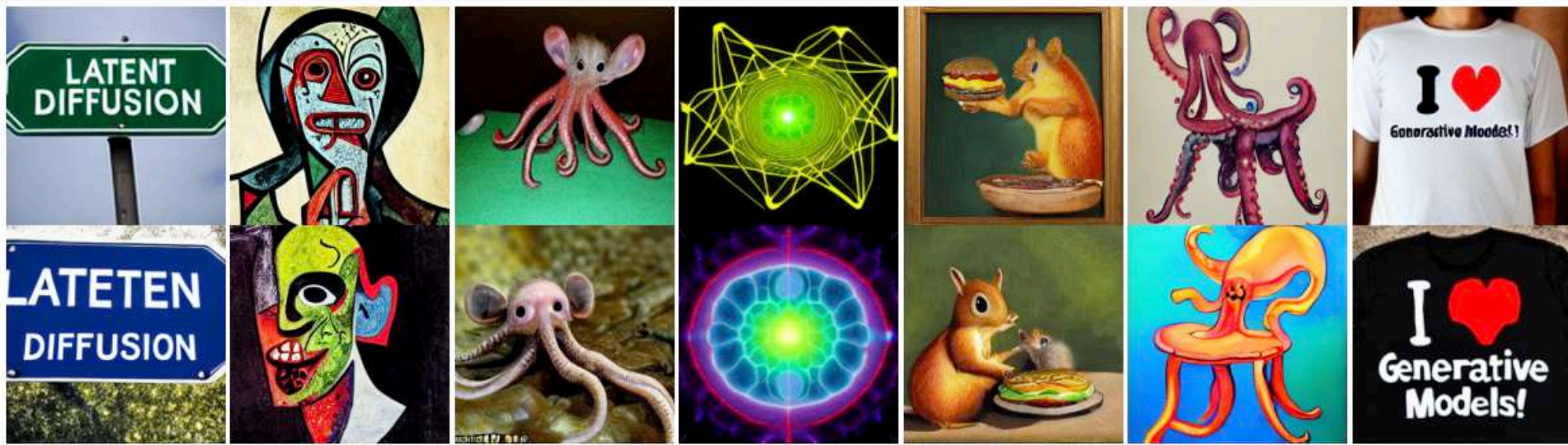
'An image of an animal  
half mouse half octopus'

'An illustration of a slightly  
conscious neural network'

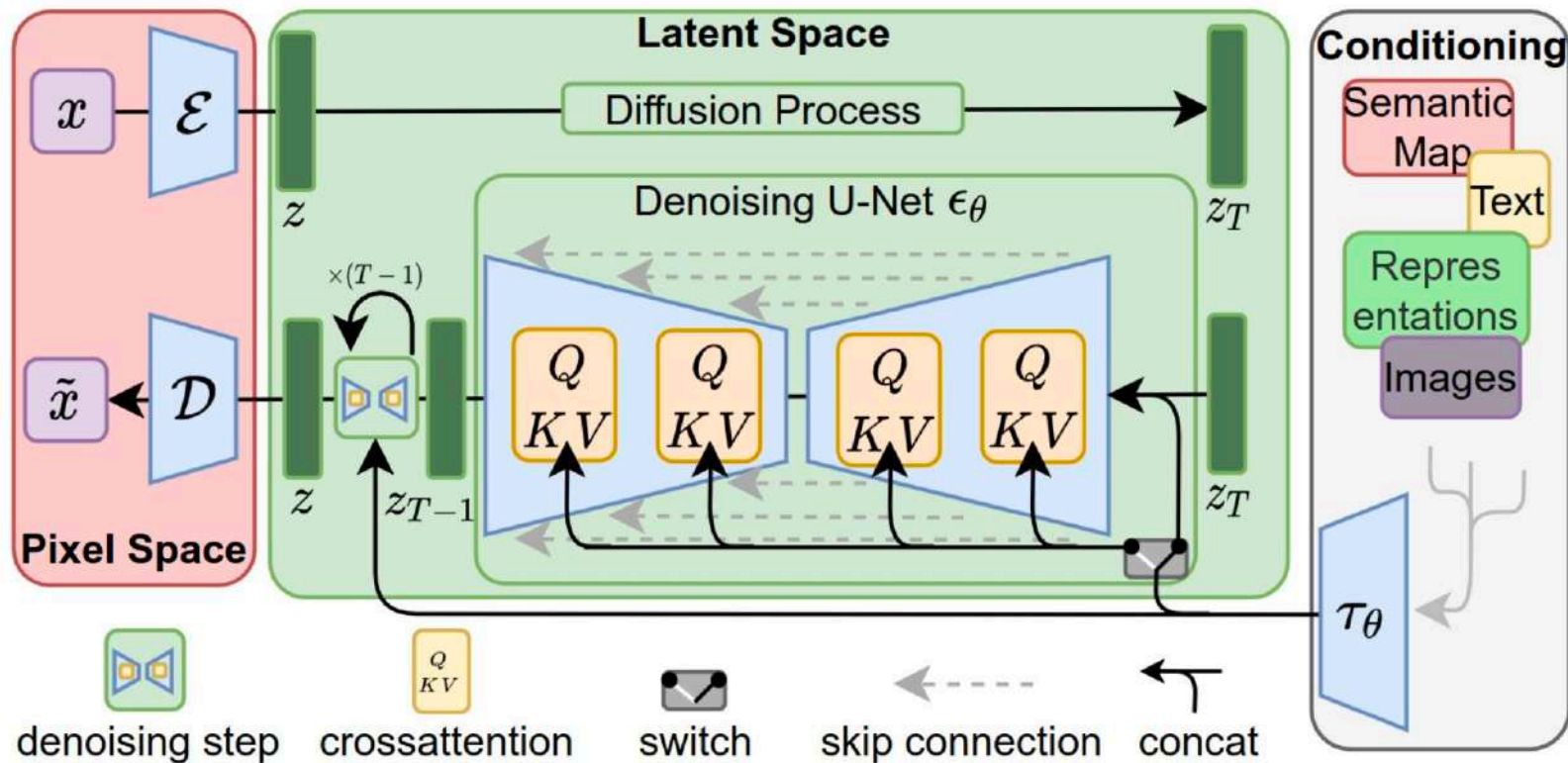
'A painting of a  
squirrel eating a burger'

'A watercolor painting of a  
chair that looks like an octopus'

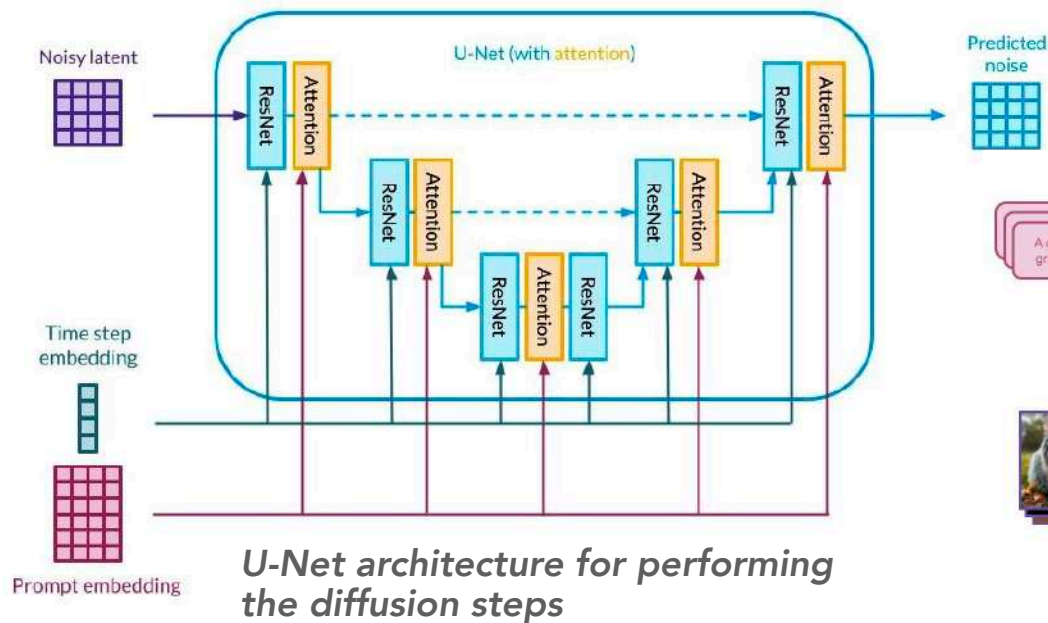
'A shirt with the inscription:  
"I love generative models!" '



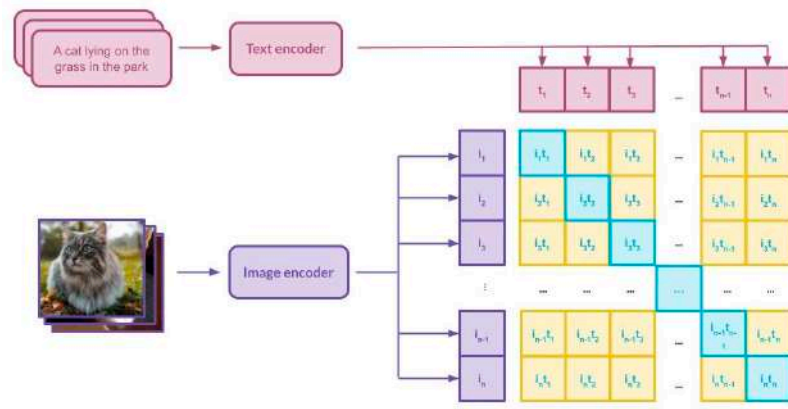
# Stable diffusion



# Stable diffusion



CLIP model for encoding the conditional prompt

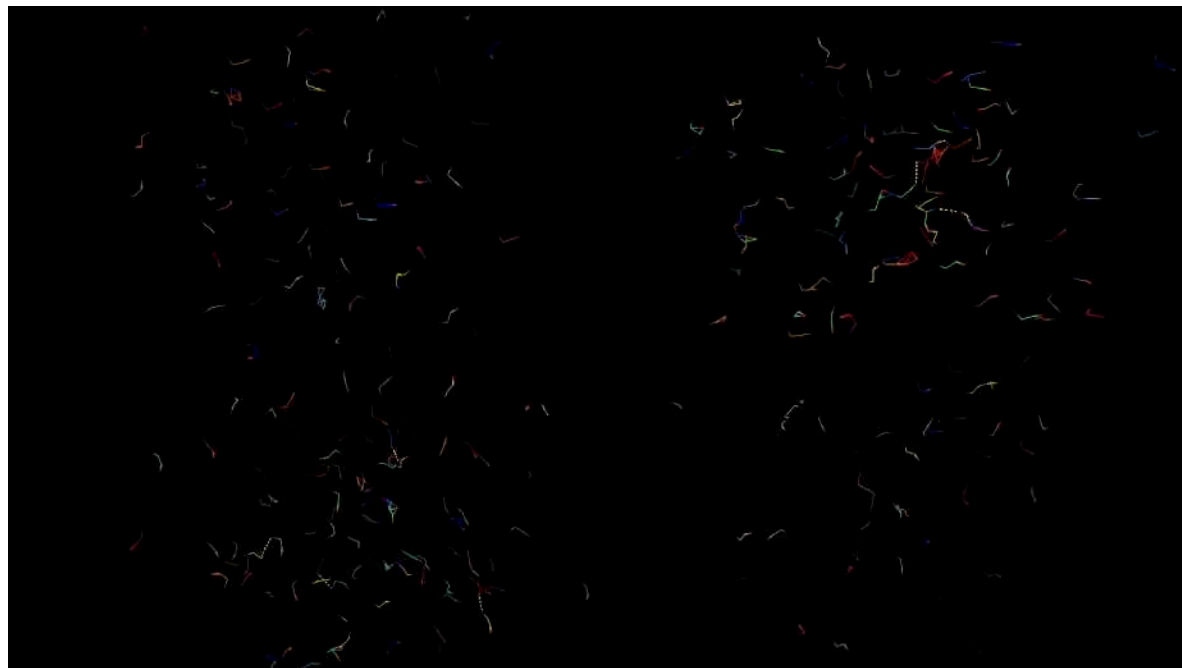


$$L_{LDM} = \mathbb{E}_{\epsilon(x), y, \epsilon \sim \mathcal{N}(0,1), t} \left[ \left\| \epsilon - \epsilon_{\theta}(z_t, t, \tau_{\theta}(y)) \right\|_2^2 \right]$$

credits: Jarosław Kochanowicz, Maciej Domagała, Dawid Stachowiak and Krzysztof Dziedzic

- From text to images (Stable Diffusion)
- **Generating proteins**
  - ▶ RFDiffusion
  - ▶ CHROMA
  - ▶ AlphaFold3...

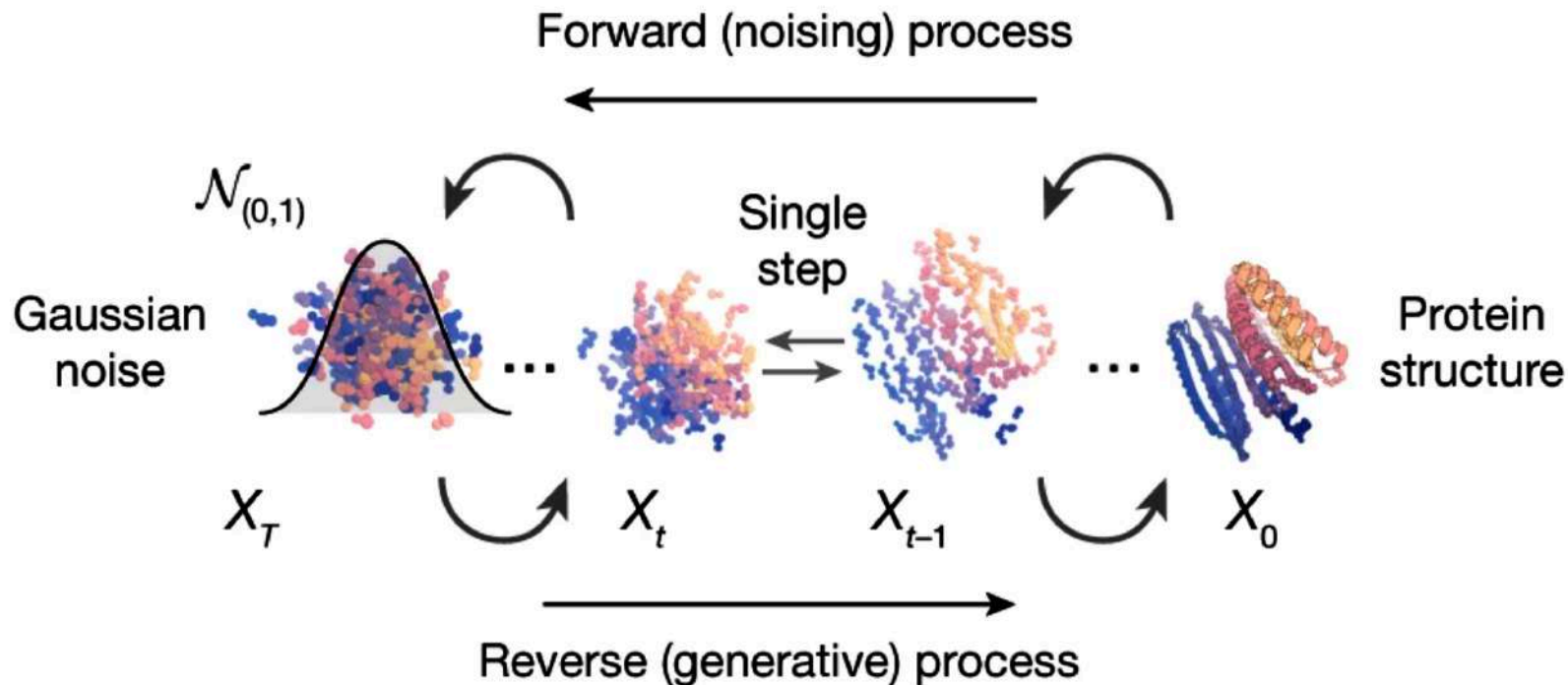
# Applications



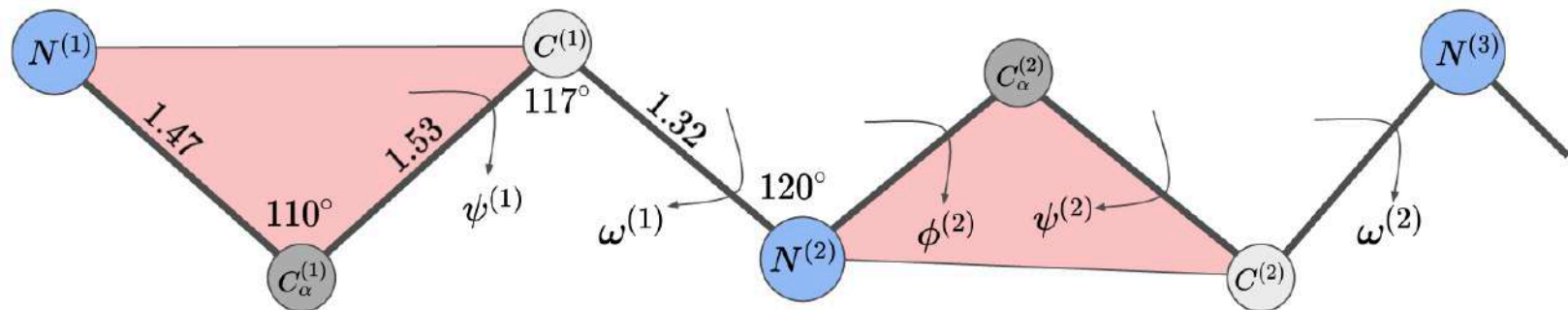
Anand, Namrata, and Tudor Achim. "Protein structure and sequence generation with equivariant denoising diffusion probabilistic models." *arXiv preprint arXiv:2205.15019* (2022).



# RosettaFold *Diffusion*



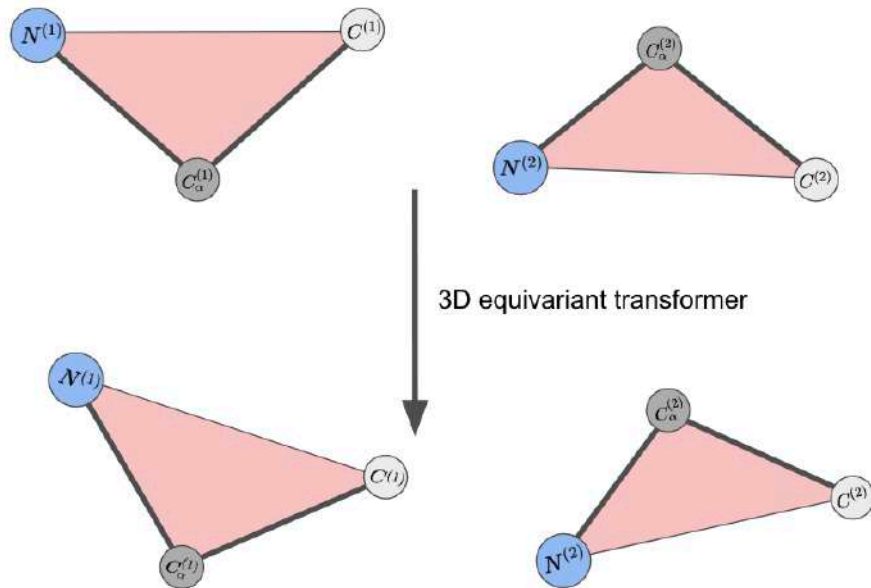
# Protein backbone



Bond lengths (in Å) and angles (in degree) between neighbouring atoms are fixed.  
The torsion angles  $\psi^{(1)}$ ,  $\omega^{(1)}$ ,  $\phi^{(2)}$ ,  $\psi^{(2)}$ ,  $\omega^{(2)}$  are the degrees of freedom.

credits: Justas. Dauparas

# Protein backbone representation



Euclidian transform from the **local frame** to the global reference frame

$$x_l = (r_l, z_l)$$

Coordinates of the  $l$ th residue

3x3 rotation matrix

Translation (coordinates of the  $C_\alpha$ )

credits: Justas. Dauparas



# Independent diffusion processes

$$x_l = (r_l, z_l)$$

$$q(z^{(t)} | z^{(0)}) = \mathcal{N}(z^{(t)}; \sqrt{\bar{\alpha}^{(t)}} z^{(0)}, (1 - \bar{\alpha}^{(t)}) I_3)$$
$$\bar{\alpha}^{(t)} = \prod_{s=1}^t \alpha^{(s)}, \quad \alpha^{(t)} = 1 - \beta^{(t)}, \quad \beta^{(1)}, \beta^{(2)}, \dots, \beta^{(T)} \in [0, 1]$$

$$p(z^{(t-1)} | x^{(t)}) = \mathcal{N}(z^{(t-1)}; \hat{\mu}(x^{(t)}), \beta^{(t)} I_3),$$
$$\text{with } \hat{\mu}(x^{(t)}) = \frac{\sqrt{\bar{\alpha}^{(t-1)}} \beta^{(t)}}{1 - \bar{\alpha}^{(t)}} \hat{z}^{(0)}(x^{(t)}) + \frac{\sqrt{\alpha^{(t)}} (1 - \bar{\alpha}^{(t-1)})}{1 - \bar{\alpha}^{(t)}} z^{(t)}$$

Trippe, Brian L., et al. "Diffusion probabilistic modeling of protein backbones in 3d for the motif-scaffolding problem." *arXiv preprint arXiv:2206.04119* (2022).

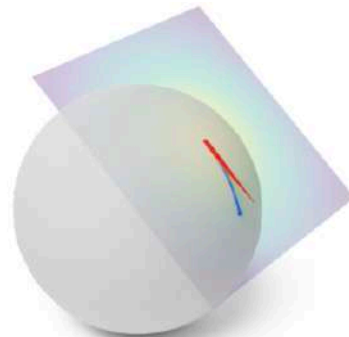
**Check out Julien's  
notebook!**

# Independent diffusion processes

$$x_l = (r_l, z_l)$$

$$q(r^{(t)} \mid r^{(0)}) = \mathcal{IG}_{SO(3)}(r^{(t)}; r^{(0)}, \sigma_t^2)$$

Which is the Isotropic Gaussian distribution on  $SO(3)$ .



(a) A single step of a Geodesic Random Walk.



(b) Many steps yield an approximate trajectory.

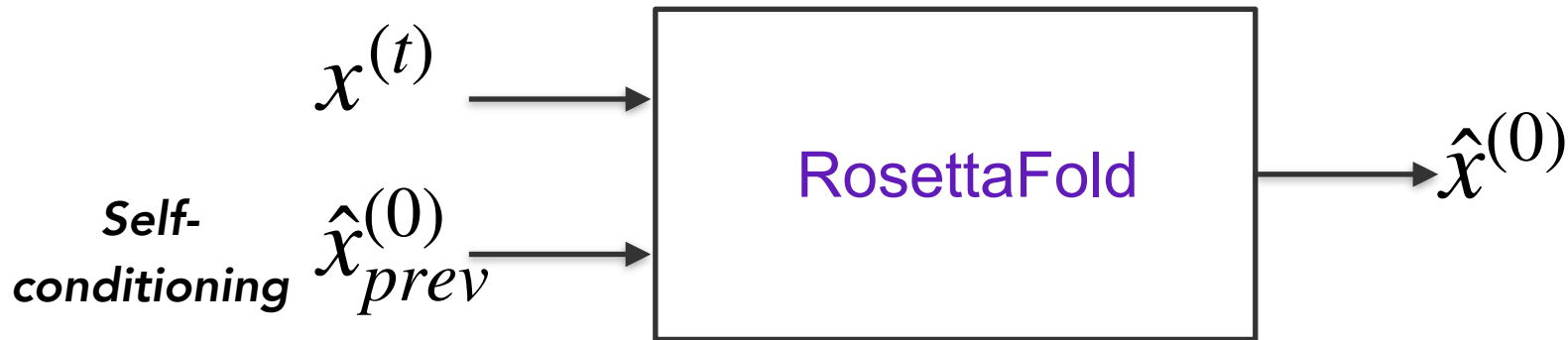
De Bortoli, Valentin, et al.  
"Riemannian score-based  
generative  
modelling." *Advances in Neural  
Information Processing  
Systems* 35 (2022): 2406-2422.

$$r^{(t-1)} = \exp_{r^{(t)}} \left\{ (\sigma_t^2 - \sigma_{t-1}^2) \nabla_{r^{(t)}} \log q(x^{(t)}) + \sqrt{\sigma_t^2 - \sigma_{t-1}^2} \sum_{d=1}^3 \epsilon_d r^{(t)} f_d \right\}$$

with  $r^{(t)} f_1, r^{(t)} f_2, r^{(t)} f_3$  the orthonormal basis of the tangent space of  $r^{(t)}$

$$\nabla_{r^{(t)}} \log q(x^{(t)}) = \nabla_{r^{(t)}} \log \mathcal{IG}_{SO(3)}(r^{(t)}; \hat{r}^{(0)}, \sigma_t^2)$$

# Denoising structures



**Loss**  $\text{MSE}_{\text{Frame}} = \frac{1}{T} \sum_{t=1}^T \mathbb{E}_q[d_{\text{frame}}(x^{(0)}, \hat{x}^{(0)}(x^{(t)}))^2],$

$$d_{\text{frame}}(x^{(0)}, \hat{x}^{(0)}) = \sqrt{\frac{1}{L} \sum_{l=1}^L \|z_l^{(0)} - \hat{z}_l^{(0)}\|_2^2 + \|I_3 - \hat{r}_l^{(0)\top} r_l^{(0)}\|_F^2},$$

# Training procedure

---

## Algorithm 3 RFdiffusion Training

---

```
1: function FORWARDNOISE( $x^{(0)}, t$ )
2:   for  $l = 1, \dots, L$  do
3:      $(r_l^{(0)}, z_l^{(0)}) = x_l^{(0)}$ 
4:      $z_l^{(t)} \sim \mathcal{N}(\sqrt{\bar{\alpha}^{(t)}} z_l^{(0)}, (1 - \bar{\alpha}^{(t)}) I_3)$ 
5:      $r_l^{(t)} \sim \mathcal{IG}_{SO(3)}(r_l^{(0)}, \sigma_t^2)$ 
6:      $x_l^{(t)} = (r_l^{(t)}, z_l^{(t)})$ 
7:   end for
8:   return  $x^{(t)}$ 
9: end function
```

<https://github.com/RosettaCommons/RFdiffusion>

Watson, Joseph L., et al. "De novo design of protein structure and function with RFdiffusion." *Nature* 620.7976 (2023): 1089-1100.

```
11: function TRAIN
12:   while not converged do
13:      $x^{(0)} \sim \text{TrainingSet}$ 
14:      $t \sim \text{Uniform}(\{1, \dots, T\})$ 
15:     if  $\text{Uniform}(0, 1.0) < 0.5$  or  $t = T$  then
16:        $\triangleright$  Train step without self-conditioning
17:        $x^{(t)} = \text{ForwardNoise}(x^{(0)}, t)$ 
18:        $\hat{x}_{\text{prev.}}^{(0)} = \vec{0}$ 
19:     else
20:        $\triangleright$  Train step with self-conditioning
21:        $x^{(t+1)} = \text{ForwardNoise}(x^{(0)}, t + 1)$ 
22:        $x^{(t)} = \text{ReverseStep}(x^{(t+1)}, x^{(0)})$ 
23:
24:        $\triangleright$  Compute self-conditioning input
25:        $\hat{x}_{\text{prev.}}^{(0)} = \text{RFdiffusion}(x^{(t+1)}, \vec{0})$ 
26:        $\hat{x}_{\text{prev.}}^{(0)} = \text{StopGradient}(\hat{x}_{\text{prev.}}^{(0)})$ 
27:     end if
28:      $\hat{x}^{(0)} = \text{RFdiffusion}(x^{(t)}, \hat{x}_{\text{prev.}}^{(0)})$ 
29:     Take gradient step on  $d_{\text{frame}}(x^{(0)}, \hat{x}^{(0)})^2$ 
30:   end while
31: end function
```

# Inference

---

## Algorithm 2 RFdiffusion generation

---

```

1: function SAMPLEREFERENCE(L)
2:   ▷ Random initial structure for  $L$  residues
3:   for  $l = 1, \dots, L$  do
4:      $r_l^{(T)} \sim \text{Uniform}(SO(3))$ 
5:      $z_l^{(T)} \sim \mathcal{N}(0, I_3)$ 
6:      $x_l^{(T)} = (r_l^{(T)}, x_l^{(T)})$ 
7:   end for
8: return  $x^{(T)}$ 
9: end function

28: function SAMPLE(L)
29:   ▷ RFdiffusion generation of  $L$ -residue backbone structure
30:    $x^{(T)} = \text{SampleReference}(L)$ 
31:    $\hat{x}_{\text{prev.}}^{(0)} = \vec{0}$ 
32:   for  $t = T, \dots, 1$  do
33:      $\hat{x}^{(0)} = \text{RFDIFFUSION}(x^{(t)}, \hat{x}_{\text{prev.}}^{(0)})$ 
34:      $x^{(t-1)} = \text{REVERSESTEP}(x^{(t)}, \hat{x}^{(0)})$ 
35:      $\hat{x}_{\text{prev.}}^{(0)} = \hat{x}^{(0)}$ 
36:   end for
37: return  $\hat{x}^{(0)}$ 
38: end function

```

```

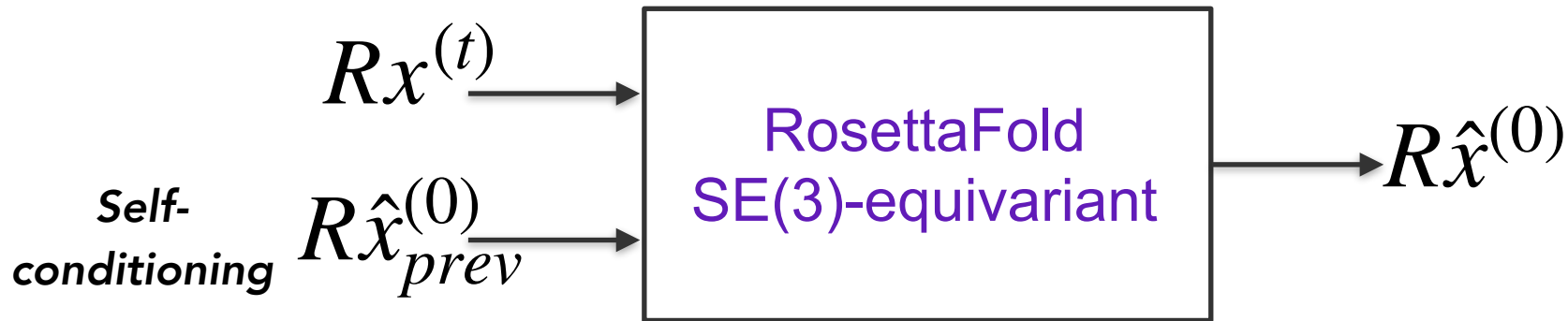
11: function REVERSESTEP( $x^{(t)}, \hat{x}^{(0)}$ )
12:   ▷ One step of reverse diffusion
13:   for  $l = 1, \dots, L$  do
14:      $(r_l^{(t)}, z_l^{(t)}) = x_l^{(t)}$ 
15:      $(\hat{r}_l^{(0)}, \hat{z}_l^{(0)}) = \hat{x}_l^{(0)}$ 
16:     ▷ Update translations
17:      $z_l^{(t-1)} \sim \mathcal{N}(\frac{\sqrt{\bar{\alpha}^{(t-1)}}\beta^{(t)}}{1-\bar{\alpha}^{(t)}}\hat{z}_l^{(0)} + \frac{\sqrt{\alpha^{(t)}(1-\bar{\alpha}^{(t-1)})}}{1-\bar{\alpha}^{(t)}}z_l^{(t)}, \beta^{(t)}I_3)$ 
18:
19:     ▷ Update rotations
20:      $s_l = \text{ROTATIONSCOREAPPROXIMATION}(r_l^{(t)}, \hat{r}_l^{(0)}, \sigma_t^2)$ 
21:      $\epsilon_{l,1}, \epsilon_{l,2}, \epsilon_{l,3} \stackrel{iid}{\sim} \mathcal{N}(0, 1)$ 
22:      $r_l^{(t-1)} = r_l^{(t)} \exp_{I_3} \left\{ (\sigma_t^2 - \sigma_{t-1}^2) r_l^{(t)\top} s_l + \sqrt{\sigma_t^2 - \sigma_{t-1}^2} \sum_{d=1}^3 \epsilon_{l,d} f_d \right\}$ 
23:      $x_l^{(t-1)} = (r_l^{(t-1)}, z_l^{(t-1)})$ 
24:   end for
25: return  $x^{(t-1)}$ 
26: end function

```

<https://github.com/RosettaCommons/RFdiffusion>

Watson, Joseph L., et al. "De novo design of protein structure and function with RFdiffusion." *Nature* 620.7976 (2023): 1089-1100.

# Invariance & equivariance in 3D



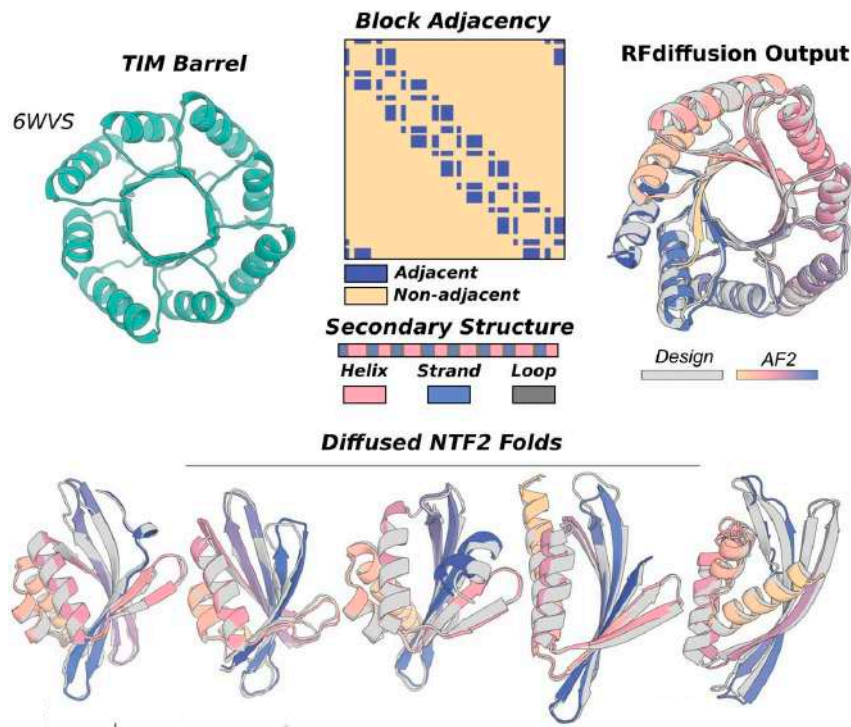
**Loss**  $\text{MSE}_{\text{Frame}} = \frac{1}{T} \sum_{t=1}^T \mathbb{E}_q[d_{\text{frame}}(x^{(0)}, \hat{x}^{(0)}(x^{(t)}))^2],$

Depends on the  
choice of the global  
reference frame!

$$d_{\text{frame}}(x^{(0)}, \hat{x}^{(0)}) = \sqrt{\frac{1}{L} \sum_{l=1}^L \|z_l^{(0)} - \hat{z}_l^{(0)}\|_2^2 + \|I_3 - \hat{r}_l^{(0)\top} r_l^{(0)}\|_F^2},$$



# Fold-specific generation

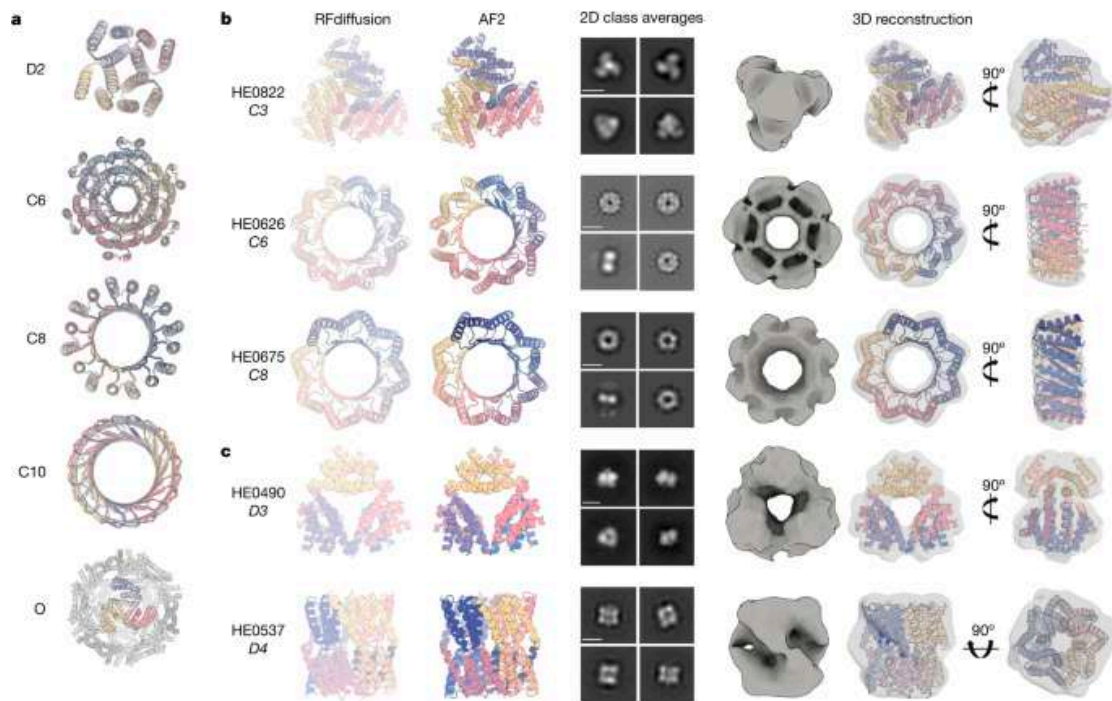


**RFdiffusion** is fine-tuned with extended input features specifying

- the secondary structure type of each residue,
- the adjacency matrix of the secondary structure blocks.

AlphaFold (AF2) is used as external validation of the plausibility or realism of the produced structures. The designs are diverse and closely recapitulated by AF2.

# Symmetric oligomers

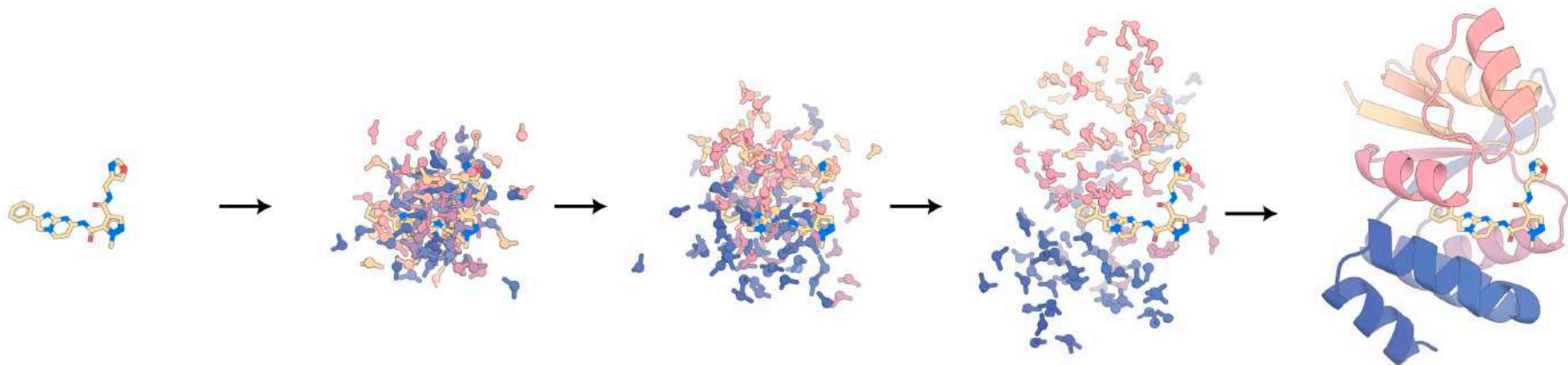


**RFDiffusion** is symmetry preserving thanks to its rotation and permutation invariance.

Explicit symmetrisation is still applied at each demonising step.

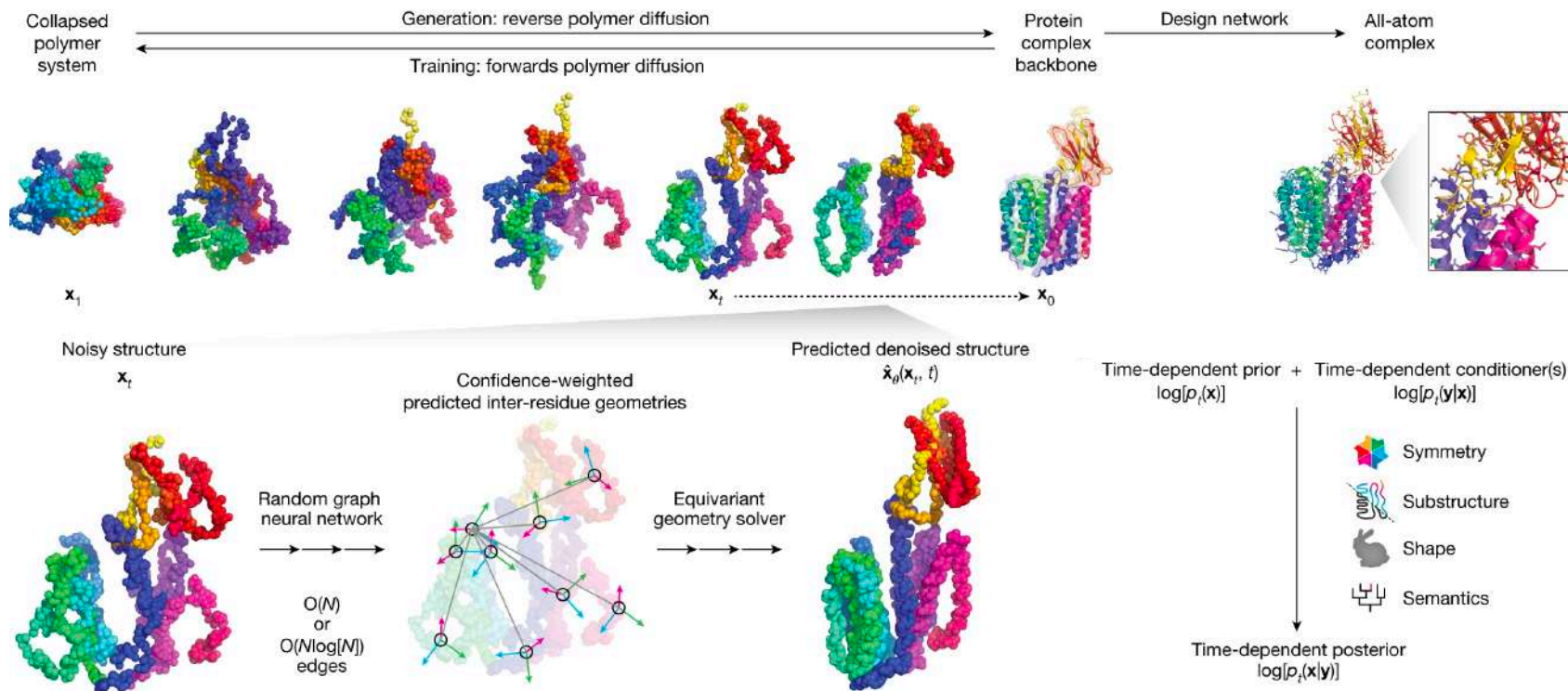


# Small molecule binder design

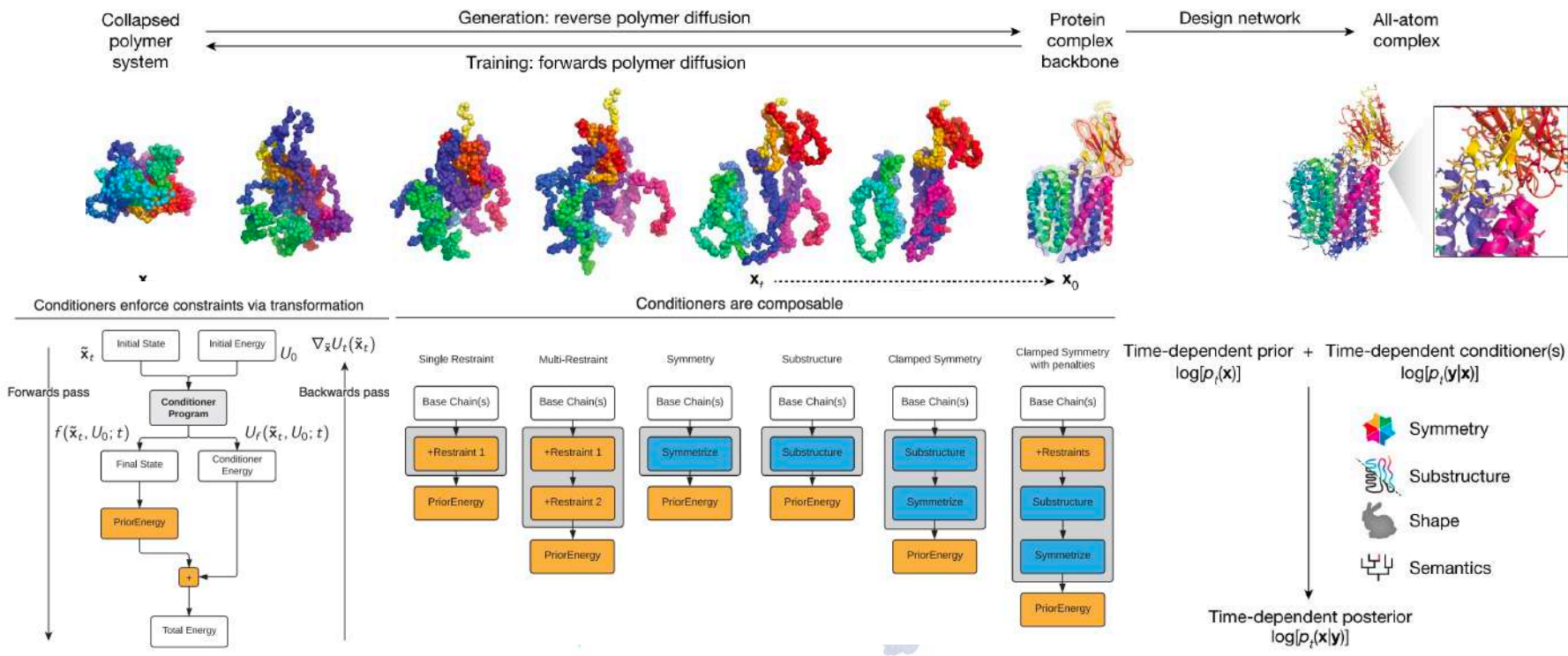


**RFDiffusion All-Atom** is explicitly trained in a conditioned way by demonising proteins bound to small molecules. The model is additionally trained to generated scaffolds conditionally on fixed motifs (with or without small molecules) by setting the time to zero for the motifs and maintaining them in original version (without any noise).

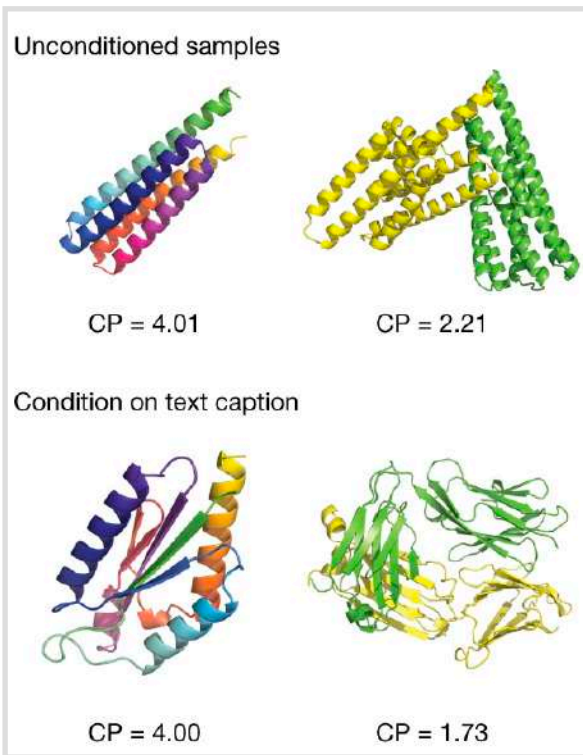
# CHROMA



# CHROMA



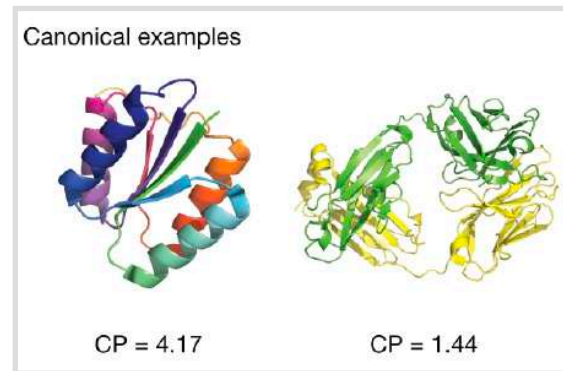
# From text to protein structures



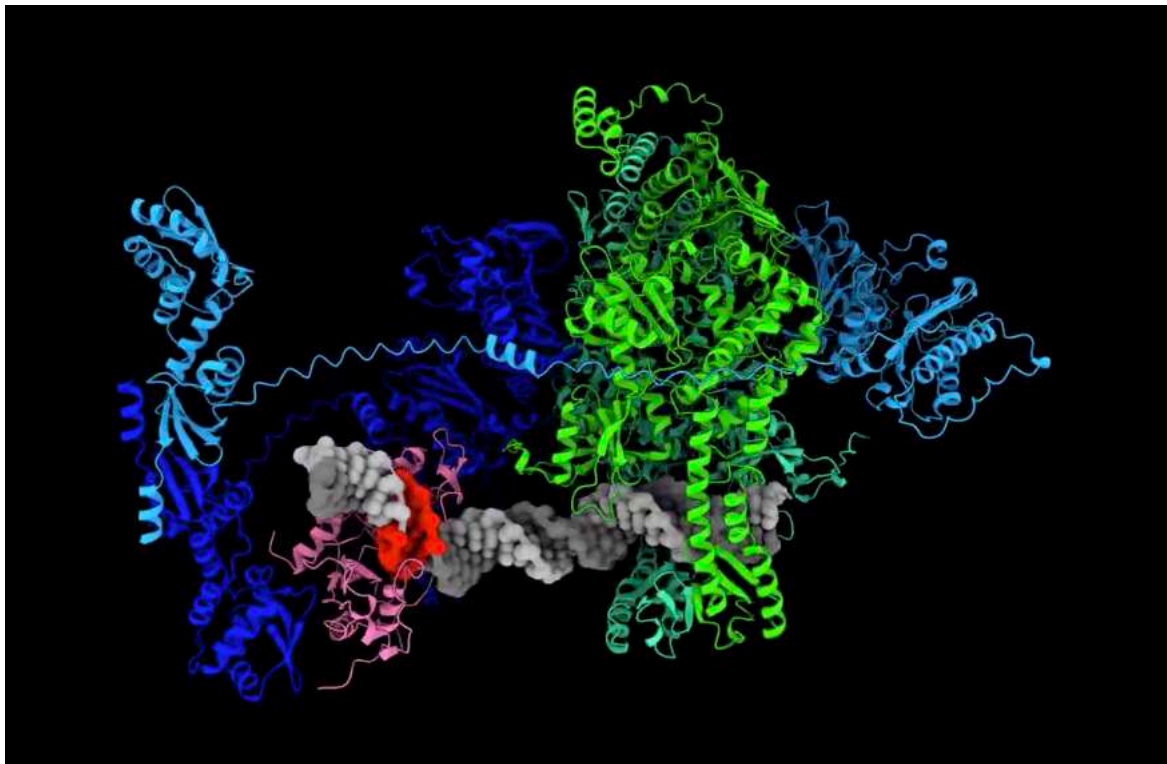
Left: *crystal structure of a Rossman fold*

Right: *crystal structure of a Fab antibody fragment*

Fine-tuning a multi-label predictor to bias a pretrained large language model into a structure caption predictor can enable natural language conditioning.



# AlphaFold 3



By Jan Kosinski:

MutS protein (green) recognizes base mismatches in DNA (gray), after which it slides off to find MutL matchmaker protein (blue), which positions and activates MutH nuclease (pink) to nick a hemimethylated GATC sequence (red) to activate removal and repair of the mismatched strand.

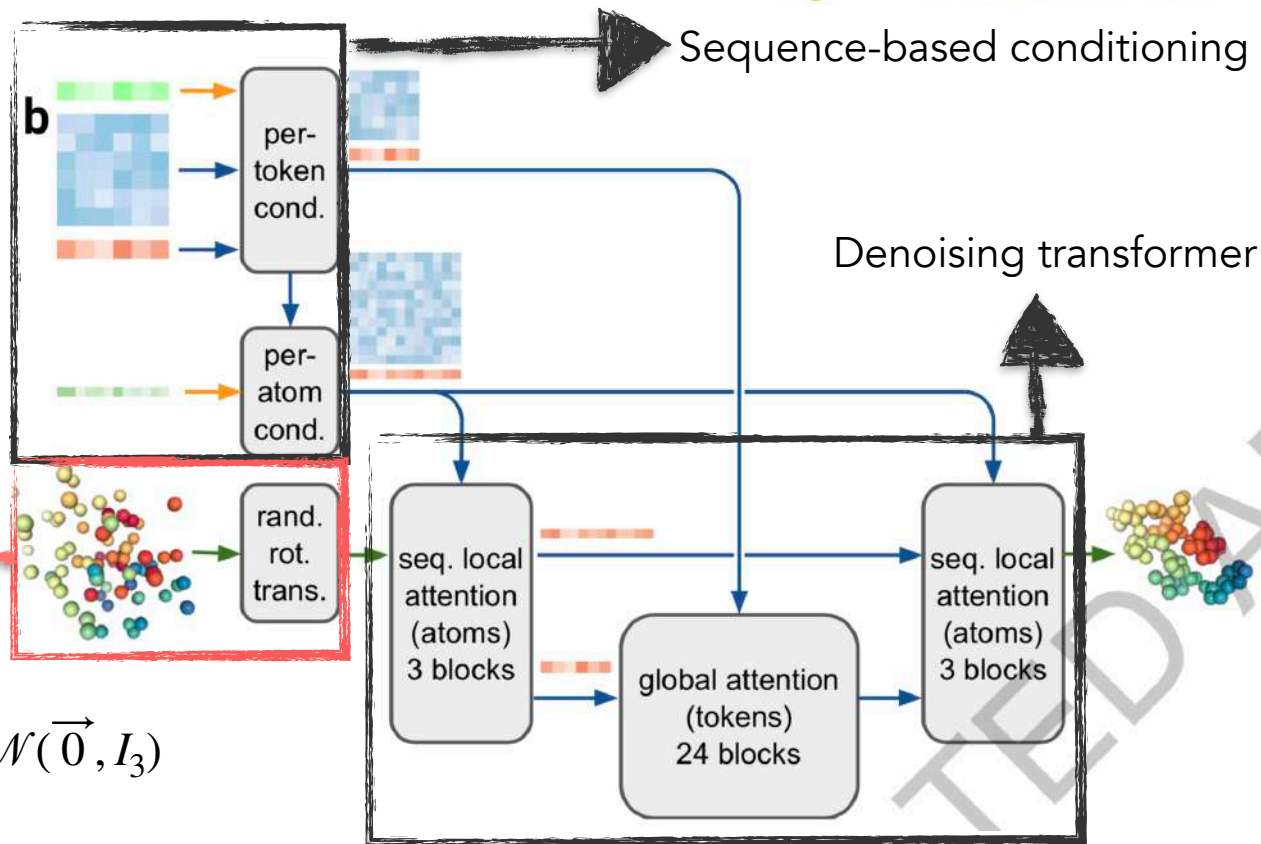


# AlphaFold 3 diffusion module

**No more frames!** The network operates directly on raw coordinates.

**Data augmentation:** get 48 samples from 1 through random rot. & trans.

$$\vec{x}_l^{aug} = R \cdot \vec{x}_l - \frac{1}{L} \sum_l \vec{x}_l + \mathcal{N}(\vec{0}, I_3)$$



# AlphaFold 3 diffusion process

---

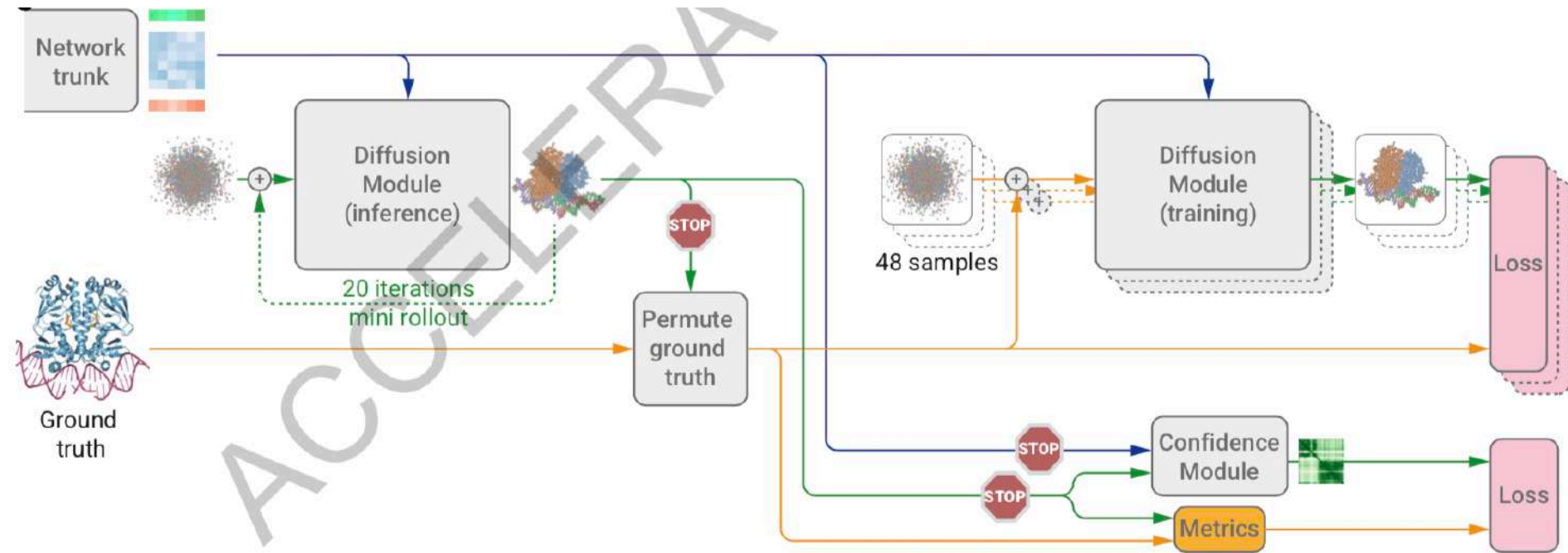
## Algorithm 18 Sample Diffusion

---

```
def SampleDiffusion( {f*}, {siinputs}, {sitrunk}, {zijtrunk}, Noise Schedule [c0, c1, ..., cT],  
                    γ0 = 0.8, γmin = 1.0, noise scale λ = 1.003, step scale η = 1.5 ):  
1:  $\vec{x}_l \sim c_0 \cdot \mathcal{N}(\vec{0}, \mathbf{I}_3) \implies$  start from pure noise  $\vec{x}_l \in \mathbb{R}^3$   
2: for all cτ ∈ [c1, ..., cT] do  $\implies$  go over noise scheduler  
3:   { $\vec{x}_l$ } ← CentreRandomAugmentation({ $\vec{x}_l$ })  $\implies$  augment the data (for batch size increase)  
4:   γ = γ0 if cτ > γmin else 0  
5:    $\hat{t} = c_{\tau-1}(\gamma + 1)$  }  $\implies$  determine the amount of noise  
6:    $\vec{\xi}_l = \lambda \sqrt{\hat{t}^2 - c_{\tau-1}^2} \cdot \mathcal{N}(\vec{0}, \mathbf{I}_3) \implies$  compute the noise vector  $\vec{\xi}_l \in \mathbb{R}^3$   
7:    $\vec{x}_l^{\text{noisy}} = \vec{x}_l + \vec{\xi}_l \implies$  corrupt the coordinates with the noise vector  
8:   { $\vec{x}_l^{\text{denoised}}$ } = DiffusionModule({ $\vec{x}_l^{\text{noisy}}$ },  $\hat{t}$ , {f*}, {siinputs}, {sitrunk}, {zijtrunk})  $\implies$  denoise conditionally on the sequence  
9:    $\vec{\delta}_l = (\vec{x}_l - \vec{x}_l^{\text{denoised}}) / \hat{t} \implies$  compute the denoising gradient  
10:  dt = cτ -  $\hat{t}$   
11:   $\vec{x}_l \leftarrow \vec{x}_l^{\text{noisy}} + \eta \cdot dt \cdot \vec{\delta}_l \implies$  update the coordinates from their noisy version in the denoising gradient direction  
12: end for  
13: return { $\vec{x}_l$ }
```

---

# AlphaFold 3 training and inference





# Taking another perspective

- Do we really need Langevin dynamics?
- What if we do not reach pure noise?
- What if we want to sample from another distribution?

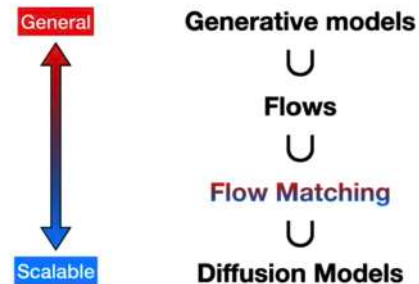
# Matching flows instead of scores

## Flow Matching

Simplifying and Generalizing Diffusion Models

Yaron Lipman

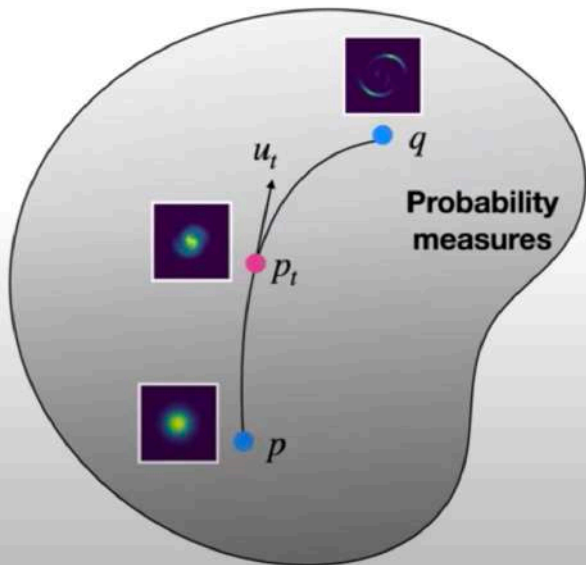
<https://www.youtube.com/watch?v=5ZSwYogAxYg>



# Matching flows instead of scores

Continuity Equation

$$\dot{p}_t = -\operatorname{div}(p_t u_t)$$



$p_t$  lies on a time-evolving probability path from

- $p$ : prior probability density (e.g. noise), to
- $q$ : target probability density (the true density of the data)

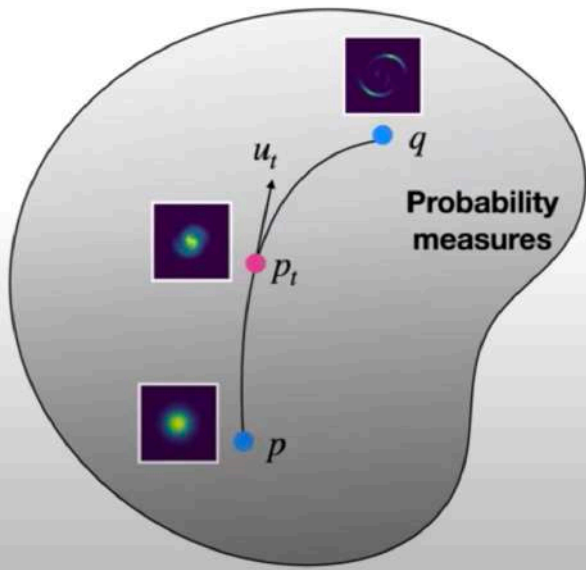
$u_t$  is the time-dependent vector field  $[0,1] \times \mathbb{R}^d \rightarrow \mathbb{R}^d$  defining the dynamics along the path, according to the Ordinary Differential Equation:  $\dot{x} = u_t(x_t)$

The **continuity equation** states that the change of probability density  $\dot{p}_t$  (e.g. seen as a fluid density) at any time point  $t$  is exactly counter-balanced by the change of flux volume in all directions  $\operatorname{div}(p_t u_t)$ .

# Compare velocities

Continuity Equation

$$\dot{p}_t = -\text{div}(p_t u_t)$$



$$L_{\text{FM}}(q||p_1) = \min \mathbb{E}_{t, p_t(x)} \|v_t(x) - u_t(x)\|^2$$

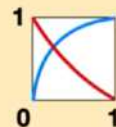
The FM loss (w. marginal vector field) is equivalent to the **conditional** FM loss where we break down to one training sample at a time.

$$L_{\text{CFM}}(q||p_1) = \min \mathbb{E}_{t, q(x_1), p_t(x|x_1)} \|v_t(x) - u_t(x|x_1)\|^2$$

$$= \min \mathbb{E}_{t, q(x_1), p(x_0)} \|v_t(x_t) - \dot{x}_t\|^2$$

$$x_t \sim p_t(x|x_1)$$

$$\dot{x}_t = u_t(x_t|x_1)$$



$$x_t = \Psi_t(x_0|x_1) = \sigma_t x_0 + \alpha_t x_1$$

$$x_0 \sim p \Rightarrow x_t \propto p\left(\frac{x - \alpha_t x_1}{\sigma_t}\right)$$

# A less restrictive framework

---

**Algorithm 1:** Flow Matching training.

---

**Input** : dataset  $q$ , noise  $p$

Initialize  $v^\theta$

**while** not converged **do**

$t \sim \mathcal{U}([0, 1])$  ▷ sample time  
 $x_1 \sim q(x_1)$  ▷ sample data  
 $x_0 \sim p(x_0)$  ▷ sample noise  
 $x_t = \Psi_t(x_0|x_1)$  ▷ conditional flow  
Gradient step with  $\nabla_\theta \|v_t^\theta(x_t) - \dot{x}_t\|^2$

**Output:**  $v^\theta$

---

$p_t(x_t|x_1)$  general  
 $p(x_0)$  is general

---

**Algorithm 2:** Diffusion training.

---

**Input** : dataset  $q$ , noise  $p$

Initialize  $s^\theta$

**while** not converged **do**

$t \sim \mathcal{U}([0, 1])$  ▷ sample time  
 $x_1 \sim q(x_1)$  ▷ sample data  
 $x_t = p_t(x_t|x_1)$  ▷ sample conditional prob  
Gradient step with  
 $\nabla_\theta \|s_t^\theta(x_t) - \nabla_{x_t} \log p_t(x_t|x_1)\|^2$

**Output:**  $v^\theta$

---

$p_t(x_t|x_1)$  closed-form from of SDE  $dx_t = f_t dt + g_t dw$

- Variance Exploding:**  $p_t(x|x_1) = \mathcal{N}(x|x_1, \sigma_{1-t}^2 I)$
- Variance Preserving:**  $p_t(x|x_1) = \mathcal{N}(x|\alpha_{1-t}x_1, (1 - \alpha_{1-t}^2)I)$   
 $\alpha_t = e^{-\frac{1}{2}T(t)}$

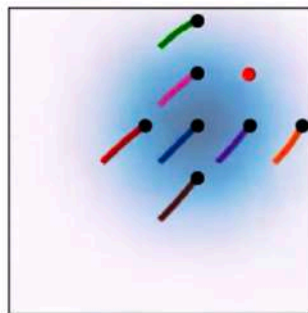
$p(x_0)$  is Gaussian

$p_0(\cdot|x_1) \approx p$

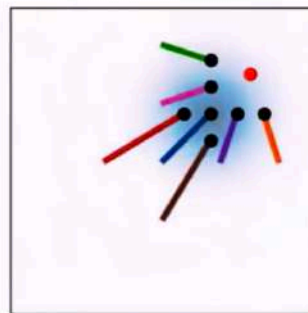
# A less restrictive framework

- Example of non-diffusion choice

$$p_t(x | x_1), u_t(x | x_1) \Leftrightarrow \Psi_t(x_0 | x_1)$$
$$\mathcal{N}(tx_1, (1-t)^2 I), \frac{x_1 - x}{1-t} \quad \Psi_t(x_0 | x_1) = (1-t)x_0 + tx_1$$



Diffusion (VP)



Cond-OT

# Take-home messages

- Diffusion models are generative models drawing inspiration from **Langevin** dynamics.
- They allow for **transforming noise** into any distribution we want to sample from. They are very flexible, allowing for conditioning on various properties.
- The learnable part of the process serves for **denoising** (estimating noise is equivalent to estimating the score) and for **conditioning** (with labels).
- Convergence to pure noise is critical, as well as defining an adequate **noise scheduler**.
- Some of the restrictions of diffusion models are alleviated by **flow-matching** models.

# Resources

- <https://iclr-blogposts.github.io/2024/blog/diffusion-theory-from-scratch/>
- <https://deepsense.ai/diffusion-models-in-practice-part-1-the-tools-of-the-trade/>
- <https://lilianweng.github.io/posts/2021-07-11-diffusion-models/>
- <https://yang-song.net/blog/2021/score/>
- <https://www.youtube.com/watch?v=a4Yfz2FxxiY&themeRefresh=1>
- <https://github.com/diff-usion/Awesome-Diffusion-Models?tab=readme-ov-file>
- <https://dauparas.github.io/post/af2/>
- Thanks to Julien Nguyen Van, Roman Klypa and Sergei Grudinin for discussions on the topic and sharing of resources.



# What now?

- Check out the Practical! [Made by Julien Nguyen Van]

