

# Neural Stochastic Flows: Solver-Free Modelling and Inference for SDE Solutions

Naoki Kiyohara<sup>1,2</sup>, Edward Johns<sup>1</sup>, Yingzhen Li<sup>1</sup>

<sup>1</sup>Imperial College London   <sup>2</sup>Canon Inc.

Project page



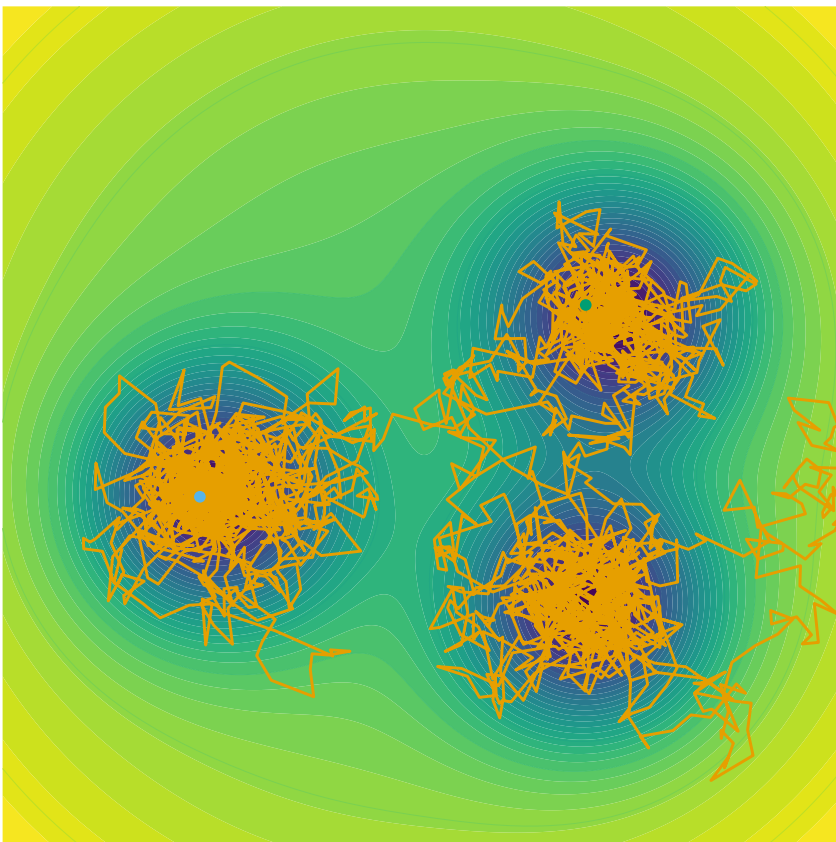
# Stochastic Differential Equations (SDEs) are Everywhere

Itô SDE:  $dx_t = \underbrace{\mu(x_t, t) dt}_{\text{Drift term (deterministic)}} + \underbrace{\sigma(x_t, t) dW_t}_{\text{Diffusion term (stochastic)}}, \quad x_0 \sim p_0$

Finance



Physics



Generative Modelling

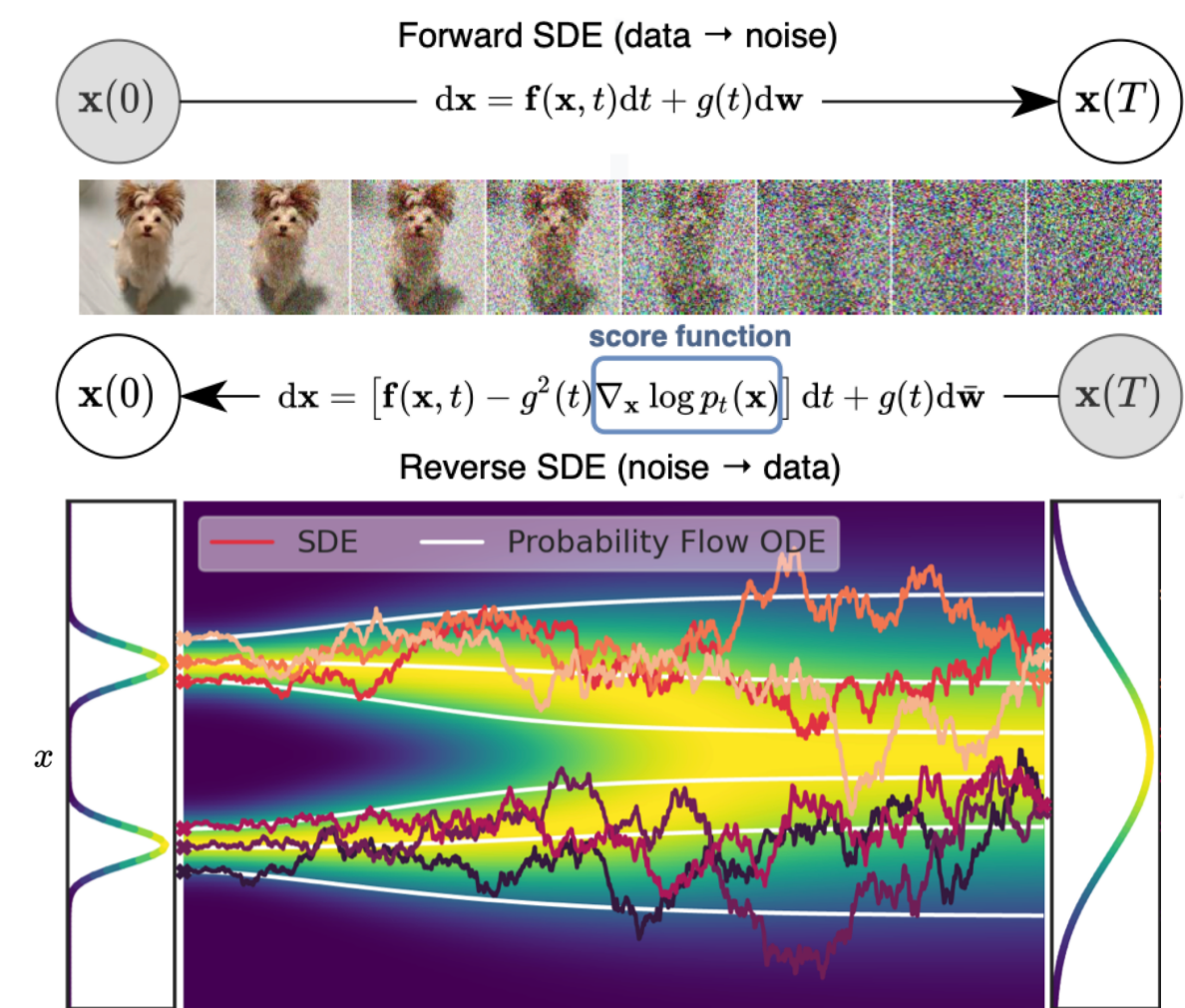


Figure adapted from [47]

# Neural SDEs

## Modelling SDE Terms using Neural Networks

Itô SDE:  $dx_t = \boxed{\mu(x_t, t)}dt + \boxed{\sigma(x_t, t)}dW_t, \quad x_0 \sim p_0$

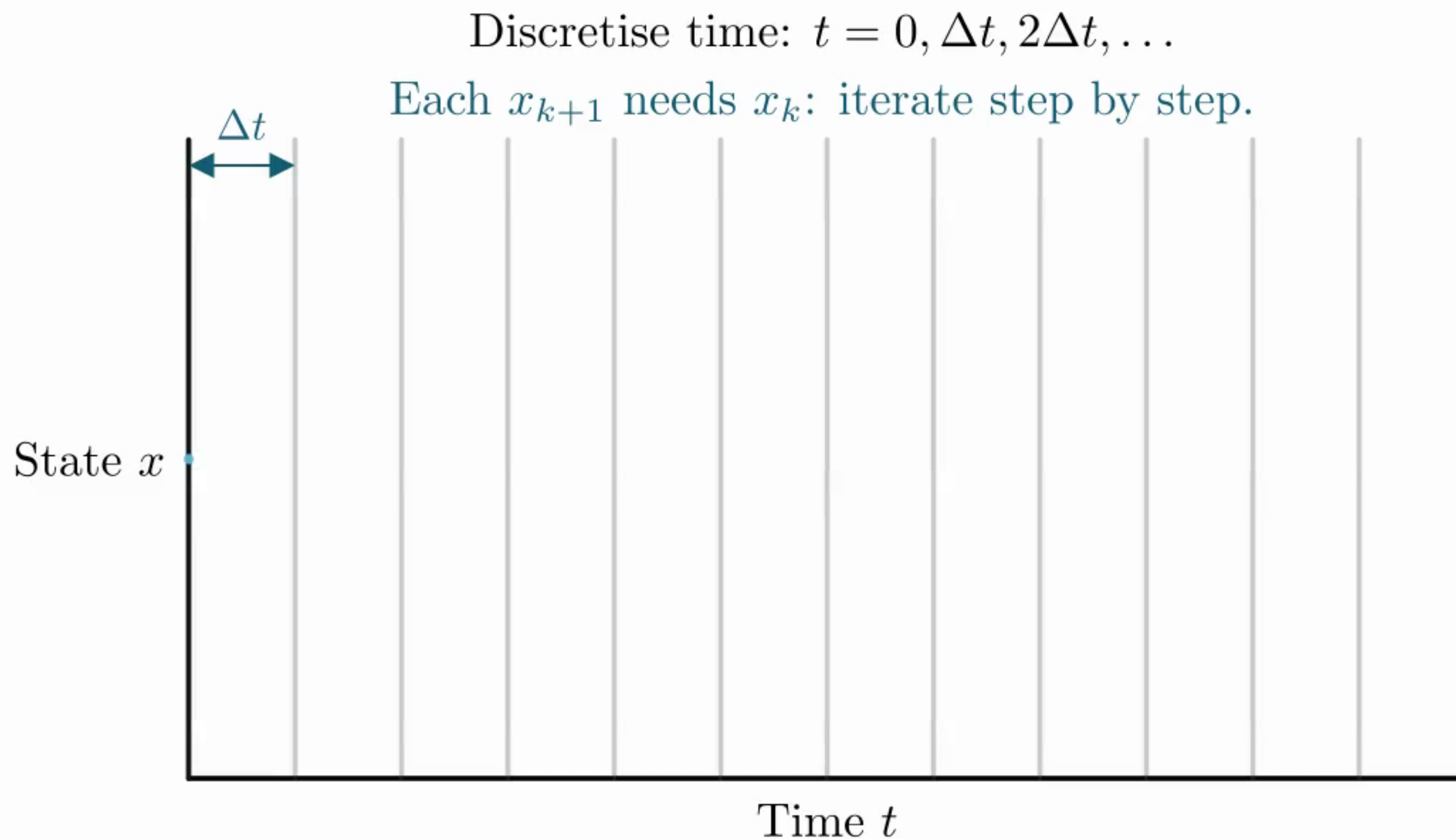
**Modelled by neural networks**



$$x_{t_{k+1}} = x_{t_k} + \mu(x_{t_k}, t_k) \Delta t + \sigma(x_{t_k}, t_k) \Delta W_k$$

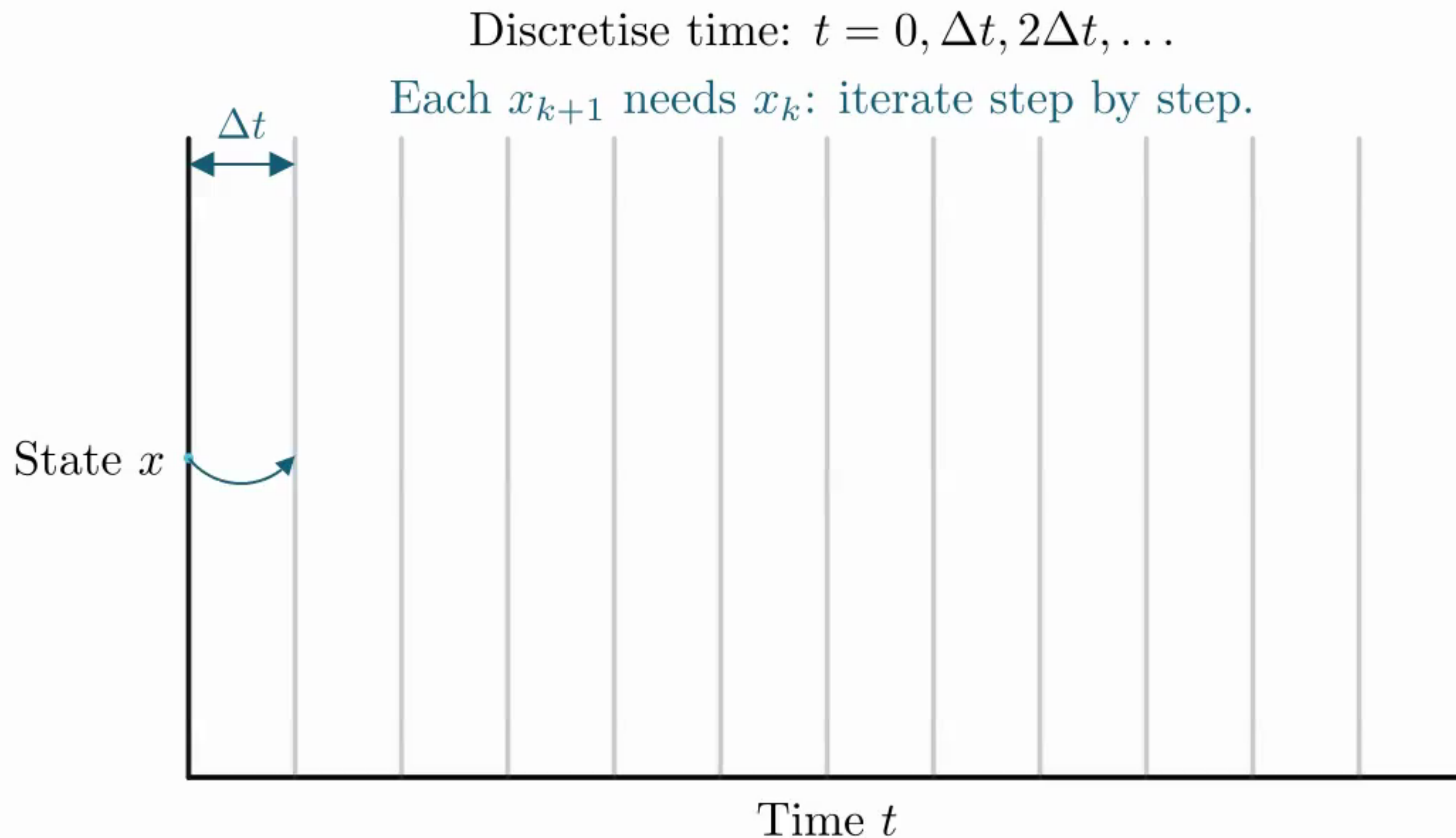
$$t_k = k\Delta t, \quad \Delta W_k = W_{t_{k+1}} - W_{t_k} \sim \mathcal{N}(0, \Delta t)$$

# Computational Challenges of Neural SDEs

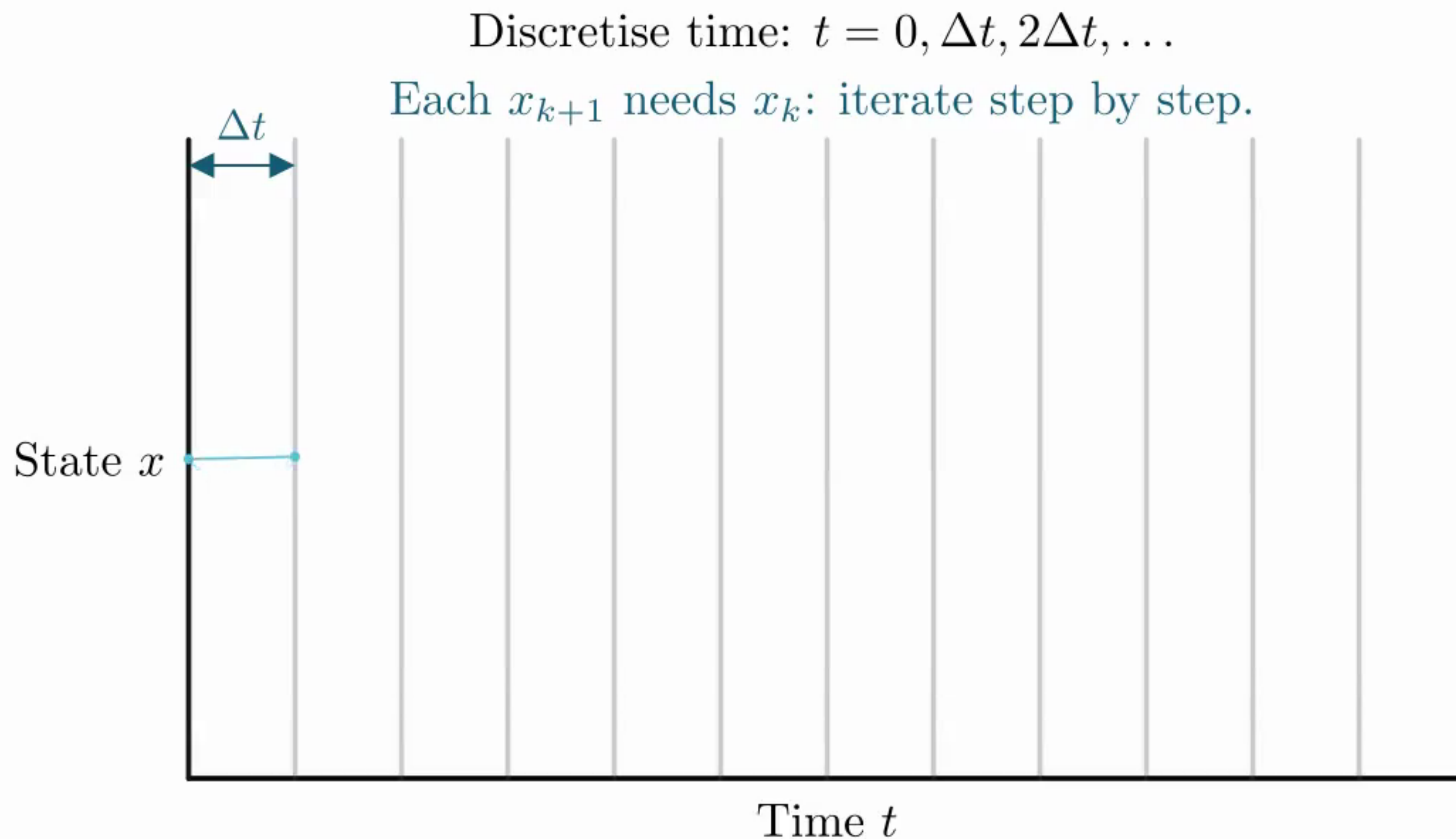




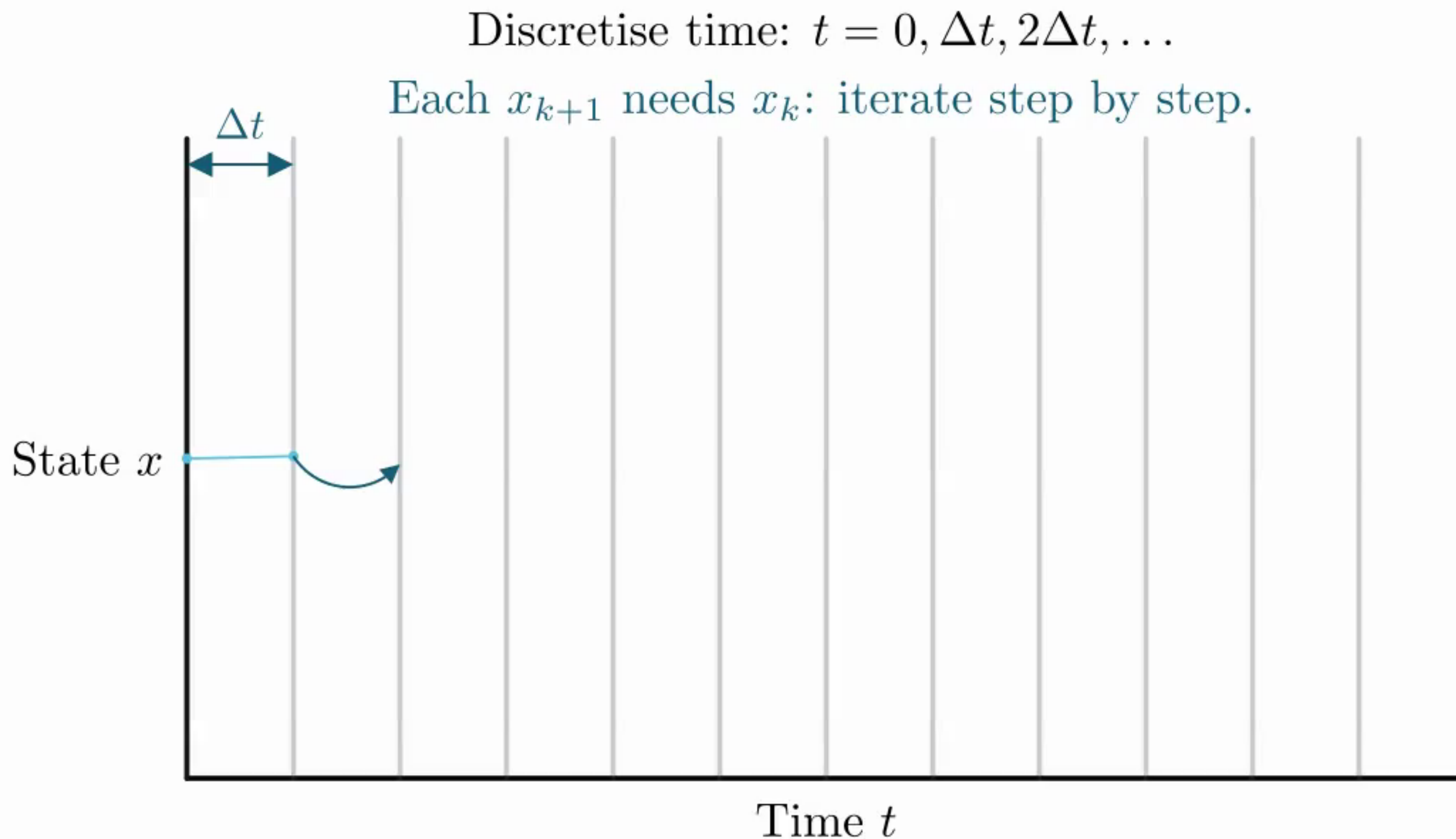
# Computational Challenges of Neural SDEs



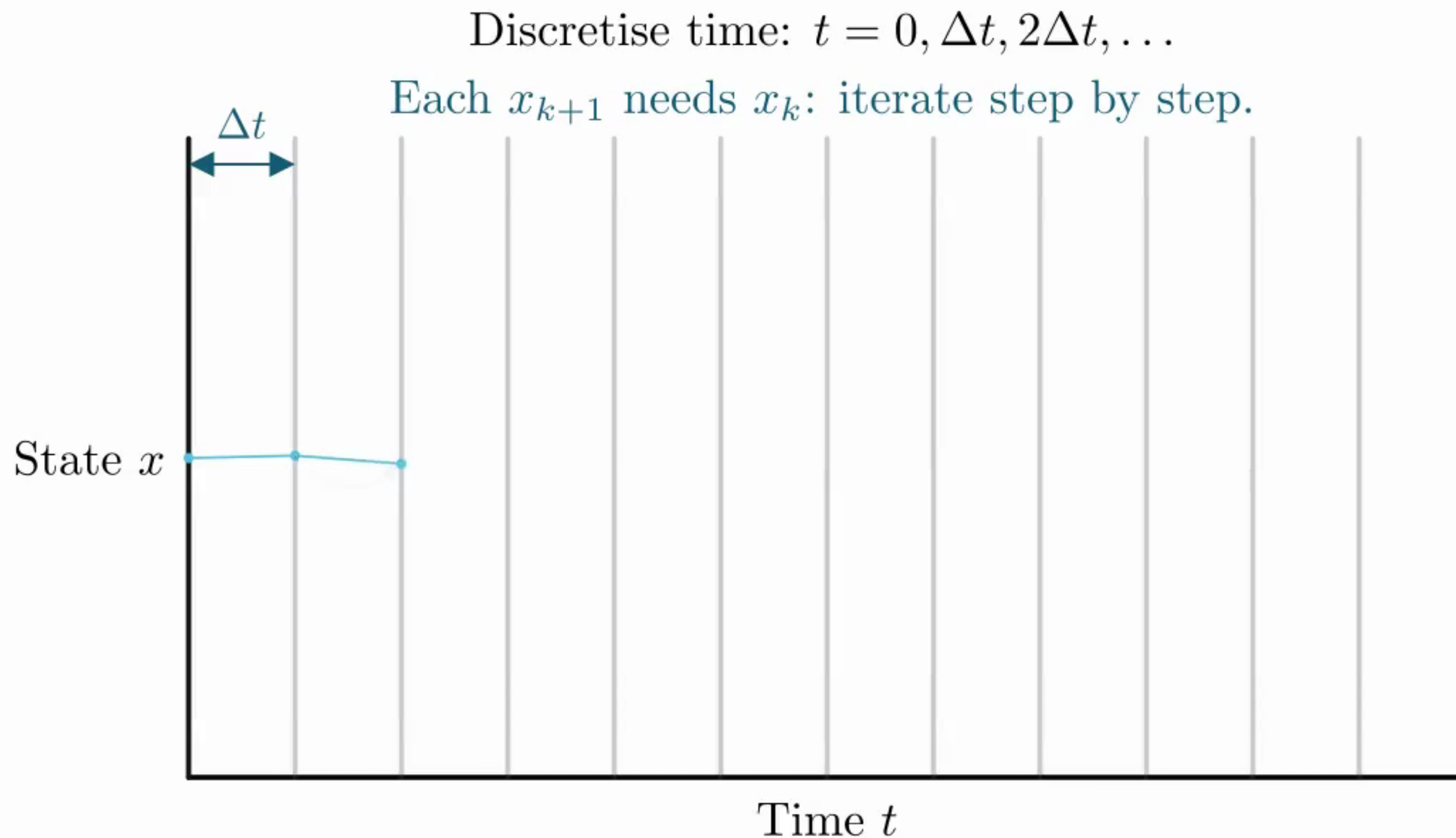
# Computational Challenges of Neural SDEs



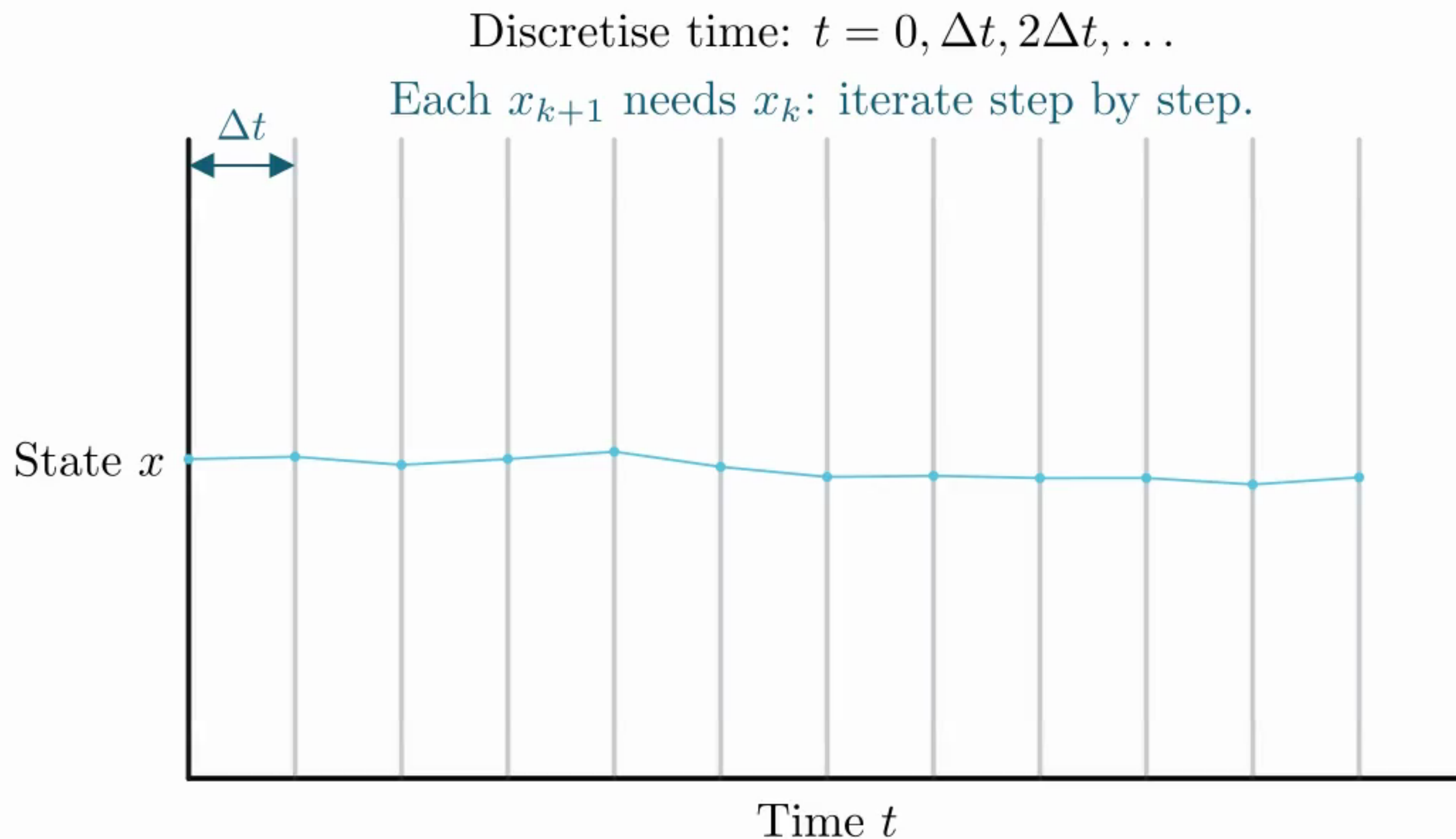
# Computational Challenges of Neural SDEs



# Computational Challenges of Neural SDEs

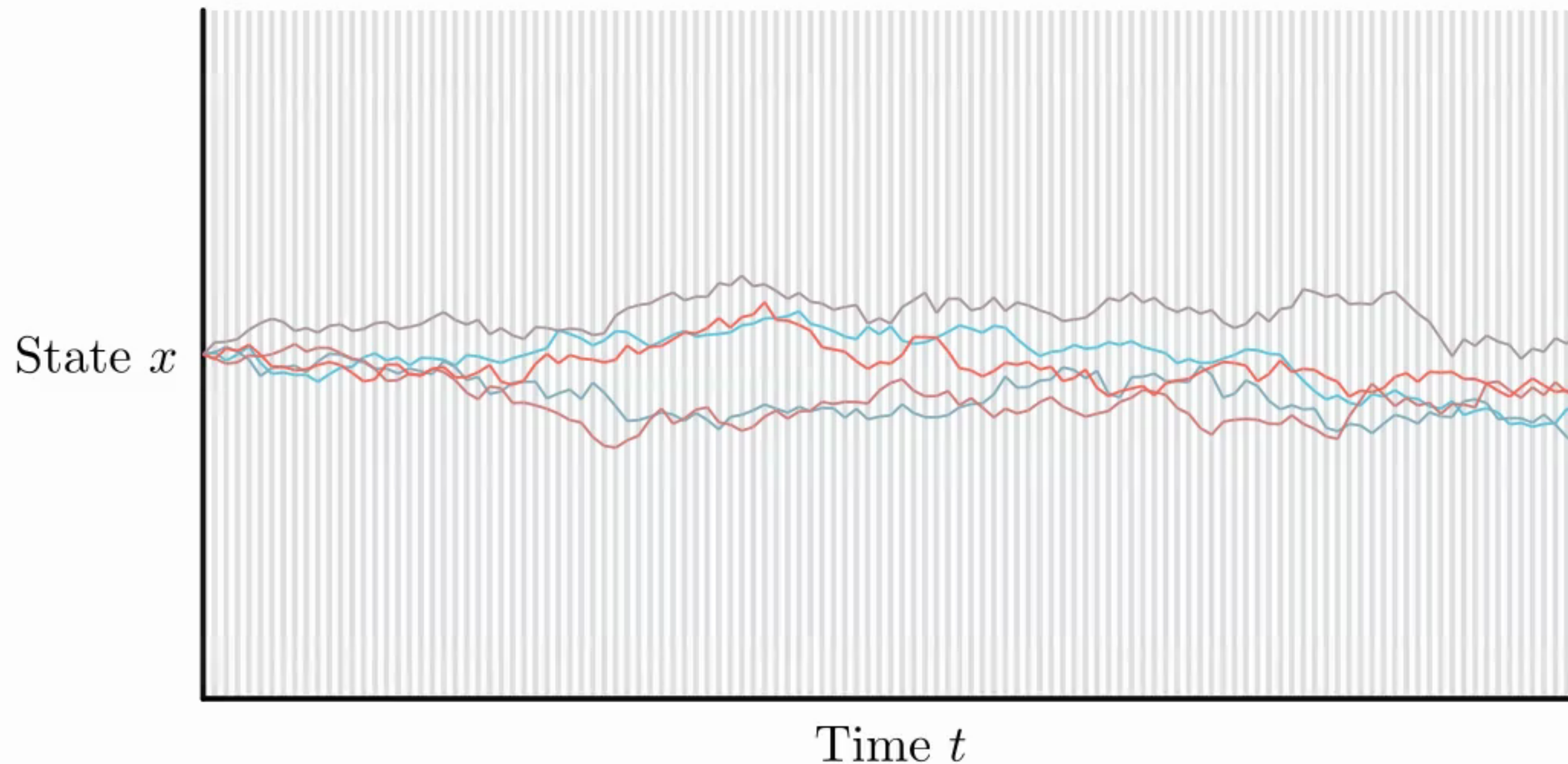


# Computational Challenges of Neural SDEs

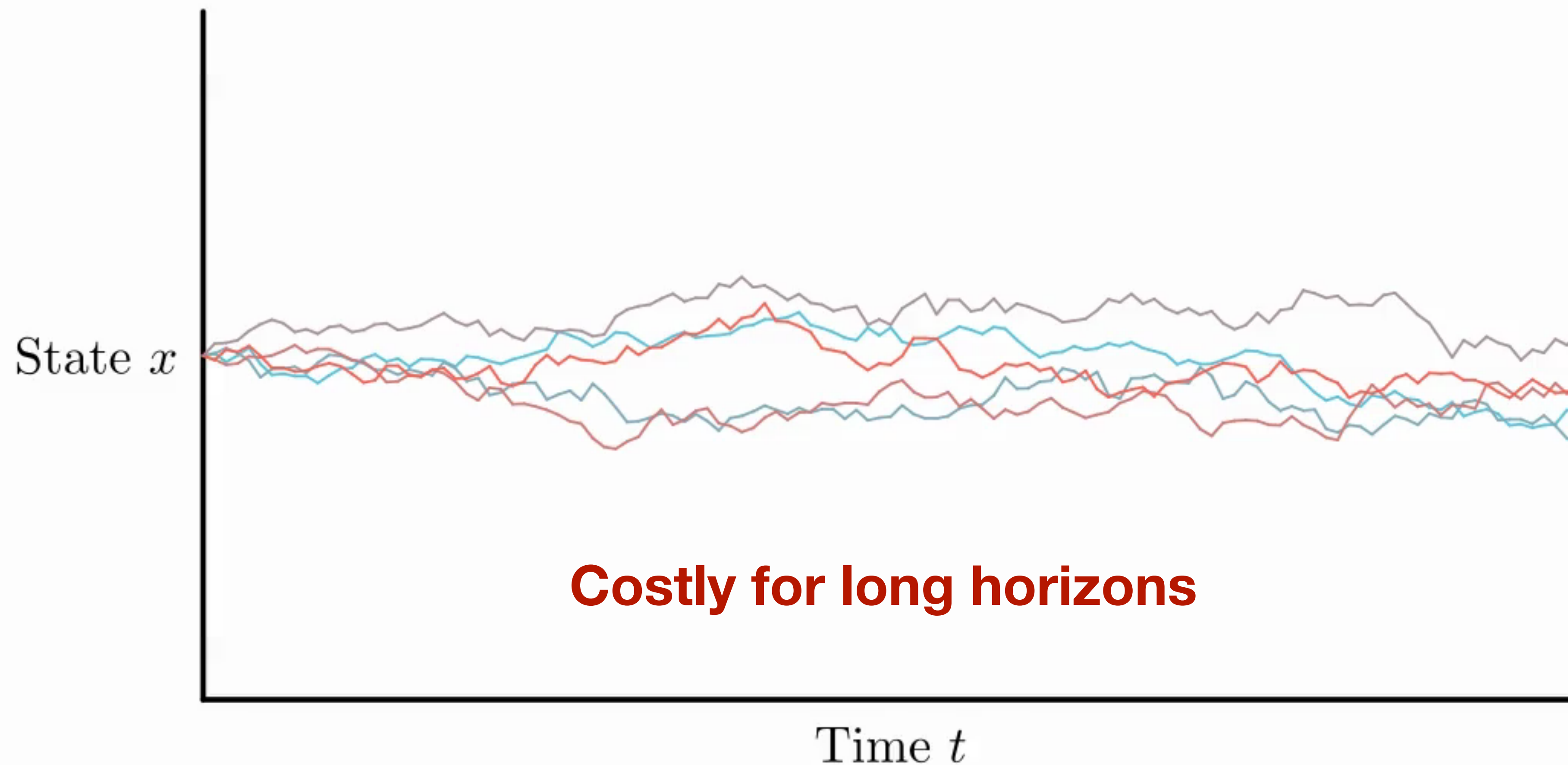




# Computational Challenges of Neural SDEs



# Computational Challenges of Neural SDEs



# Strong Solution and Weak Solution of the SDEs

$$\text{Itô SDE: } dx_t = \mu(x_t, t) dt + \sigma(x_t, t) dW_t, \quad x_0 \sim p_0$$

## Strong solution

A solution  $x_t$  defined on a given probability space with a given Brownian motion. You do not change the noise;  $x_t$  must adapt to the fixed  $W_t$ .

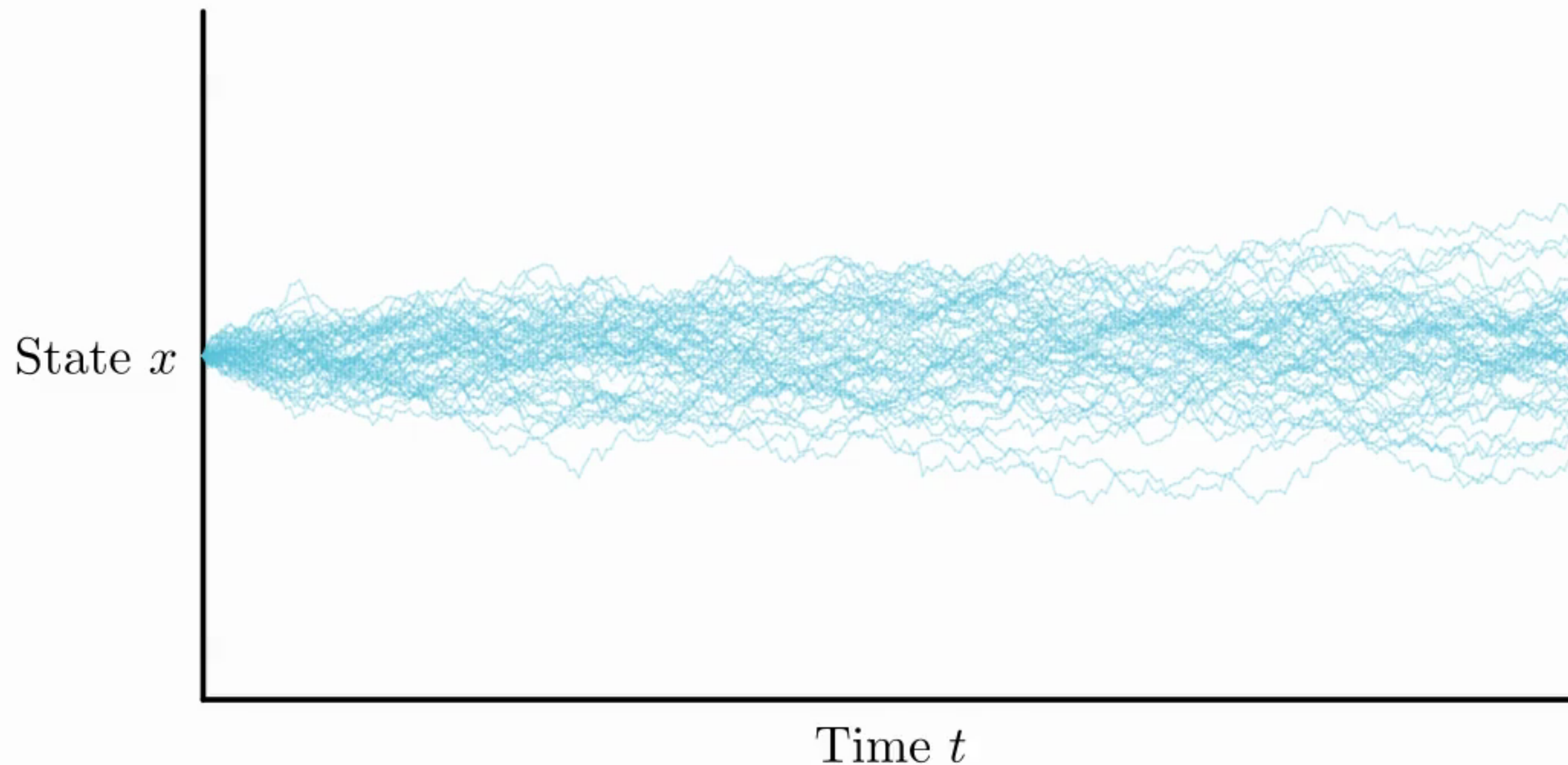
- Neural SDE: model sample paths w.r.t. a fixed Brownian motion.

## Weak solution

A solution  $x_t$  where the probability space and the Brownian motion are part of what you can choose. You only require that  $x_t$  and some Brownian motion satisfy the SDE and have the correct distribution.

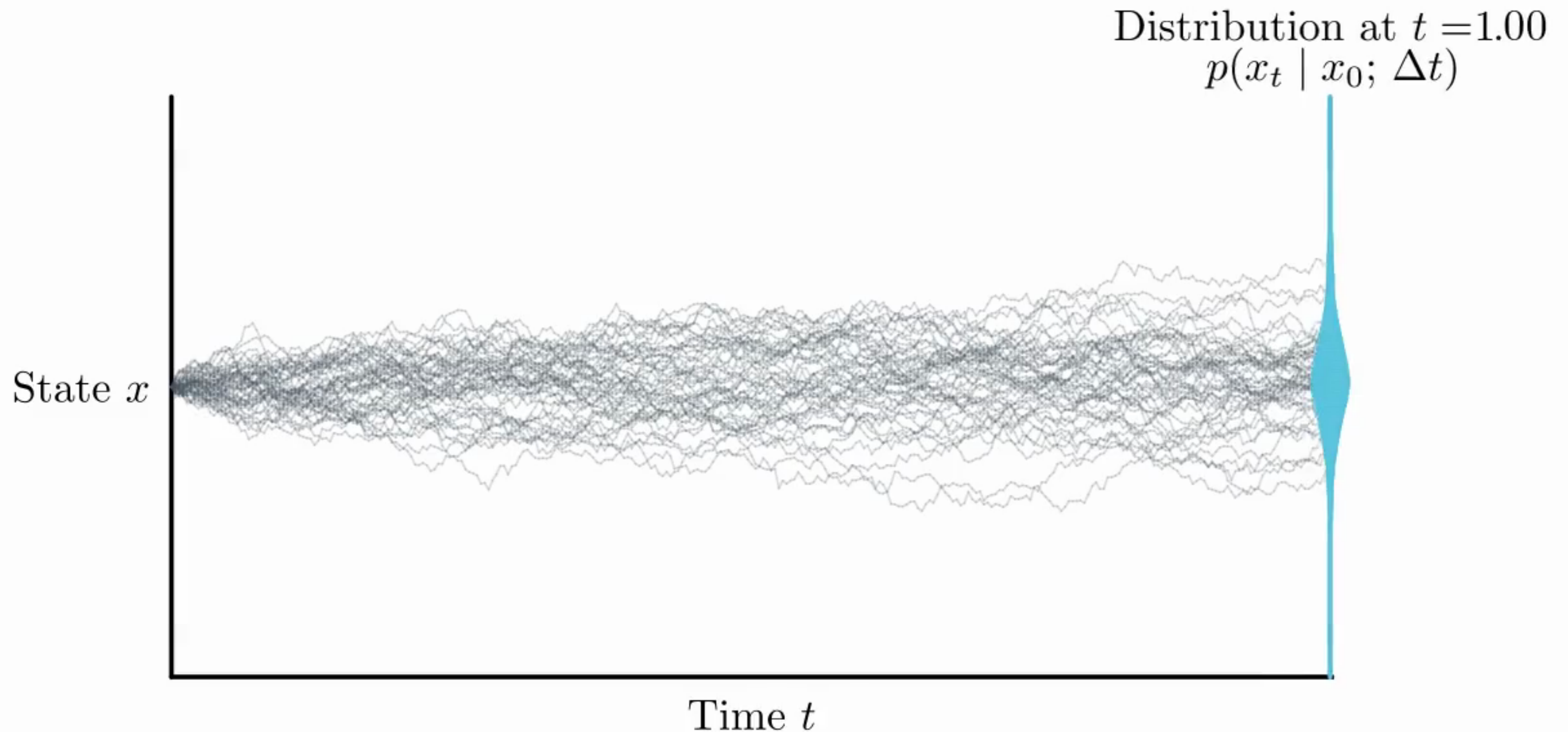
- **Our idea: SDE as a Markov model—characterised by the transition law  $p(x_t | x_s)$ .**

# Direct Modelling of Transition Distributions



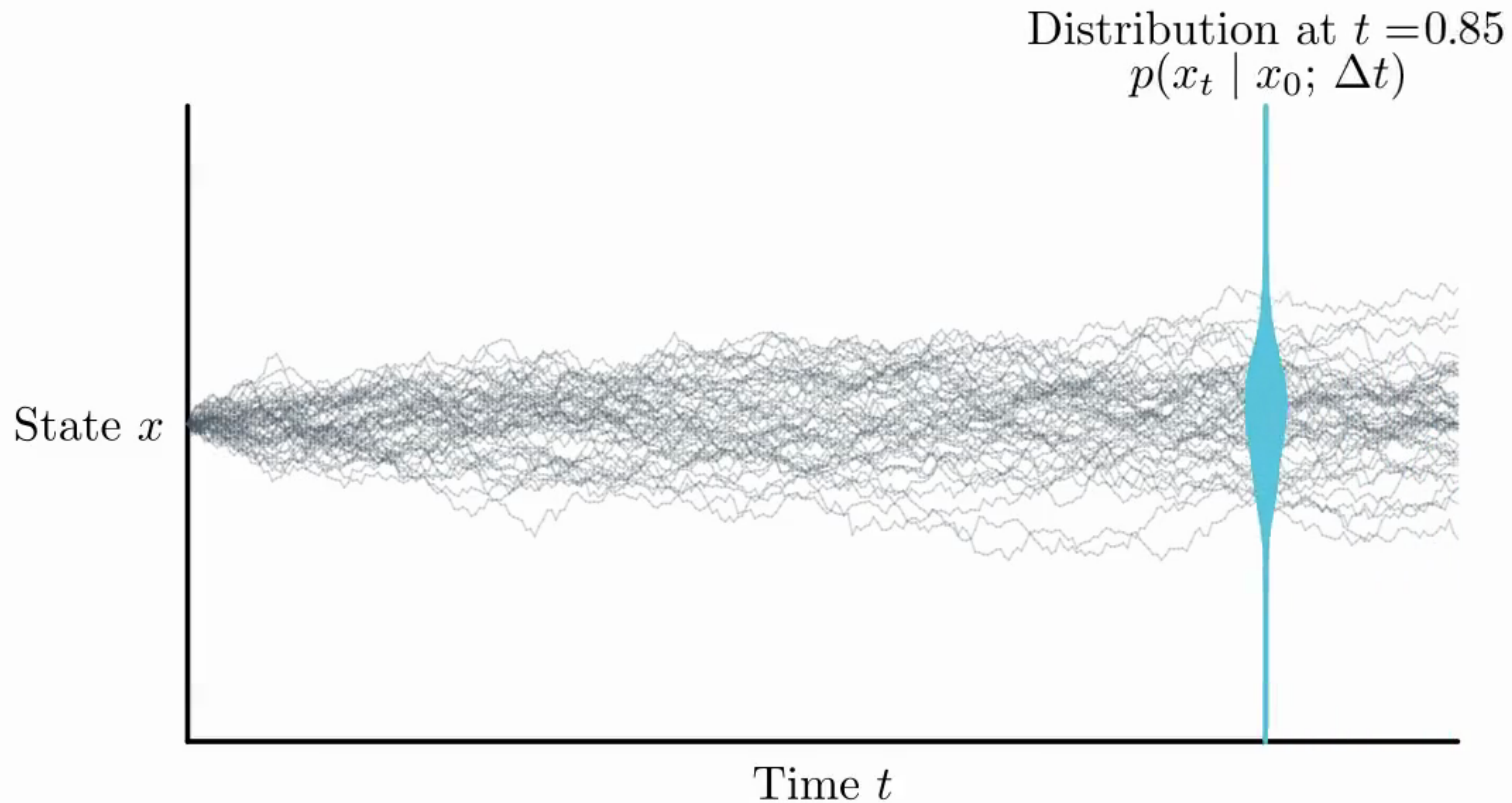


# Direct Modelling of Transition Distributions

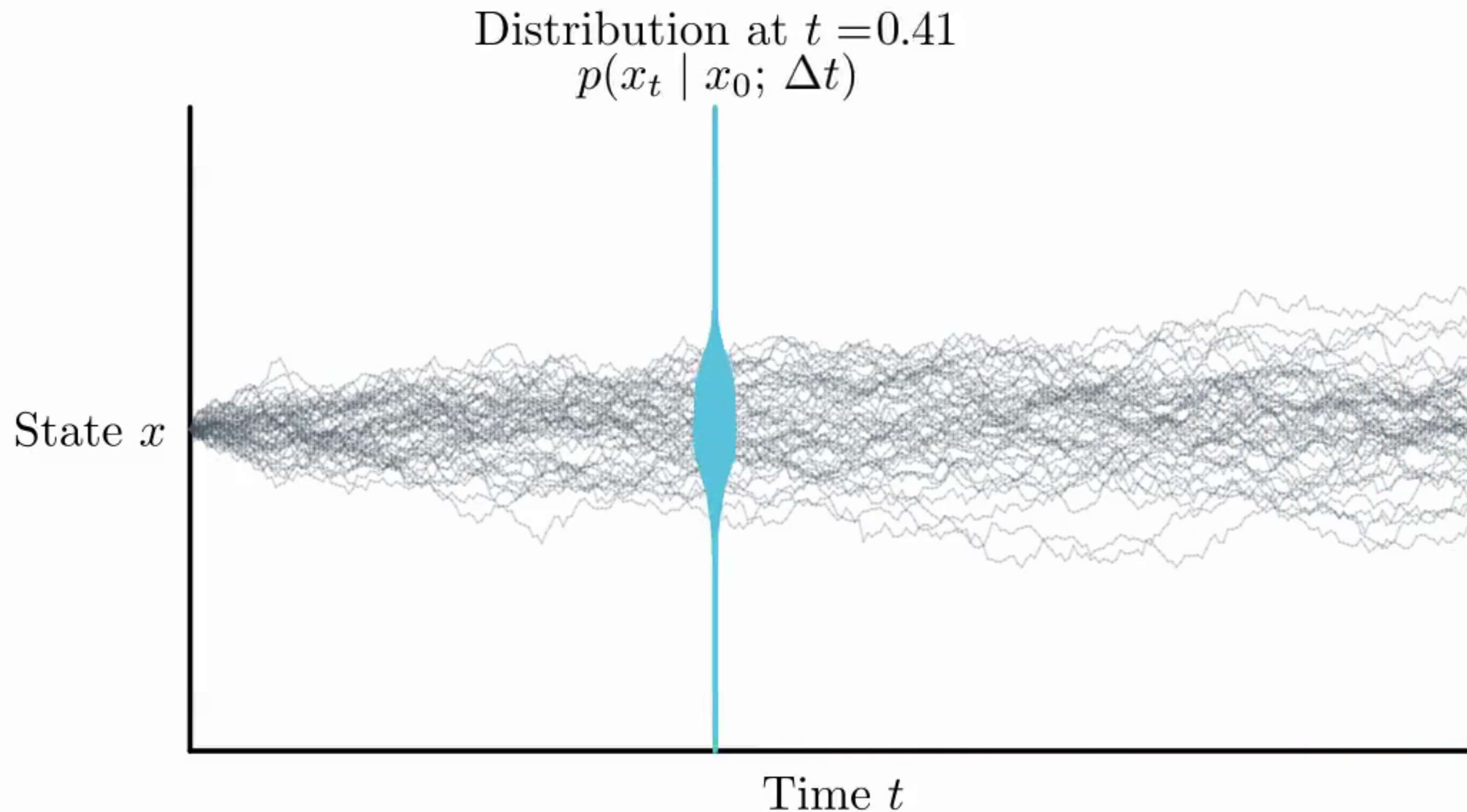




# Direct Modelling of Transition Distributions



# Direct Modelling of Transition Distributions



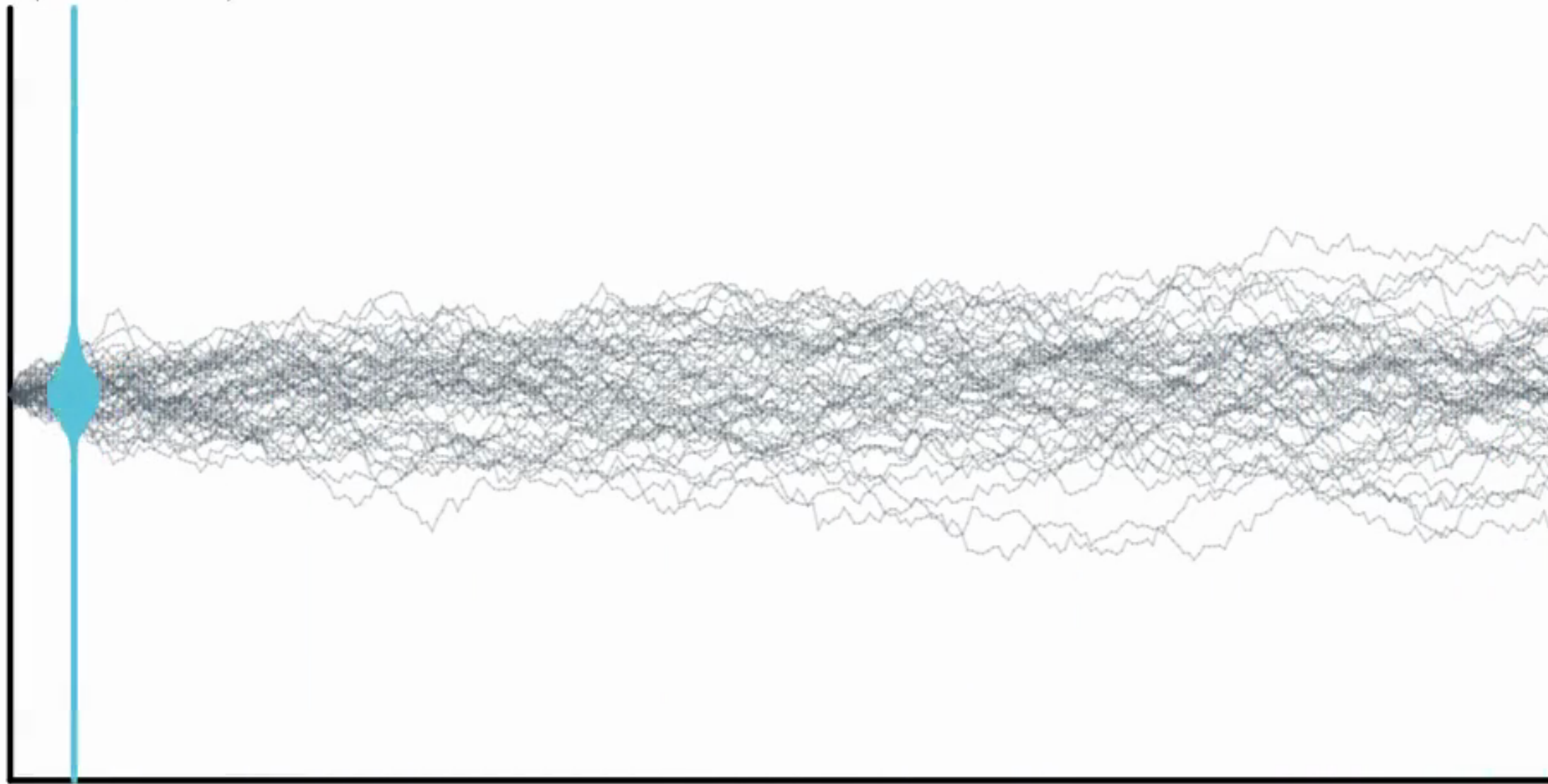


# Direct Modelling of Transition Distributions

Distribution at  $t = 0.04$   
 $p(x_t | x_0; \Delta t)$

State  $x$

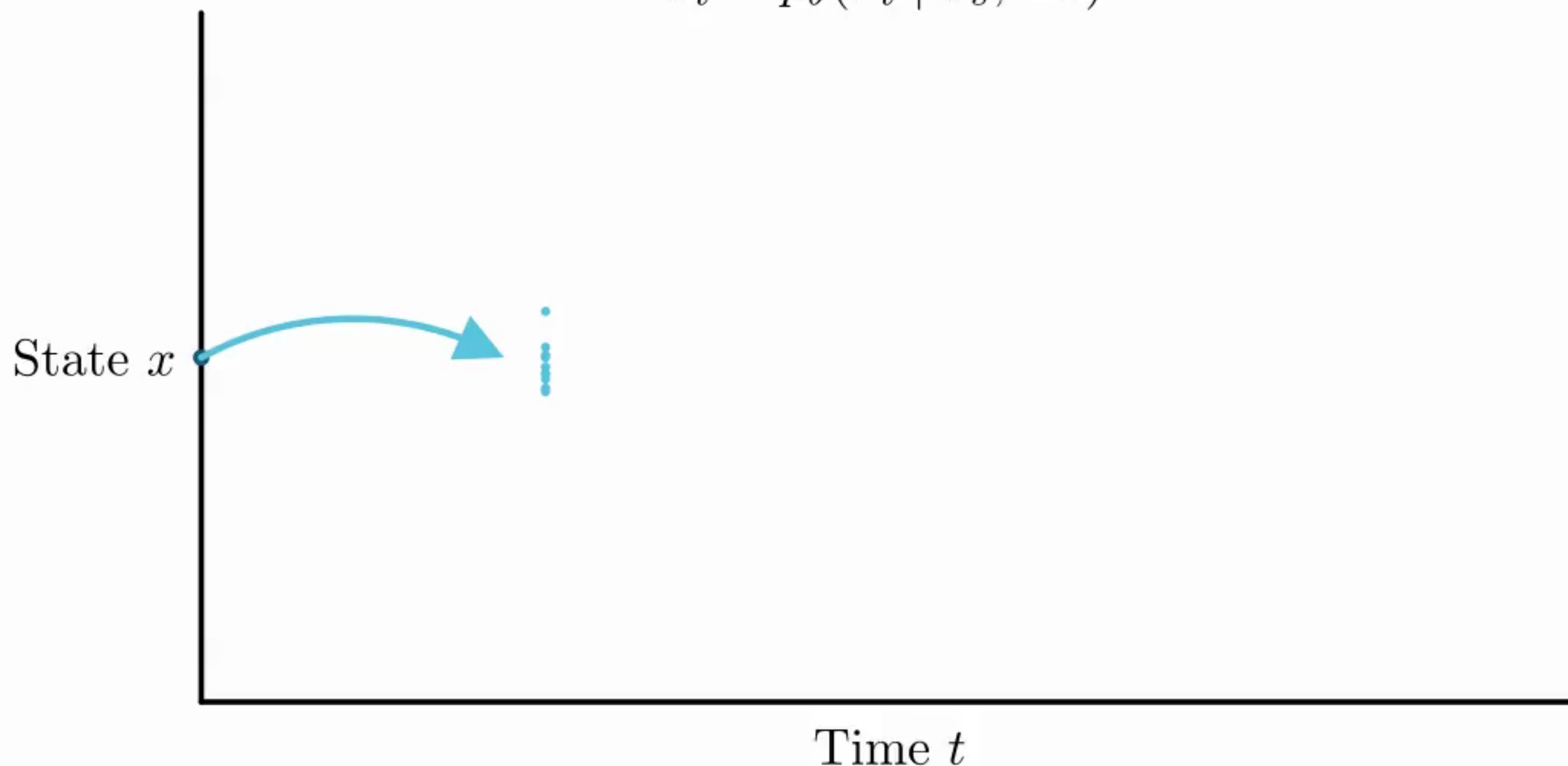
Time  $t$



# Direct Modelling of Transition Distributions

One-step sampling at arbitrary time

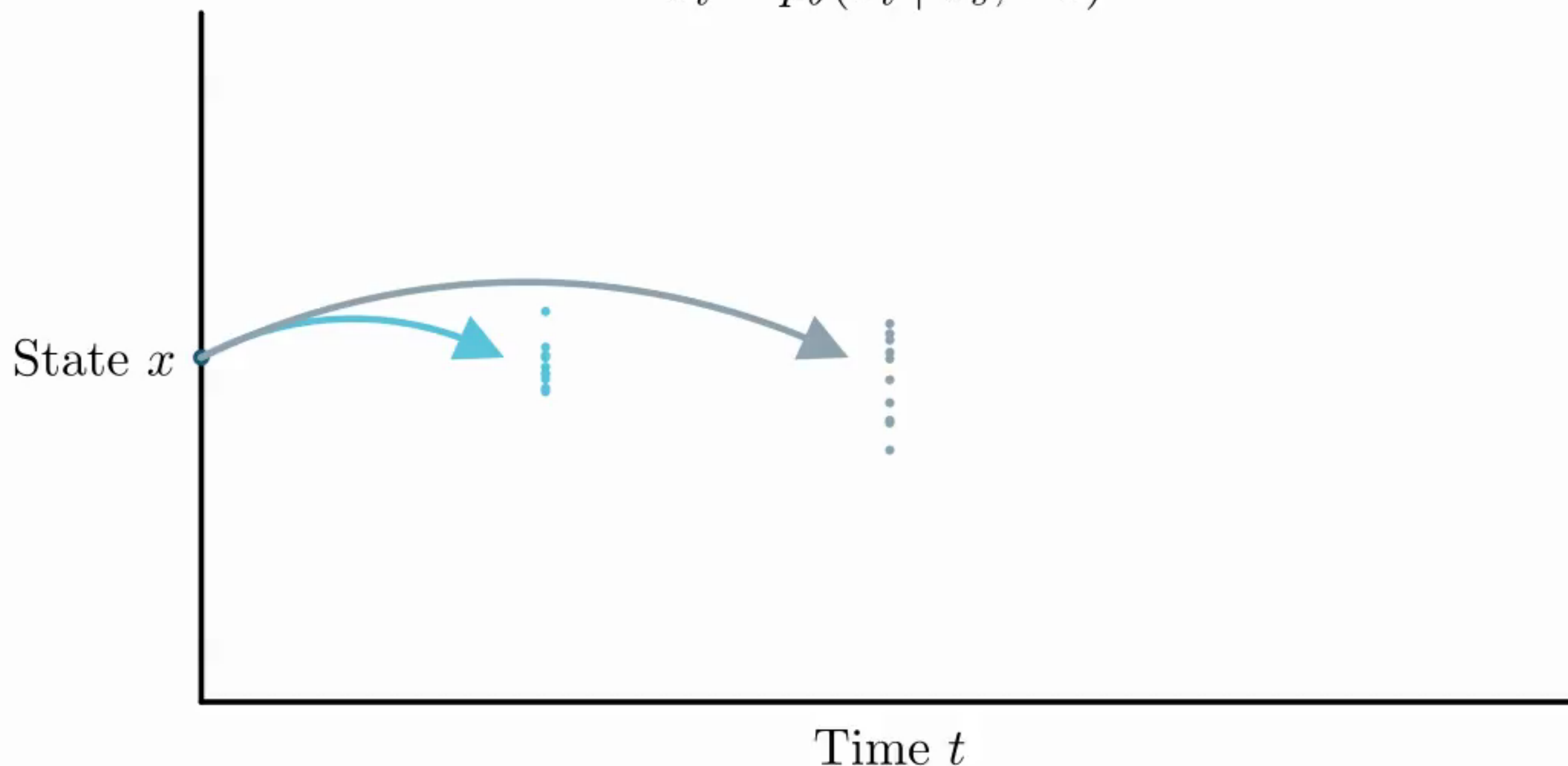
$$x_t \sim p_\theta(x_t \mid x_s; \Delta t)$$



# Direct Modelling of Transition Distributions

One-step sampling at arbitrary time

$$x_t \sim p_\theta(x_t \mid x_s; \Delta t)$$

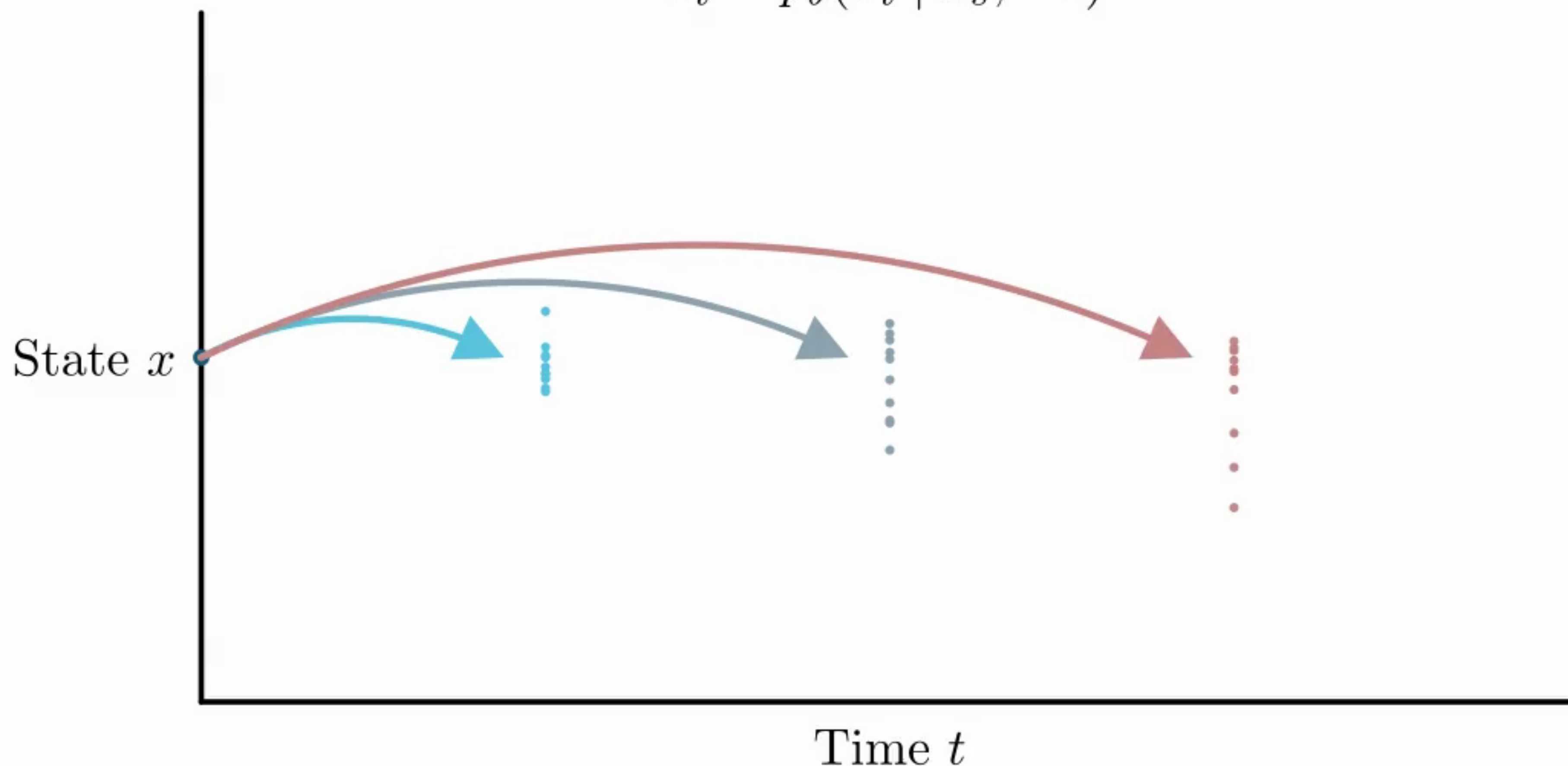




# Direct Modelling of Transition Distributions

One-step sampling at arbitrary time

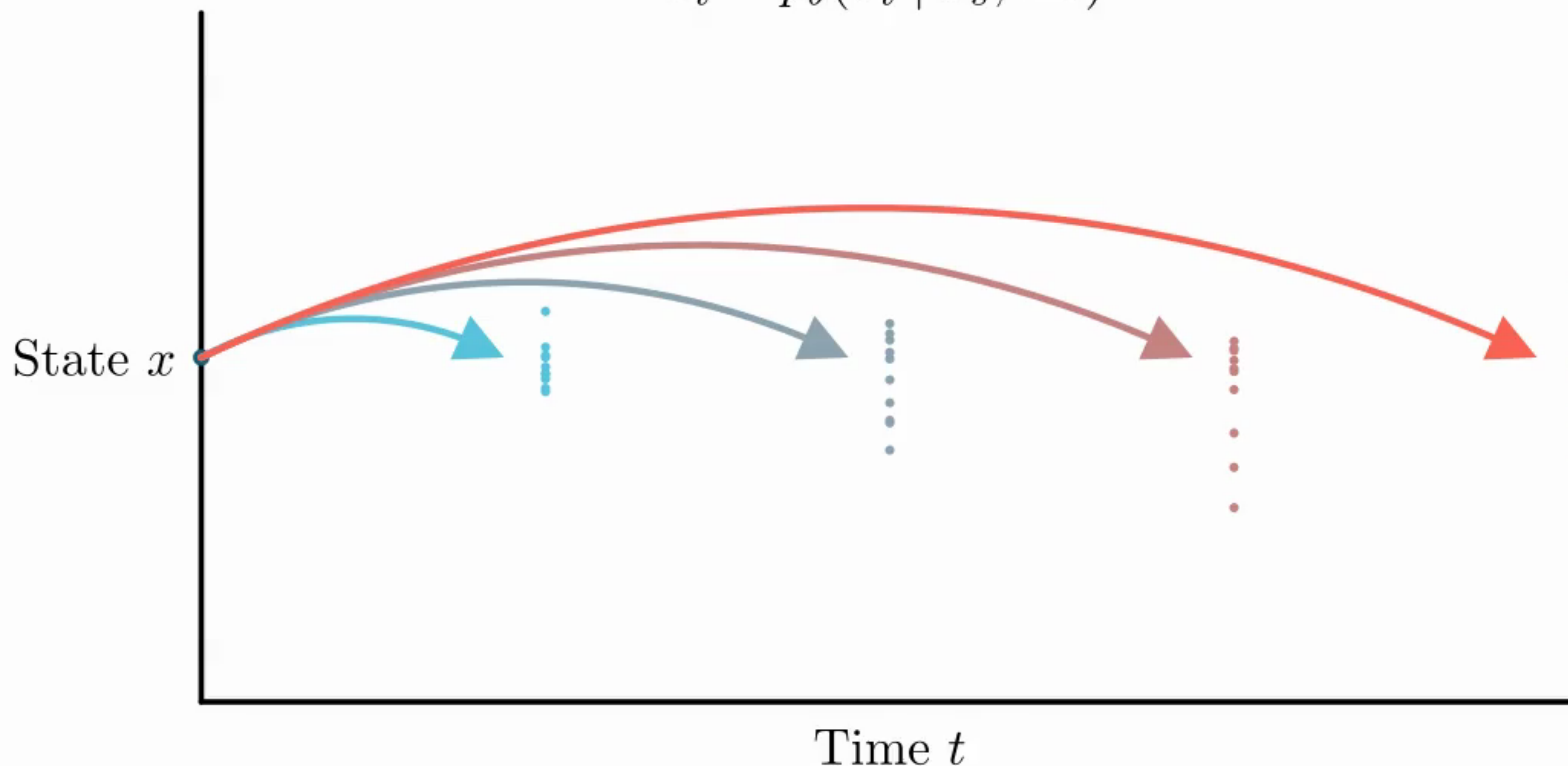
$$x_t \sim p_\theta(x_t \mid x_s; \Delta t)$$



# Direct Modelling of Transition Distributions

One-step sampling at arbitrary time

$$x_t \sim p_\theta(x_t \mid x_s; \Delta t)$$



# Mathematical Requirements for the Transition Kernel

## Properties of $p_{\theta}(\mathbf{x}_t \mid \mathbf{x}_s; s, \Delta t = t - s)$

### 1. Independence.

For any  $0 \leq t_1 \leq \dots \leq t_n$ ,  $p_{\theta}(\mathbf{x}_{t_{k+1}} \mid \mathbf{x}_{t_k}; t_k, t_{k+1} - t_k)$  are independent.

### 2. Identity.

When  $t = s$ ,  $p_{\theta}(\mathbf{x}_t \mid \mathbf{x}_s; s, 0) = \delta(\mathbf{x}_t - \mathbf{x}_s)$

### 3. Flow property.

For any  $0 \leq t_i \leq t_j \leq t_k$ ,  $p_{\theta}(\mathbf{x}_{t_k} \mid \mathbf{x}_{t_i}; t_i, t_k - t_i) = \int p_{\theta}(\mathbf{x}_{t_k} \mid \mathbf{x}_{t_j}; t_j, t_k - t_j) p_{\theta}(\mathbf{x}_{t_j} \mid \mathbf{x}_{t_i}; t_i, t_j - t_i) d\mathbf{x}_{t_j}$

### 4. Stationarity (for autonomous SDEs).

$$p_{\theta}(\mathbf{x}_{t_j} \mid \mathbf{x}_s; t_i, t_j - t_i) = p_{\theta}(\mathbf{x}_{t_j+r} \mid \mathbf{x}_s; t_i + r, t_j - t_i)$$

# Mathematical Requirements for the Transition Kernel

## Properties of $p_{\theta}(\mathbf{x}_t \mid \mathbf{x}_s; s, \Delta t = t - s)$

### 1. Independence. ✓ Realised by non-overlapping sampling.

For any  $0 \leq t_1 \leq \dots \leq t_n$ ,  $p_{\theta}(\mathbf{x}_{t_{k+1}} \mid \mathbf{x}_{t_k}; t_k, t_{k+1} - t_k)$  are independent.

### 2. Identity.

When  $t = s$ ,  $p_{\theta}(\mathbf{x}_t \mid \mathbf{x}_s; s, 0) = \delta(\mathbf{x}_t - \mathbf{x}_s)$

### 3. Flow property.

For any  $0 \leq t_i \leq t_j \leq t_k$ ,  $p_{\theta}(\mathbf{x}_{t_k} \mid \mathbf{x}_{t_i}; t_i, t_k - t_i) = \int p_{\theta}(\mathbf{x}_{t_k} \mid \mathbf{x}_{t_j}; t_j, t_k - t_j) p_{\theta}(\mathbf{x}_{t_j} \mid \mathbf{x}_{t_i}; t_i, t_j - t_i) d\mathbf{x}_{t_j}$

### 4. Stationarity (for autonomous SDEs).

$$p_{\theta}(\mathbf{x}_{t_j} \mid \mathbf{x}_s; t_i, t_j - t_i) = p_{\theta}(\mathbf{x}_{t_j+r} \mid \mathbf{x}_s; t_i + r, t_j - t_i)$$

# Architectural Design

**Conditional normalising flow for  $p_{\theta}(x_{t_j} \mid x_{t_i}; t_i, \Delta t)$**

## Base distribution

Using conditioning  $c := (x_{t_i}, \Delta t, t_i)$ ,

$$z = x_{t_i} + \Delta t \cdot \text{MLP}_{\mu}(c; \theta_{\mu}) + \sqrt{\Delta t} \cdot \text{MLP}_{\sigma}(c; \theta_{\sigma}) \odot \varepsilon, \quad \varepsilon \sim \mathcal{N}(\mathbf{0}, I)$$

## Conditional affine coupling layers [4, 13, 28]

$$x_{t_j} = f_{\theta}(z, c) = f_L(\cdot; c, \theta_L) \circ f_{L-1}(\cdot; c, \theta_{L-1}) \circ \cdots \circ f_1(z; c, \theta_1)$$

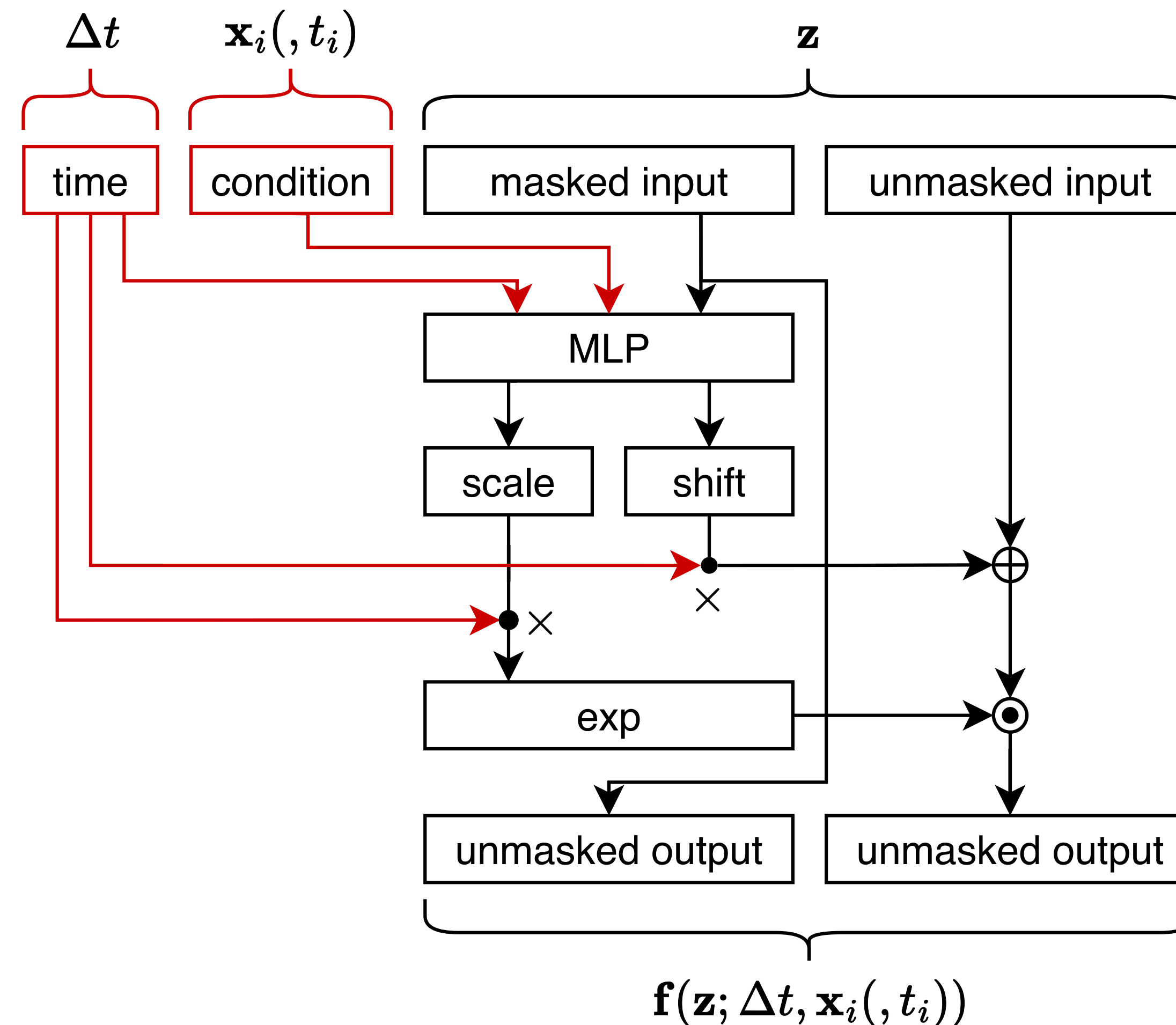
With  $z_A, z_B = \text{Split}(z)$ , each layer  $f$  is designed as:

$$\text{Concat} \left( z_A, z_B \odot \exp \left( \Delta t \text{MLP}_{\text{scale}}^{(i)}(z_A, c; \theta_{\text{scale}}^{(i)}) \right) + \Delta t \text{MLP}_{\text{shift}}^{(i)}(z_A, c; \theta_{\text{shift}}^{(i)}) \right)$$



# Architectural Design

Conditional affine coupling part for  $p_{\theta}(\mathbf{x}_{t_j} \mid \mathbf{x}_{t_i}; t_i, \Delta t)$



# Architectural Design

Conditional normalising flow for  $p_{\theta}(x_{t_j} \mid x_{t_i}; t_i, \Delta t)$

## Base distribution

Using conditioning  $c := (x_{t_i}, \Delta t, t_i)$ ,

$$z = x_{t_i} + \Delta t \cdot \text{MLP}_{\mu}(c; \theta_{\mu}) + \sqrt{\Delta t} \cdot \text{MLP}_{\sigma}(c; \theta_{\sigma}) \odot \varepsilon, \quad \varepsilon \sim \mathcal{N}(\mathbf{0}, I)$$

## Conditional affine coupling layers [4, 13, 28]

$$x_{t_j} = f_{\theta}(z, c) = f_L(\cdot; c, \theta_L) \circ f_{L-1}(\cdot; c, \theta_{L-1}) \circ \cdots \circ f_1(z; c, \theta_1)$$

With  $z_A, z_B = \text{Split}(z)$ , each layer  $f$  is designed as:

$$\text{Concat} \left( z_A, z_B \odot \exp \left( \Delta t \text{MLP}_{\text{scale}}^{(i)}(z_A, c; \theta_{\text{scale}}^{(i)}) \right) + \Delta t \text{MLP}_{\text{shift}}^{(i)}(z_A, c; \theta_{\text{shift}}^{(i)}) \right)$$

✓ Identity when  $\Delta t = 0$ .    ✓ Stationary when we drop  $t_i$  from  $c$ .

# Mathematical Requirements for the Transition Kernel

## Properties of $p_{\theta}(\mathbf{x}_t \mid \mathbf{x}_s; s, \Delta t = t - s)$

### 1. Independence. ✓ Realised by non-overlapping sampling.

For any  $0 \leq t_1 \leq \dots \leq t_n$ ,  $p_{\theta}(\mathbf{x}_{t_{k+1}} \mid \mathbf{x}_{t_k}; t_k, t_{k+1} - t_k)$  are independent.

### 2. Identity. ✓ Realised by architecture.

When  $t = s$ ,  $p_{\theta}(\mathbf{x}_t \mid \mathbf{x}_s; s, 0) = \delta(\mathbf{x}_t - \mathbf{x}_s)$

### 3. Flow property.

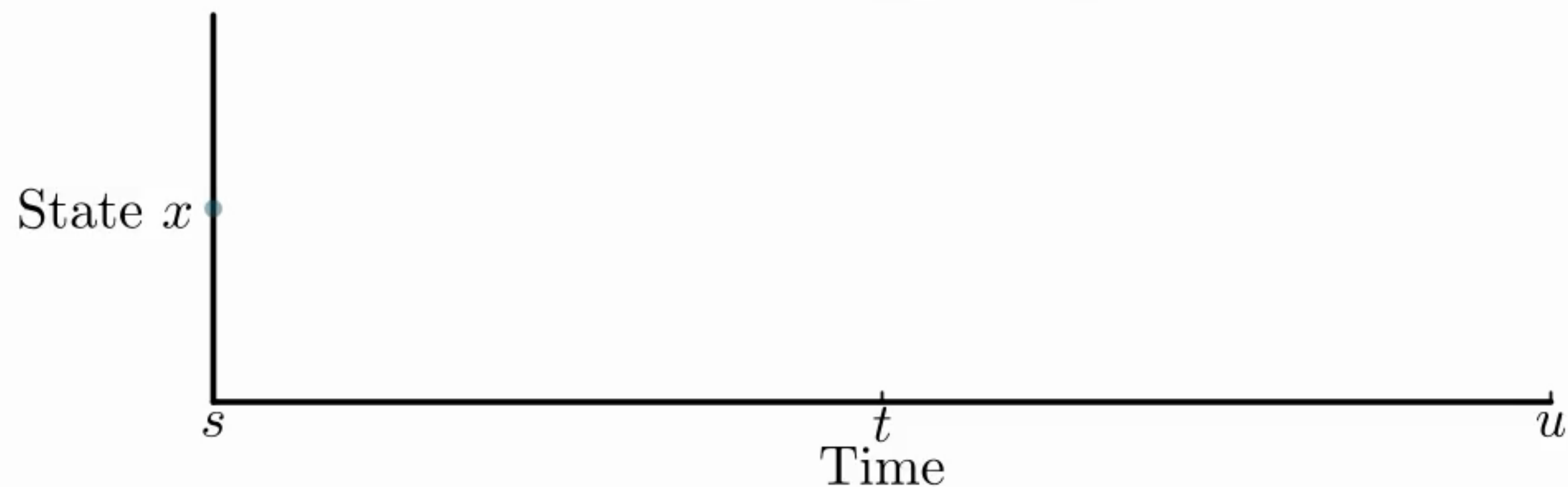
For any  $0 \leq t_i \leq t_j \leq t_k$ ,  $p_{\theta}(\mathbf{x}_{t_k} \mid \mathbf{x}_{t_i}; t_i, t_k - t_i) = \int p_{\theta}(\mathbf{x}_{t_k} \mid \mathbf{x}_{t_j}; t_j, t_k - t_j) p_{\theta}(\mathbf{x}_{t_j} \mid \mathbf{x}_{t_i}; t_i, t_j - t_i) d\mathbf{x}_{t_j}$

### 4. Stationarity (for autonomous SDEs). ✓ Realised by architecture optionally.

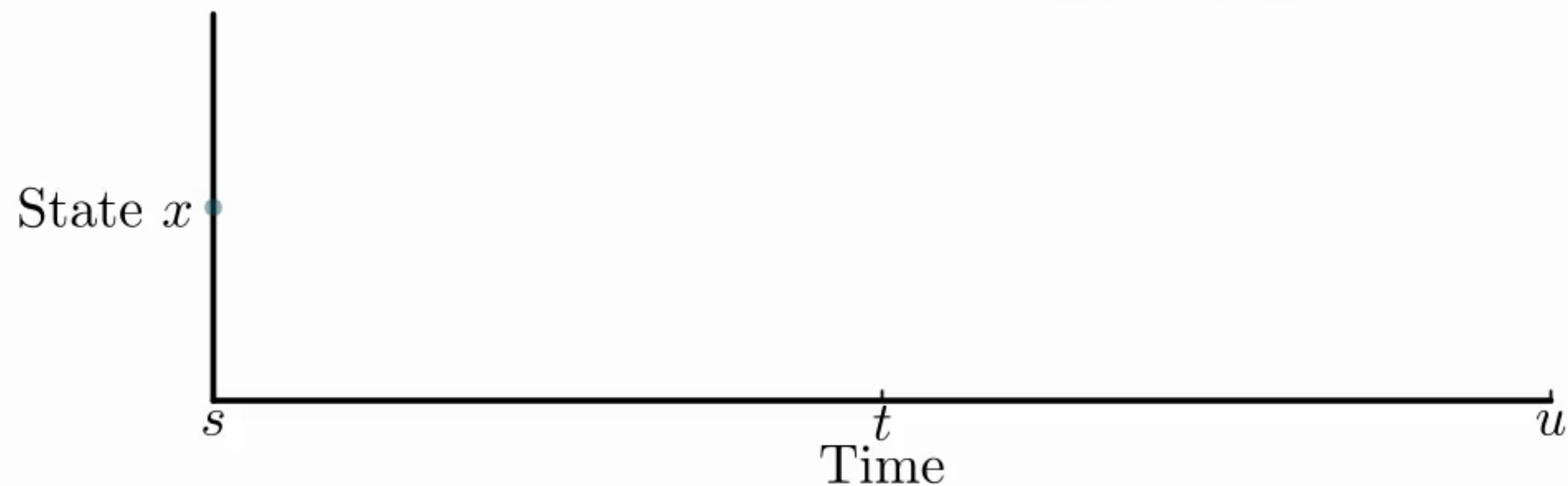
$$p_{\theta}(\mathbf{x}_{t_j} \mid \mathbf{x}_s; t_i, t_j - t_i) = p_{\theta}(\mathbf{x}_{t_j+r} \mid \mathbf{x}_s; t_i + r, t_j - t_i)$$

# Regularisation for Flow Consistency

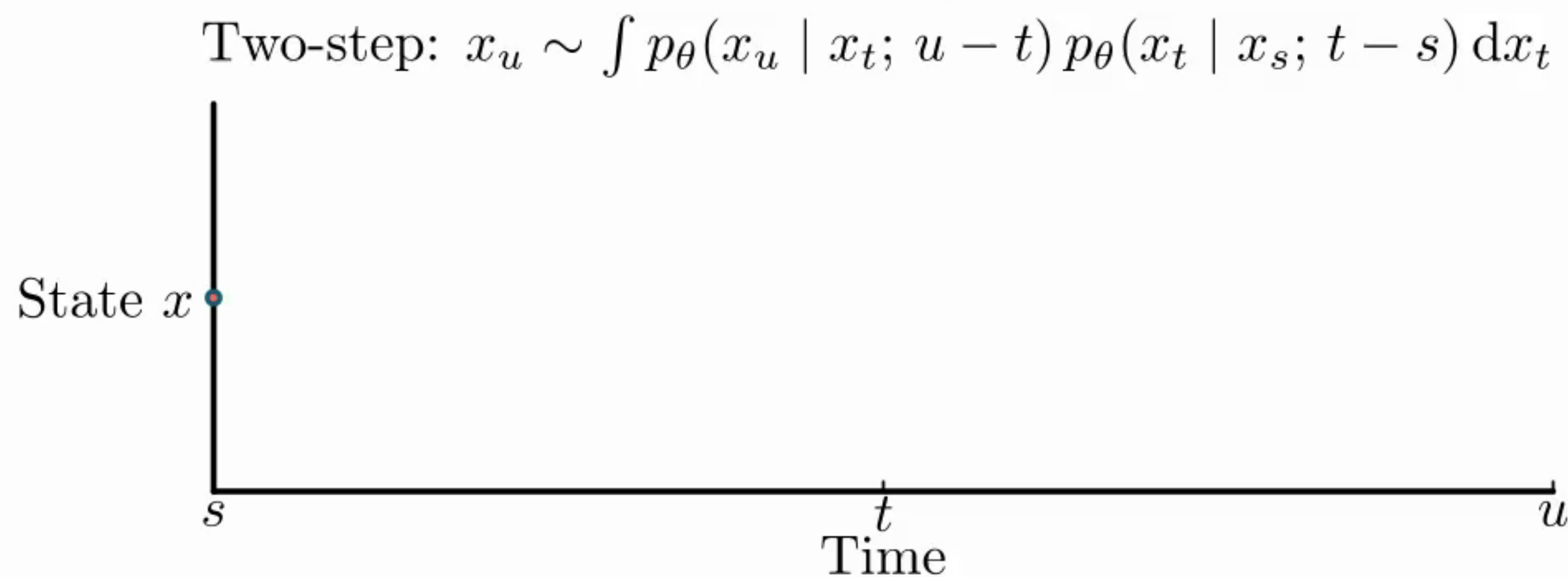
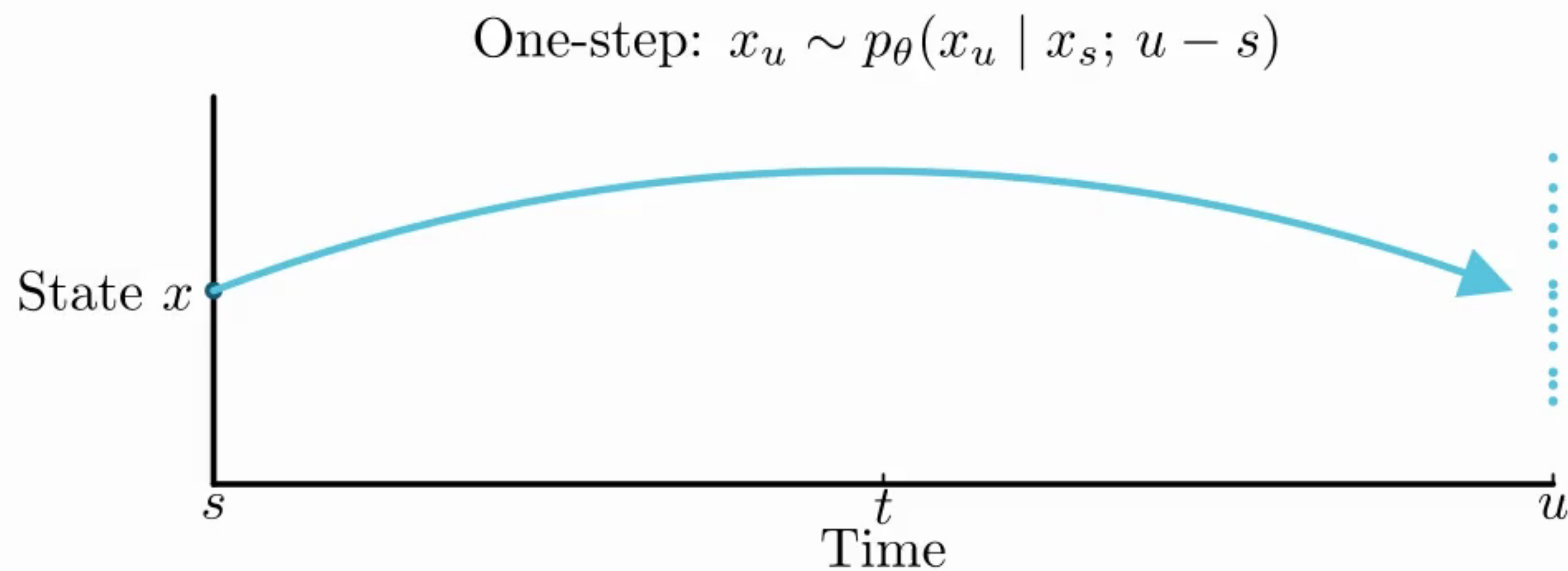
One-step:  $x_u \sim p_\theta(x_u \mid x_s; u - s)$



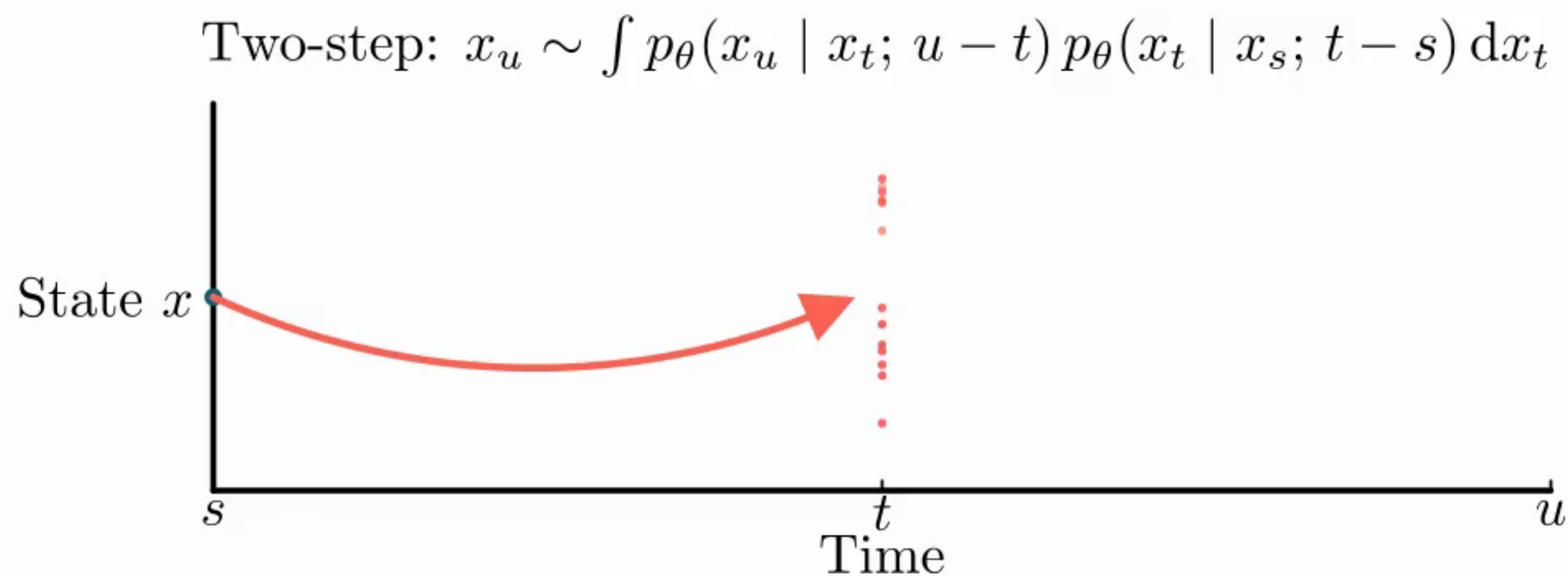
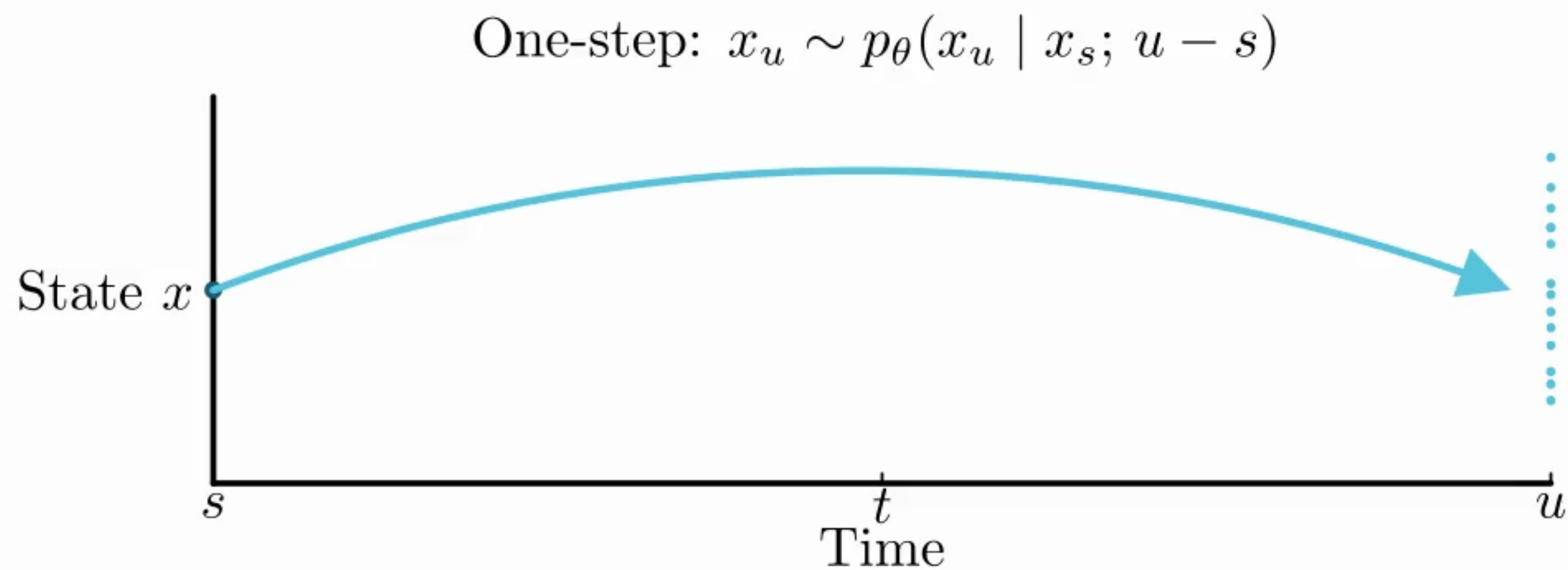
Two-step:  $x_u \sim \int p_\theta(x_u \mid x_t; u - t) p_\theta(x_t \mid x_s; t - s) dx_t$



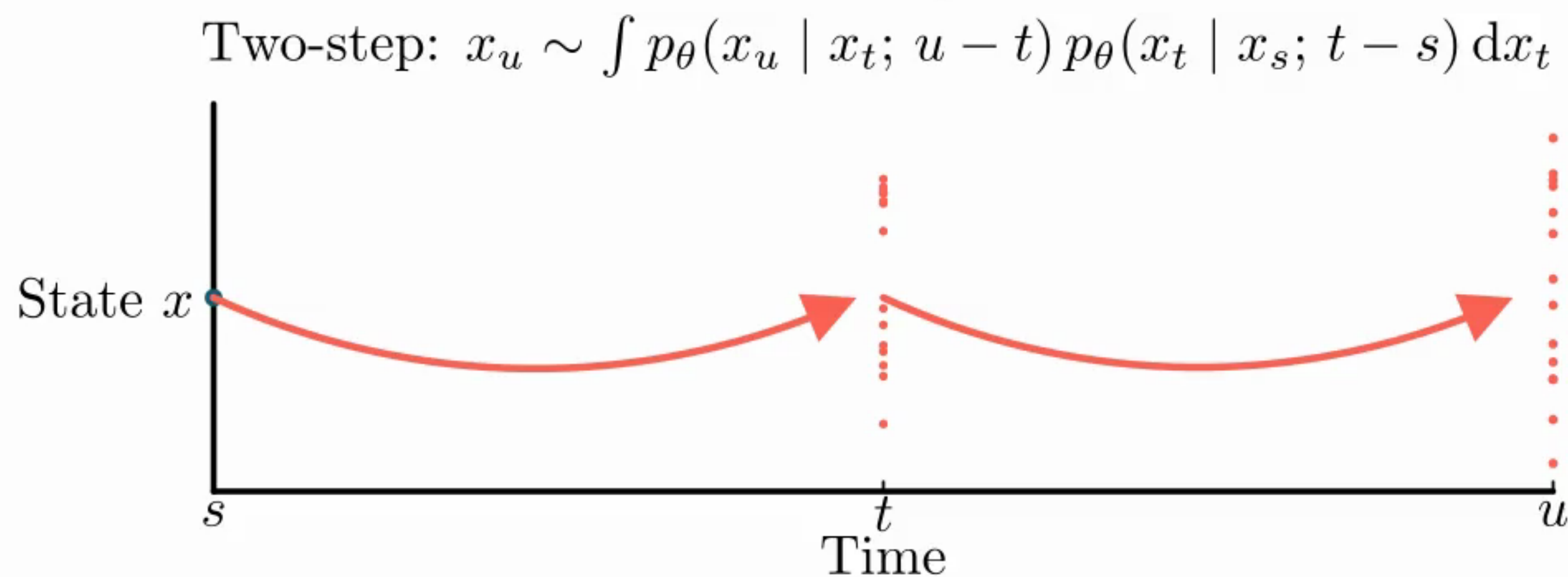
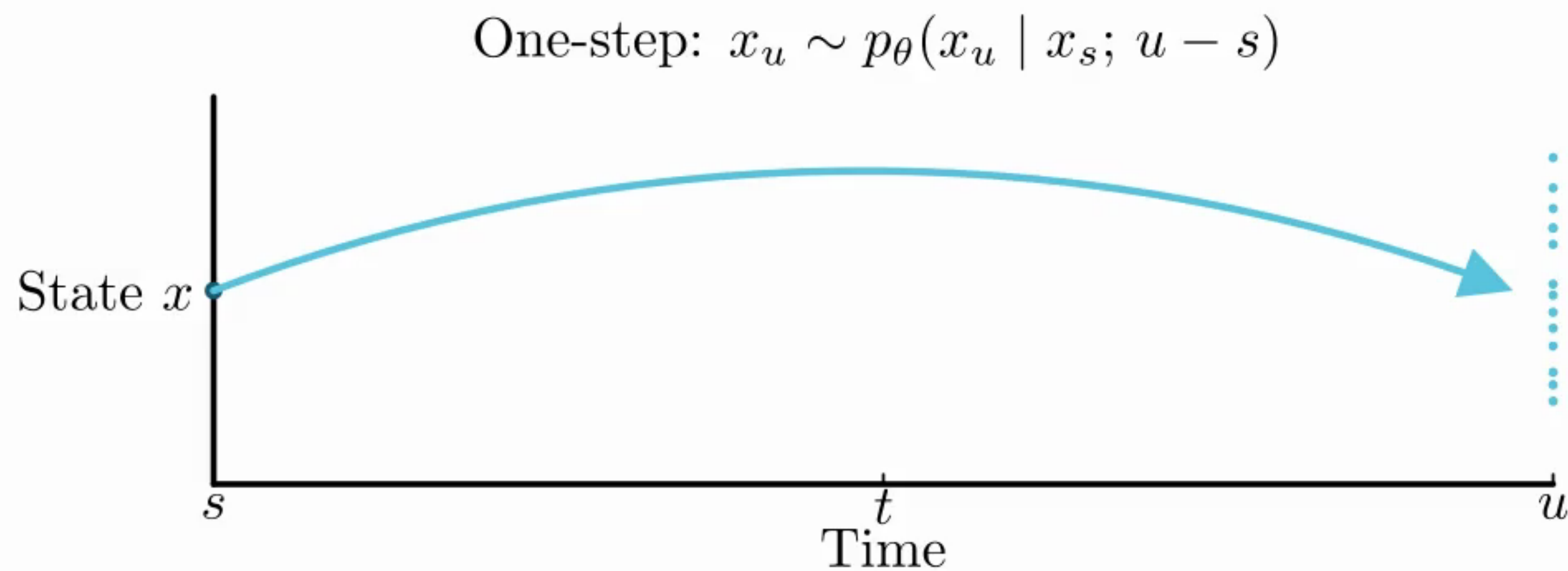
# Regularisation for Flow Consistency



# Regularisation for Flow Consistency

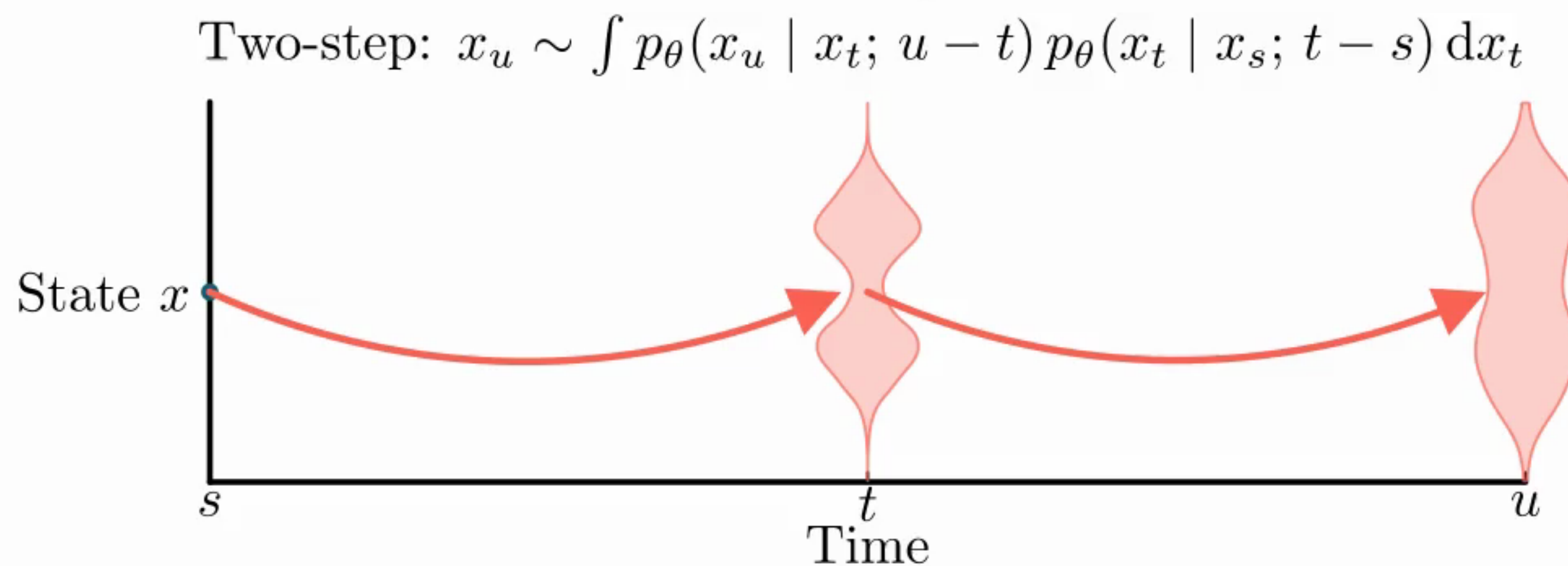
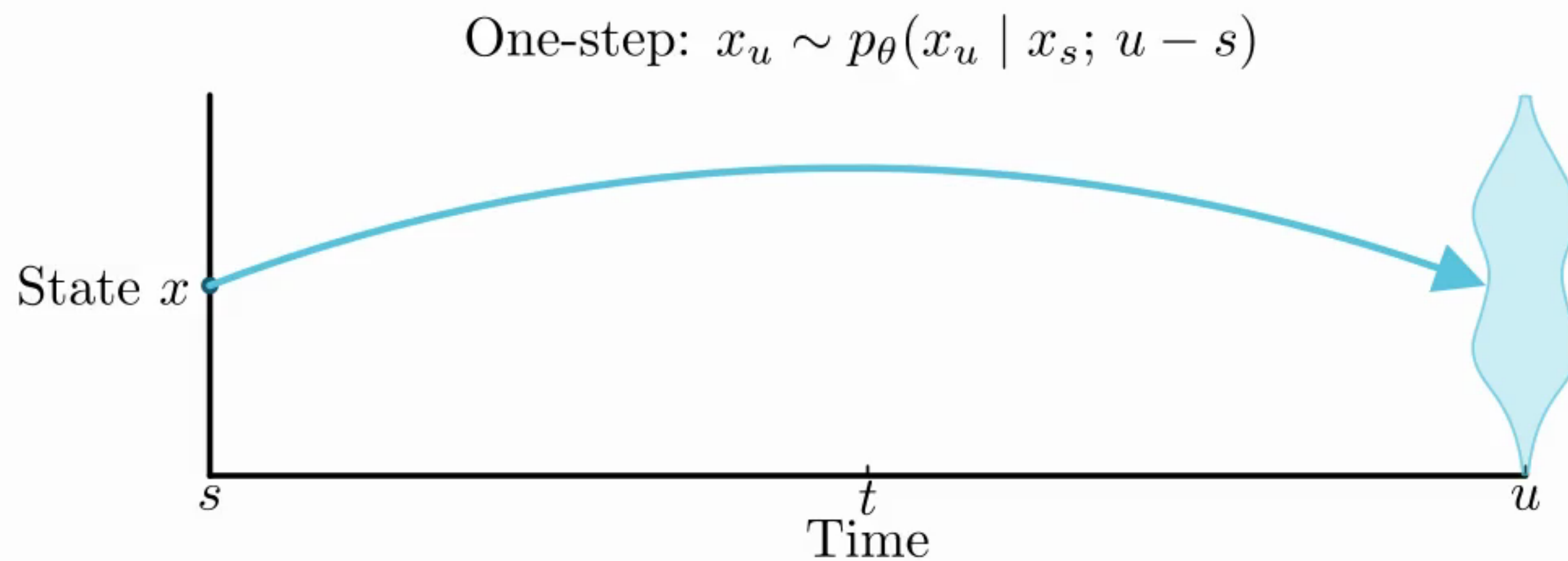


# Regularisation for Flow Consistency

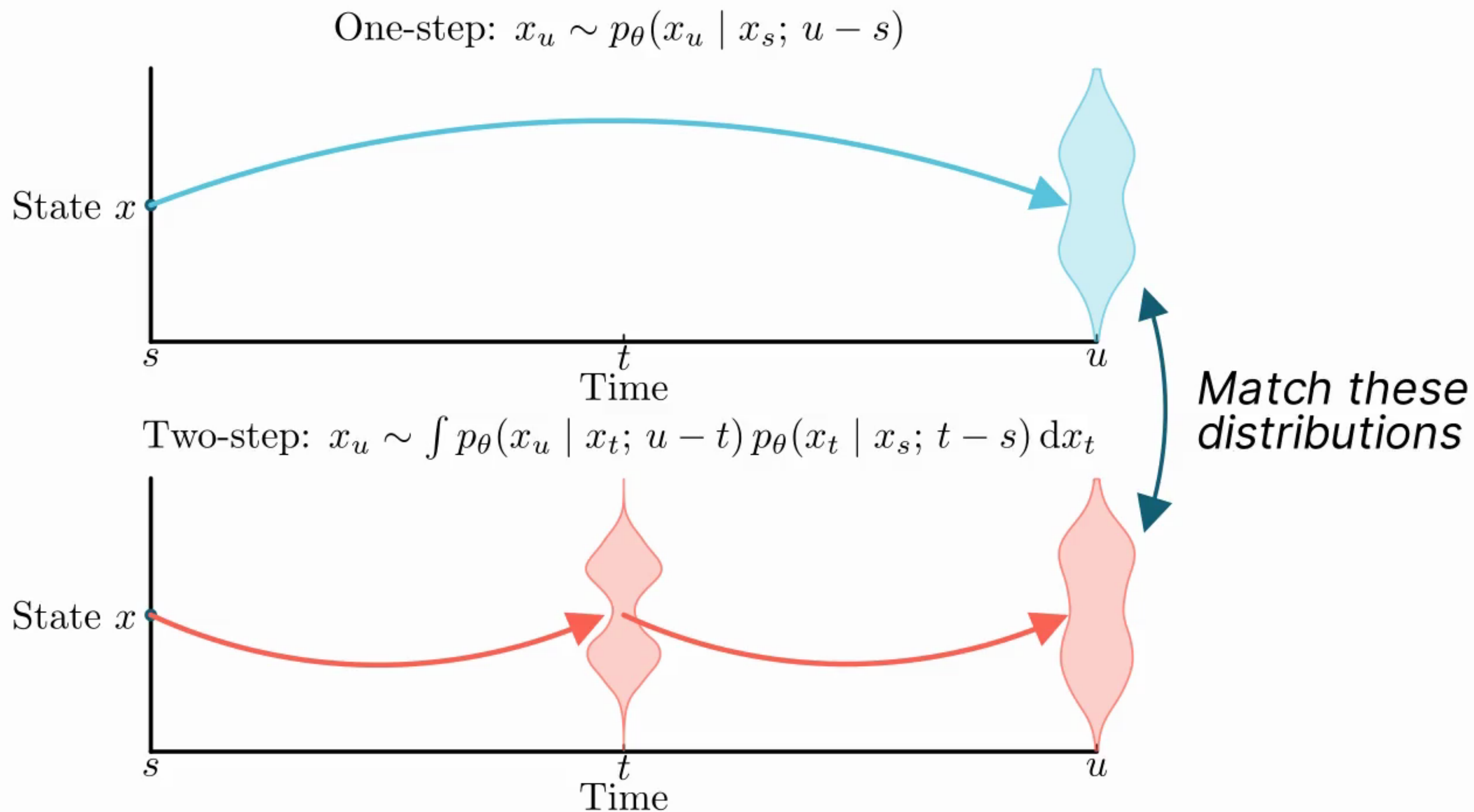




# Regularisation for Flow Consistency



# Regularisation for Flow Consistency



# Regularisation for Flow Consistency

$$\mathcal{L}_{\text{flow}} = D \left[ p_{\theta}^1 (\mathbf{x}_{t_k} \mid \mathbf{x}_{t_i}), \int p_{\theta}^2 (\mathbf{x}_{t_k} \mid \mathbf{x}_{t_j}) p_{\theta}^2 (\mathbf{x}_{t_j} \mid \mathbf{x}_{t_i}) d\mathbf{x}_{t_j} \right]$$

## Challenges in computing the flow loss

- The marginalisation is generally intractable  
→ **Direct access to the marginalised density is not available**
- Naive Monte Carlo estimation would require many samples



**We derive the upper bounds for KL divergences**

# Regularisation for Flow Consistency

## Upper bounds for KL divergences

### 1-step → 2-step KL divergence

$$\begin{aligned}
 D_{\text{KL}} \left( p_{\boldsymbol{\theta}}(\mathbf{x}_{t_k} \mid \mathbf{x}_{t_i}) \parallel \int p_{\boldsymbol{\theta}}(\mathbf{x}_{t_k} \mid \mathbf{x}_{t_j}) p_{\boldsymbol{\theta}}(\mathbf{x}_{t_j} \mid \mathbf{x}_{t_i}) d\mathbf{x}_{t_j} \right) \\
 \leq \mathbb{E}_{\mathbf{x}_{t_k} \sim p_{\boldsymbol{\theta}}(\cdot \mid \mathbf{x}_{t_i})} \left[ \mathbb{E}_{\mathbf{x}_{t_j} \sim b_{\boldsymbol{\xi}}(\cdot \mid \mathbf{x}_{t_i}, \mathbf{x}_{t_k})} \left[ \log \left( \frac{p_{\boldsymbol{\theta}}(\mathbf{x}_{t_k} \mid \mathbf{x}_{t_i}) b_{\boldsymbol{\xi}}(\mathbf{x}_{t_j} \mid \mathbf{x}_{t_i}, \mathbf{x}_{t_k})}{p_{\boldsymbol{\theta}}(\mathbf{x}_{t_k} \mid \mathbf{x}_{t_j}) p_{\boldsymbol{\theta}}(\mathbf{x}_{t_j} \mid \mathbf{x}_{t_i})} \right) \right] \right] =: \mathcal{L}_{\text{flow}, 1\text{-to-2}}(\boldsymbol{\theta}, \boldsymbol{\xi}; t_i, t_j, t_k).
 \end{aligned}$$

### 2-step → 1-step KL divergence

$$\begin{aligned}
 D_{\text{KL}} \left( \int p_{\boldsymbol{\theta}}(\mathbf{x}_{t_k} \mid \mathbf{x}_{t_j}) p_{\boldsymbol{\theta}}(\mathbf{x}_{t_j} \mid \mathbf{x}_{t_i}) d\mathbf{x}_{t_j} \parallel p_{\boldsymbol{\theta}}(\mathbf{x}_{t_k} \mid \mathbf{x}_{t_i}) \right) \\
 \leq \mathbb{E}_{\mathbf{x}_{t_j} \sim p_{\boldsymbol{\theta}}(\cdot \mid \mathbf{x}_{t_i})} \left[ \mathbb{E}_{\mathbf{x}_{t_k} \sim p_{\boldsymbol{\theta}}(\cdot \mid \mathbf{x}_{t_j})} \left[ \log \left( \frac{p_{\boldsymbol{\theta}}(\mathbf{x}_{t_j} \mid \mathbf{x}_{t_i}) p_{\boldsymbol{\theta}}(\mathbf{x}_{t_k} \mid \mathbf{x}_{t_j})}{b_{\boldsymbol{\xi}}(\mathbf{x}_{t_j} \mid \mathbf{x}_{t_i}, \mathbf{x}_{t_k}) p_{\boldsymbol{\theta}}(\mathbf{x}_{t_k} \mid \mathbf{x}_{t_i})} \right) \right] \right] =: \mathcal{L}_{\text{flow}, 2\text{-to-1}}(\boldsymbol{\theta}, \boldsymbol{\xi}; t_i, t_j, t_k).
 \end{aligned}$$

$$\mathcal{L}_{\text{flow}}(\boldsymbol{\theta}, \boldsymbol{\xi}) := \mathbb{E}_{p(t_i, t_j, t_k)} [\mathcal{L}_{\text{flow}, 1\text{-to-2}}(\boldsymbol{\theta}, \boldsymbol{\xi}; t_i, t_j, t_k) + \mathcal{L}_{\text{flow}, 2\text{-to-1}}(\boldsymbol{\theta}, \boldsymbol{\xi}; t_i, t_j, t_k)]$$

# Mathematical Requirements for the Transition Kernel

## Properties of $p_{\theta}(\mathbf{x}_t \mid \mathbf{x}_s; s, \Delta t = t - s)$

### 1. Independence. ✓ Realised by non-overlapping sampling.

For any  $0 \leq t_1 \leq \dots \leq t_n$ ,  $p_{\theta}(\mathbf{x}_{t_{k+1}} \mid \mathbf{x}_{t_k}; t_k, t_{k+1} - t_k)$  are independent.

### 2. Identity. ✓ Realised by architecture.

When  $t = s$ ,  $p_{\theta}(\mathbf{x}_t \mid \mathbf{x}_s; s, 0) = \delta(\mathbf{x}_t - \mathbf{x}_s)$

### 3. Flow property. ✓ Approximately realised by flow loss.

For any  $0 \leq t_i \leq t_j \leq t_k$ ,  $p_{\theta}(\mathbf{x}_{t_k} \mid \mathbf{x}_{t_i}; t_i, t_k - t_i) = \int p_{\theta}(\mathbf{x}_{t_k} \mid \mathbf{x}_{t_j}; t_j, t_k - t_j) p_{\theta}(\mathbf{x}_{t_j} \mid \mathbf{x}_{t_i}; t_i, t_j - t_i) d\mathbf{x}_{t_j}$

### 4. Stationarity (for autonomous SDEs). ✓ Realised by architecture optionally.

$$p_{\theta}(\mathbf{x}_{t_j} \mid \mathbf{x}_s; t_i, t_j - t_i) = p_{\theta}(\mathbf{x}_{t_j+r} \mid \mathbf{x}_s; t_i + r, t_j - t_i)$$

# Training Objective

$$\mathcal{L} = \mathbb{E}_{\mathbf{x}_{t_j}, \mathbf{x}_{t_i} \sim \mathcal{D}} [\log p_{\boldsymbol{\theta}}(\mathbf{x}_{t_j} \mid \mathbf{x}_{t_i})] + \lambda \mathcal{L}_{\text{flow}}(\boldsymbol{\theta}, \boldsymbol{\xi})$$

**Negative log likelihood**

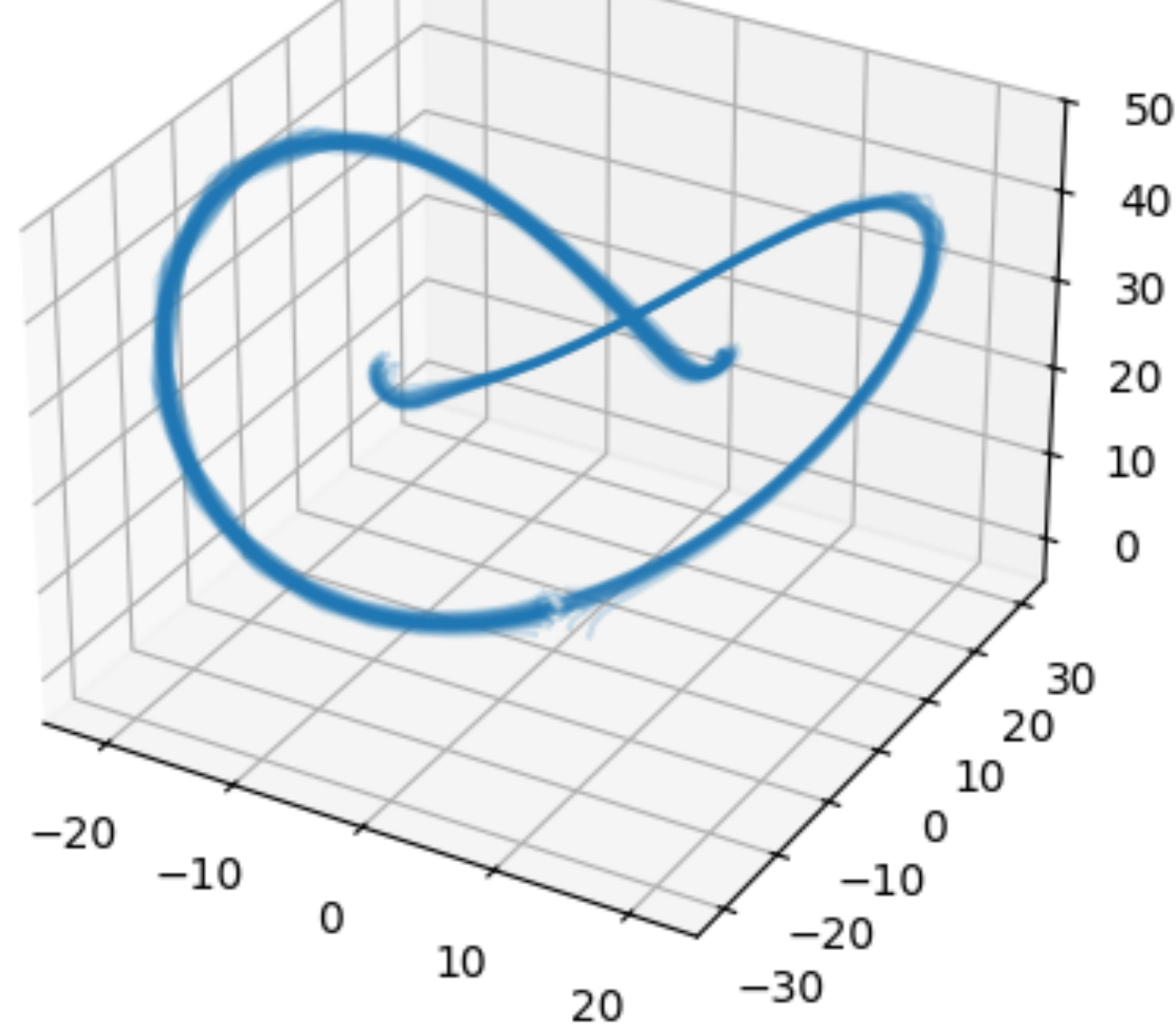
**Flow loss**



# Experimental Configurations

## Chaotic dynamics—Stochastic Lorenz Attractor

Ground truth trajectories



$$dx = \sigma(y - x)dt + \alpha_x dW_1$$

$$dy = (x(\rho - z) - y)dt + \alpha_y dW_2$$

$$dz = (xy - \beta z)dt + \alpha_z dW_3$$

$$\sigma = 10, \rho = 28, \beta = 8/3, \alpha = (0.15, 0.15, 0.15)$$

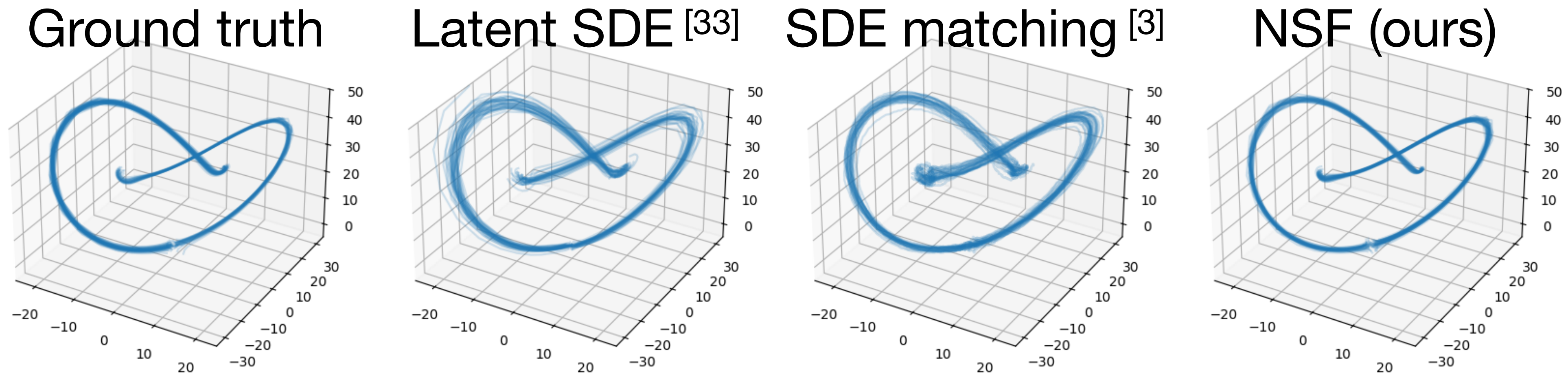
Initial state:  $\mathcal{N}(0, I)$ ; Duration:  $t \in [0, 1]$

**Trained on 1,024 trajectories, tested on 1,024 trajectories.**



# Experimental Results

## With a comparable network size



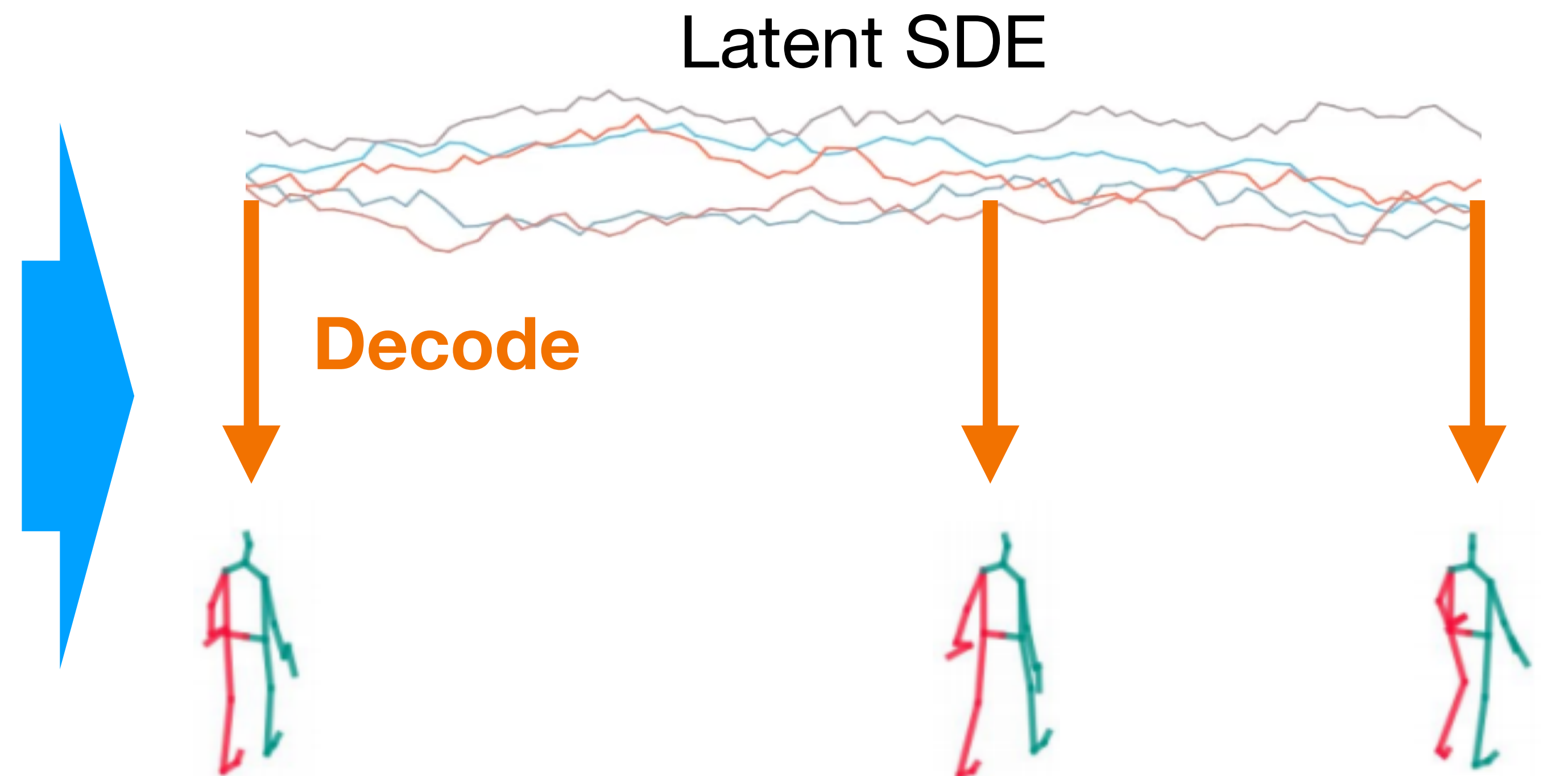
Method	t = 0.25		t = 0.5		t = 0.75		t = 1.0	
	KL ↓	kFLOPs ↓	KL ↓	kFLOPs ↓	KL ↓	kFLOPs ↓	KL ↓	kFLOPs ↓
<b>Latent SDE [33]</b>	$2.1 \pm 0.9$	959	$1.8 \pm 0.1$	1,917	$0.9 \pm 0.3$	2,839	$1.5 \pm 0.5$	3,760
<b>SDE matching [3]</b>	$6.3 \pm 0.4$	1,917	$11.7 \pm 0.5$	3,834	$7.9 \pm 0.3$	5,677	$6.0 \pm 0.3$	7,520
<b>NSF (ours)</b>	<b><math>0.8 \pm 0.7</math></b>	<b>53</b>	<b><math>1.3 \pm 0.1</math></b>	<b>53</b>	<b><math>0.6 \pm 0.3</math></b>	<b>53</b>	<b><math>0.2 \pm 0.6</math></b>	<b>53</b>

# Latent SDEs and Latent Neural Stochastic Flows

Extension for partially-observed/high-dimensional data

Real-world time-series data is often:

- Partially-observed
- Irregularly-sampled
- High-dimensional



(Figure adapted from [github.com/jutanke/mocap](https://github.com/jutanke/mocap))

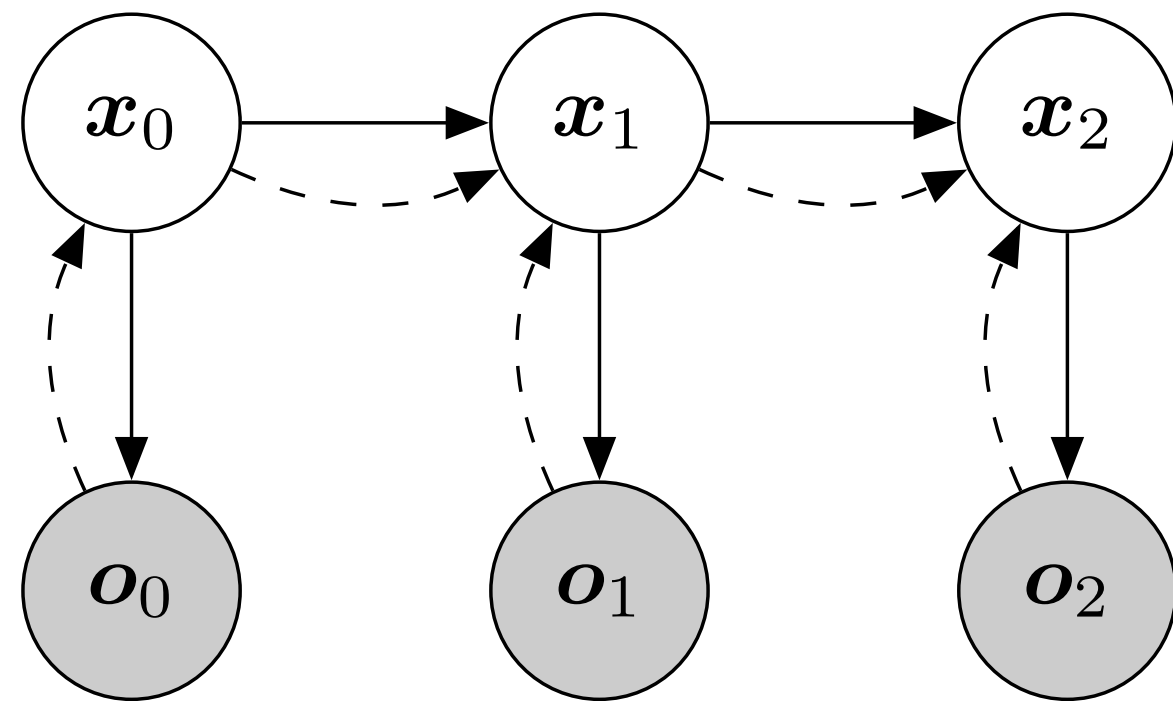
**Latent SDEs are powerful tools, but they require high computational costs.**

**We replace SDEs with Neural Stochastic Flows.**

# Latent Neural Stochastic Flows

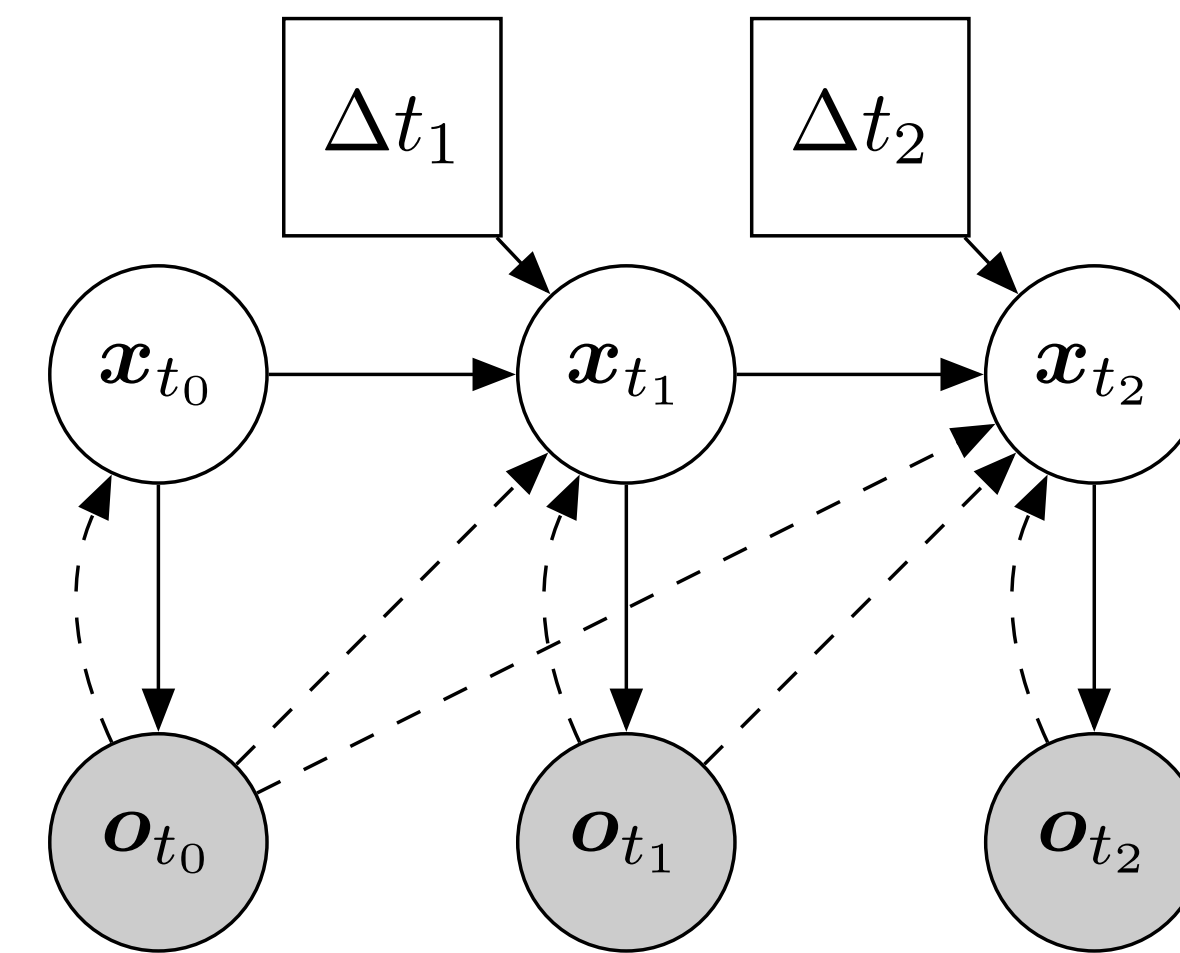
Extension for partially-observed/high-dimensional data

**Variational State Space Model [8, 21, 22, 31]**



A proven template for sequential modelling, including RL (PlaNet<sup>[16]</sup>, Dreamer, etc.)

**Latent NSF (ours)**



- ✓ **VSSM structure remains the same by NSF's closed-form log density**
- ✓ **One-step transition for arbitrary  $\Delta t$**



# Experimental Results with Latent NSF

Real-world/high-dimensional sequential prediction

CMU Motion Capture (50-dims.)

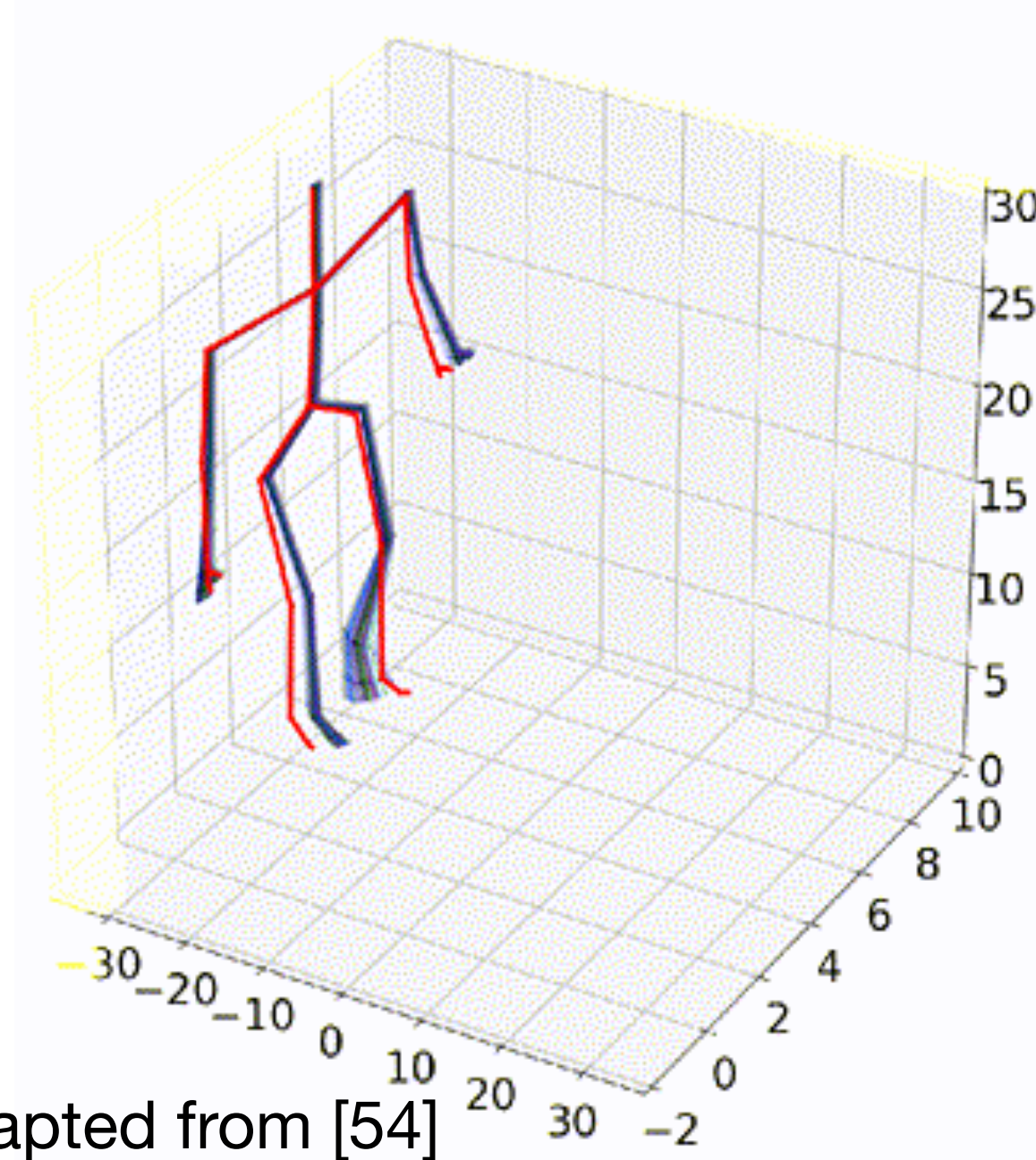


Figure adapted from [54]

Setup 1: The same length for train/test

Setup 2: Training length < Test length

Stochastic Moving MNIST (64x64 px)



# Experimental Results with Latent NSF

## Real-world/high-dimensional sequential prediction

### CMU Motion Capture (50-dims.)

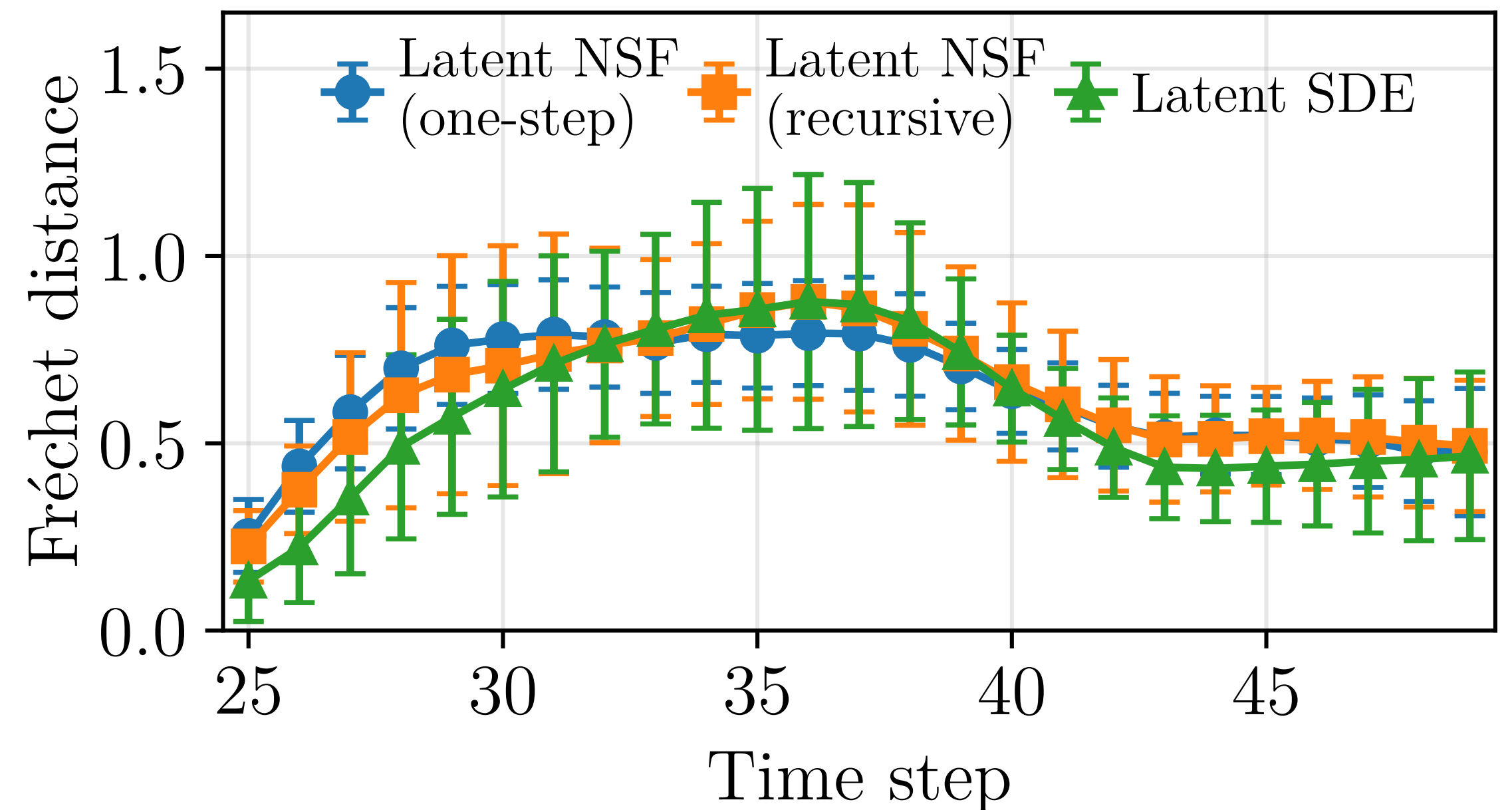
Compared by mean square error

Methods	Setup 1	Setup 2
npODE [17]	22.96 <sup>†</sup>	—
Neural ODE [6]	22.49 ± 0.88 <sup>†</sup>	—
ODE2VAE-KL [54]	8.09 ± 1.95 <sup>†</sup>	—
Latent ODE [43]	5.98 ± 0.28 <sup>*</sup>	31.62 ± 0.05 <sup>§</sup>
Latent SDE [33]	12.91 ± 2.90 <sup>1</sup>	9.52 ± 0.21 <sup>§</sup>
Latent Approx SDE [46]	7.55 ± 0.05 <sup>§</sup>	10.50 ± 0.86
ARCTA [9]	7.62 ± 0.93 <sup>‡</sup>	9.92 ± 1.82
NCDSSM [2]	5.69 ± 0.01 <sup>§</sup>	4.74 ± 0.01 <sup>§</sup>
SDE matching [3]	<b>5.20 ± 0.43<sup>2</sup></b>	4.26 ± 0.35
Latent NSF (ours)	8.62 ± 0.32	<b>3.41 ± 0.27</b>

✓ **Best/tied-best performance**

✓ **Two orders of magnitude faster**

### Stochastic Moving MNIST (64x64 px)



✓ **Comparable performance**

✓ **2x faster**



# Position of Our Work

Method(s)	Target dynamics	Solver-free training	Solver-free inference
Neural ODEs [6]	General ODEs	✗	✗
Neural flows [4]	General ODEs	✓	✓
Score-based diffusion via reverse SDEs/PF-ODEs [47]; flow matching [34]	Pre-defined diffusion SDEs/ODEs	✓	✗
Progressive distillation [44]	Pre-defined diffusion SDEs/ODEs	✗	✓
Consistency models [26, 48]; rectified flows [35]	Pre-defined diffusion SDEs/ODEs	✓	✓
Neural (latent) SDEs [25, 33, 38, 46, 50]	<b>General Itô SDEs</b>	✗	✗
ARCTA [9]; SDE matching [3]	<b>General Itô SDEs</b>	✓	✗
<b>Neural Stochastic Flows (ours)</b>	<b>General Itô SDEs</b>	✓	✓

# Summary

- We proposed **Neural Stochastic Flows**, which enable solver-free in both training and inference, with closed-form log-densities.
- Extended to **Latent NSF**: a variational state-space model using NSF as the latent transition for partially-observed/high-dimensional data.

## Key Benefits:

- **One-step prediction** for arbitrary time gaps.
- **Speed**: up to two orders of magnitude faster.
- **Stable training** with tractable densities.

## Limitations:

- Chapman–Kolmogorov relation is enforced **approximately**.
- Affine-coupling design constrains **expressivity**.

Project  
page



Paper  
&  
Refs.

