

---

# Scalable Gradients for Stochastic Differential Equations

---

**Xuechen Li\***  
Google Research

**Ting-Kam Leonard Wong**  
University of Toronto

**Ricky T. Q. Chen**  
University of Toronto  
Vector Institute

**David Duvenaud**  
University of Toronto  
Vector Institute

Published at AISTAT 2020

# Content

- Recap of SDE
- Recap of Neural ODE and adjoint method
- Stochastic adjoint method
- Latent-SDE and Bayesian learning(VI)
- Other

# Stochastic Differential Equations

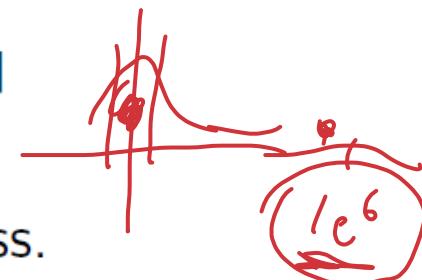
- Stochastic differential equations (SDEs) generalize ODEs to incorporate a noise component driven by a Brownian motion

$$Z_T = Z_0 + \int_0^T \underbrace{b(Z_t, t)}_{\text{drift}} dt +$$

$$\underbrace{\int_0^T \underbrace{\sigma(Z_t, t)}_{\text{diffusion}} dW_t}_{\text{Itô integral}}$$

$dW_t$ : White noise  
(a random sample from zero-mean Gauss)

- $b : \mathbb{R}^d \times \mathbb{R} \rightarrow \mathbb{R}^d$  and  $\sigma : \mathbb{R}^d \times \mathbb{R} \rightarrow \mathbb{R}^{d \times m}$  are the drift and diffusion coefficients, respectively.
- $W_t$  is a (multi-dimensional) Brownian motion/Wiener process.
- SDEs are widely used in mathematical finance, population genetics, stochastic optimal control, etc.



$$\frac{dz_t}{dt} = b(z_t, t) + \sigma(z_t, t) dw_t$$

# Lebesgue–Stieltjes Integral vs Itô Integral

Lebesgue–Stieltjes

$$\int_a^b f(t) dg(t) = \lim_{|\Pi| \rightarrow 0} \sum f(t_i^*) (g(t_{i+1}) - g(t_i))$$

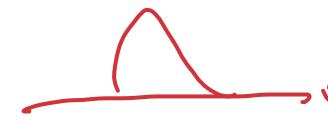
- $t_i^* \in [t_i, t_{i+1})$
- $g(t)$  of bounded variation
- Obeys usual integration by parts rule
- Convergence is in usual metric on  $\mathbb{R}$

Itô

$$\int_a^b f(t) dW_t = \lim_{|\Pi| \rightarrow 0} \sum f(t_i^*) (W_{t_{i+1}} - W_{t_i})$$

$$dW_t$$
  
$$|W_t - W_{t+\delta}|$$

- $t_i^* = t_i$
- $W_t$  not of bounded variation
- Doesn't obey usual integration by parts rule
- Convergence is in  $L^2$  and probability



# What kind of solutions do SDEs have?

$x(t)$

- **Path of solution:** Draw random path  $w(t)$  (or  $\beta(t)$ ) and solve the equation using it as the input.

- Monte Carlo simulation of SDE solutions.
- Used in particle filtering and smoothing methods.



- **Distribution of solution:** Given many random  $w(t)$ 's, what is the distribution of the state  $p(x(t))$ ?

- Solution is given by the Fokker-Planck-Kolmogorov PDE.
- Used in grid based and basis function methods (FEM, BEM).



- **Moments:** What are the mean and covariance of  $x(t)$ ?

- Ordinary differential equations for the mean and covariance.
- Used in non-linear Kalman (Gaussian) filters and smoothers.



## What does a solution of SDE look like? (cont.)

**Example 3.3** (Solution of Ornstein–Uhlenbeck equation). *The complete solution to the scalar SDE*

$$dx = (-\lambda x \, dt + d\beta), \quad x(0) = x_0, \quad (3.32)$$

where  $\lambda > 0$  is a given constant and  $\beta(t)$  is a Brownian motion is

$$x(t) = \exp(-\lambda t) x_0 + \int_0^t \exp(-\lambda(t-\tau)) \, d\beta(\tau). \quad (3.33)$$

The solution, called the Ornstein–Uhlenbeck process, is illustrated in Figure 3.2.

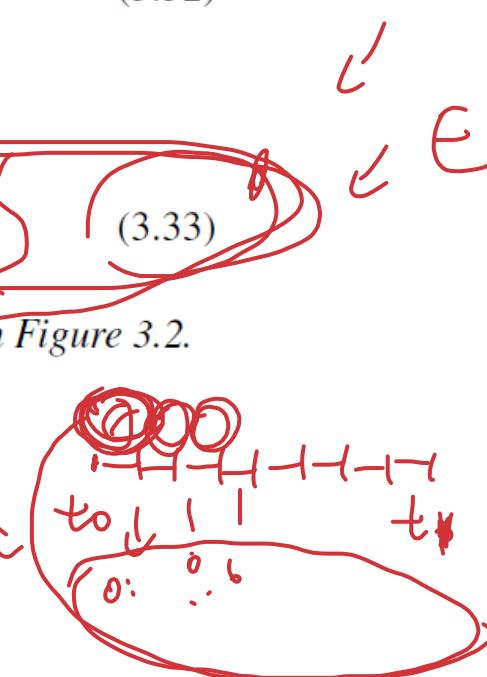
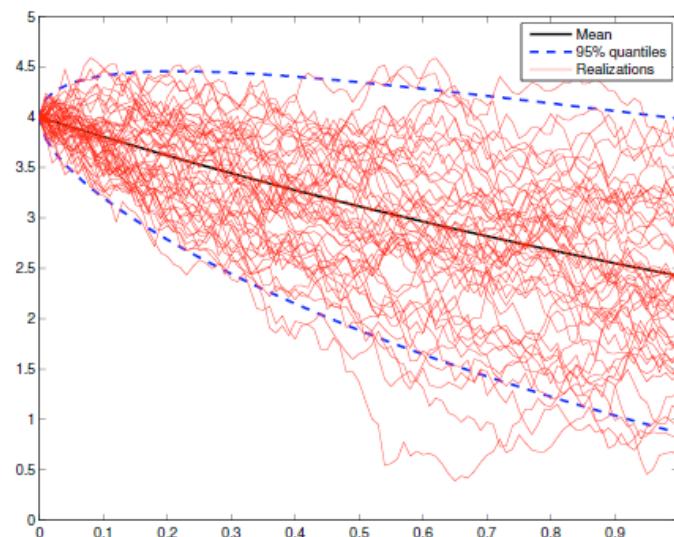
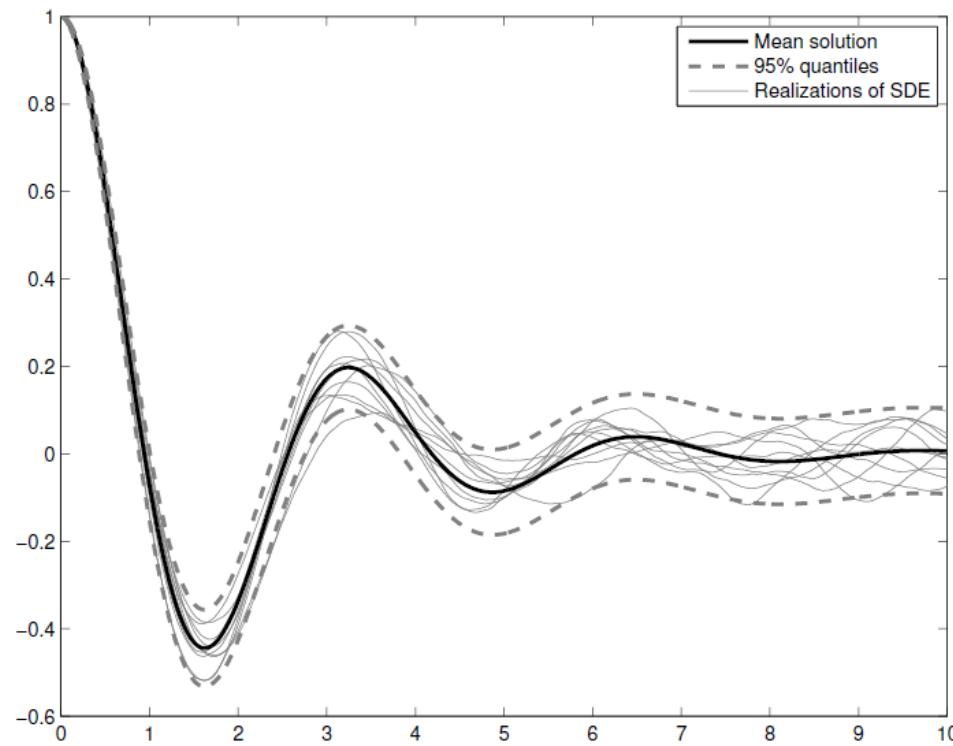


Figure 3.2: Realizations, mean, and 95% quantiles of Ornstein–Uhlenbeck process.

## What does a solution of SDE look like? (cont.)



$$\frac{d^2x(t)}{dt^2} + \gamma \frac{dx(t)}{dt} + \nu^2 x(t) = w(t).$$

$x_1 = x$

$x_2 = \frac{dx}{dt}$

**Example 2.8** (Stochastic Spring model). If in the spring model of Equation (1.4), we replace the input force with a white noise with spectral density  $q$ , we get the following LTI SDE:

$$\underbrace{\begin{pmatrix} \frac{dx_1(t)}{dt} \\ \frac{dx_2(t)}{dt} \end{pmatrix}}_{\mathbf{dx}(t)/dt} = \underbrace{\begin{pmatrix} 0 & 1 \\ -\nu^2 & -\gamma \end{pmatrix}}_{\mathbf{F}} \underbrace{\begin{pmatrix} x_1(t) \\ x_2(t) \end{pmatrix}}_{\mathbf{x}} + \underbrace{\begin{pmatrix} 0 \\ 1 \end{pmatrix}}_{\mathbf{L}} w(t). \quad (2.39)$$

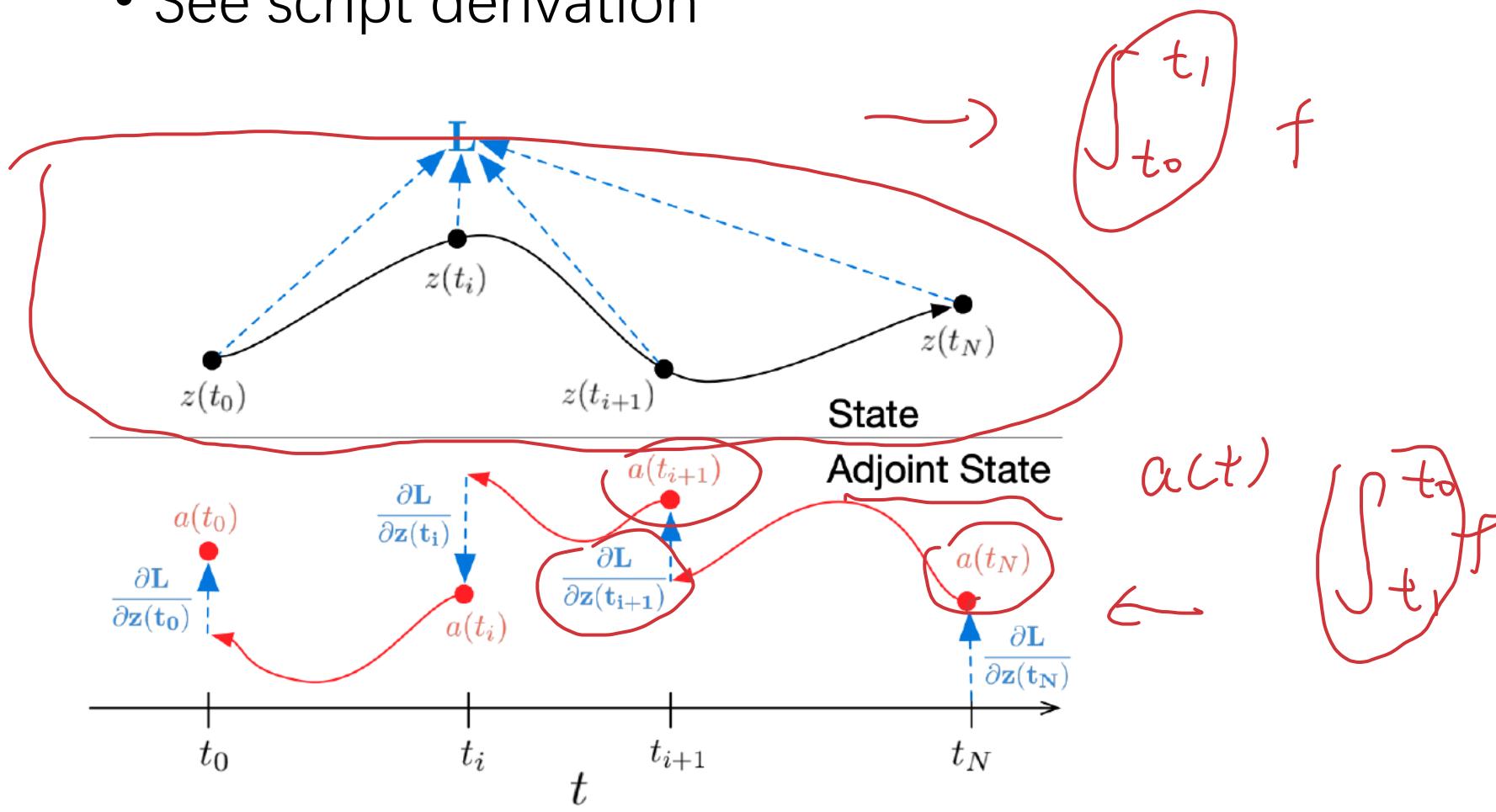
The equations for the mean and covariance are thus given as

$$\begin{aligned} \begin{pmatrix} \frac{dm_1}{dt} \\ \frac{dm_2}{dt} \end{pmatrix} &= \begin{pmatrix} 0 & 1 \\ -\nu^2 & -\gamma \end{pmatrix} \begin{pmatrix} m_1 \\ m_2 \end{pmatrix} \\ \begin{pmatrix} \frac{dP_{11}}{dt} & \frac{dP_{12}}{dt} \\ \frac{dP_{21}}{dt} & \frac{dP_{22}}{dt} \end{pmatrix} &= \begin{pmatrix} 0 & 1 \\ -\nu^2 & -\gamma \end{pmatrix} \begin{pmatrix} P_{11} & P_{12} \\ P_{21} & P_{22} \end{pmatrix} \\ &+ \begin{pmatrix} P_{11} & P_{12} \\ P_{21} & P_{22} \end{pmatrix} \begin{pmatrix} 0 & 1 \\ -\nu^2 & -\gamma \end{pmatrix}^T + \begin{pmatrix} 0 & 0 \\ 0 & q \end{pmatrix} \end{aligned} \quad (2.40)$$

When  $\mathbf{F}$  and  $\mathbf{L}$  are non-constant functions, always hard to get such a solution.

# Recap Neural ODE and adjoint method

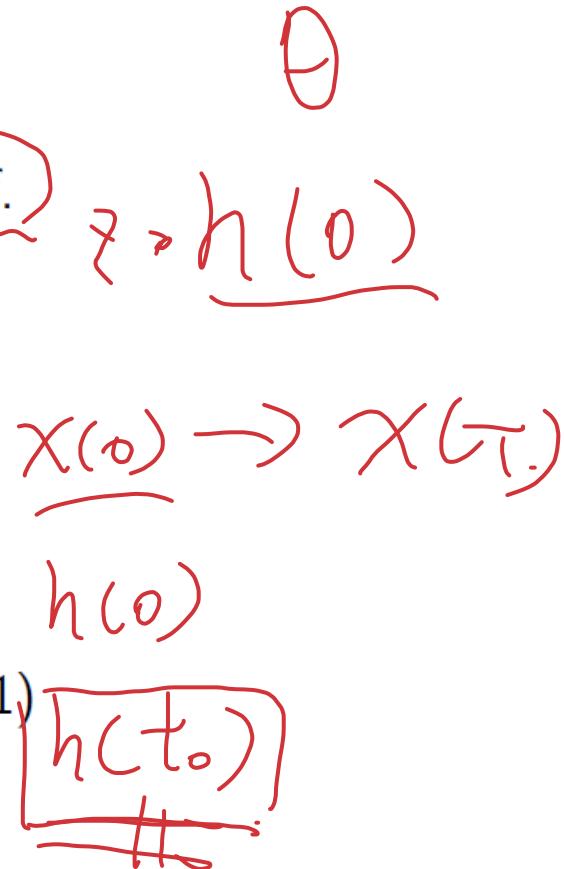
- See script derivation



# Combining Adjoints with Reverse-mode Autodiff

- Chen et al. [1] combined adjoints with reverse-mode autodiff.
- Define loss that is a function of the end state  $\mathcal{L} = \mathcal{L}(z(T))$ .
- Define vector field  $f(z, t, \theta)$ , where parameter  $\theta$  is fixed.
- $\theta$  is the control in the maximum principle framework.
- Define adjoint  $a(t) = \partial \mathcal{L} / \partial z(t)$ ; it satisfies another ODE:

$$\frac{\partial \mathcal{L}}{\partial z(0)} = a(0) = a(T) - \int_T^0 a(t) \underbrace{\frac{\partial f(z, t, \theta)}{\partial z}}_{\text{vJp}} dt.$$

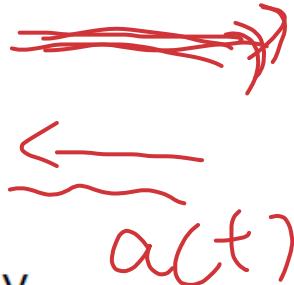


- Term in blue is evaluated as vector-Jacobian product (vJp).
- vJps are almost as cheap as evaluating the function itself; much cheaper than evaluating the full Jacobian!

Nature extension:

- Use **NNs** to parameterize a **SDE** ?-easy and trivial
- Still use **original adjoint** to optimize? – non-trivial, very difficult

# Generalizing Adjoints to SDEs is Difficult



- To machine learners, we want a method that scales gracefully in terms of both memory and time.
- For ODEs, the adjoint sensitivity method achieved this by
  - i) solving adjoint dynamics in reverse time (memory), and
  - ii) relying on vJps (time) in numerical solve.
- For SDEs, it is less clear what “solving (pathwise) in reverse time” means.
- The classical approach to SDEs is built upon Itô calculus,
  - a) doesn't obey usual rules of integration (e.g. the usual integration by parts), and
  - b) rarely considers sample paths individually.

# Challenge 1: Simulating SDEs Backward in Time

- One alternative way to think about SDEs is to reason about the pathwise behavior of solutions and consider the induced *stochastic flow*.  
A
- Just like how ODE flows can be diffeomorphisms, so can SDE flows be, once given a realization of Brownian motion.
- When the SDE is in *Stratonovich form*, Kunita [4] showed a simple *time-reverse SDE*<sup>1</sup> has a solution that corresponds to the *inverse flow* of the original SDE.  
A
- SDEs in Stratonovich form can be converted into Itô form, and vice versa.

# Challenge 1: Simulating SDEs Based on Wiener Process

Some advanced definitions in advanced random process and differential geometry, no need to go into such deep..

- One alternative way to think about SDEs is the pathwise behavior of solutions and consider the *stochastic flow*.
- Just like how ODE flows can be diffeomorphisms, stochastic flows be, once given a realization of Brownian motion.
- When the SDE is in *Stratonovich form*, Kunita [1981] shows that a simple time-reverse SDE<sup>1</sup> has a solution that corresponds to the inverse flow of the original SDE.
- SDEs in Stratonovich form can be converted into Itô form, and vice versa.

solution. Let  $\Phi_{s,t}(z) := Z_t^{s,z}$  be the solution at time  $t$  when the process is started at  $z$  at time  $s$ . Given a realization of the Wiener process, this defines a collection of continuous maps  $\mathcal{S} = \{\Phi_{s,t}\}_{s \leq t; s, t \in \mathbb{T}}$  from  $\mathbb{R}^d$  to itself.

**Theorem 2.1** ([41, Theorem 3.7.1]). (a) With probability 1, the collection  $\mathcal{S} = \{\Phi_{s,t}\}_{s \leq t; s, t \in \mathbb{T}}$  satisfies the flow property

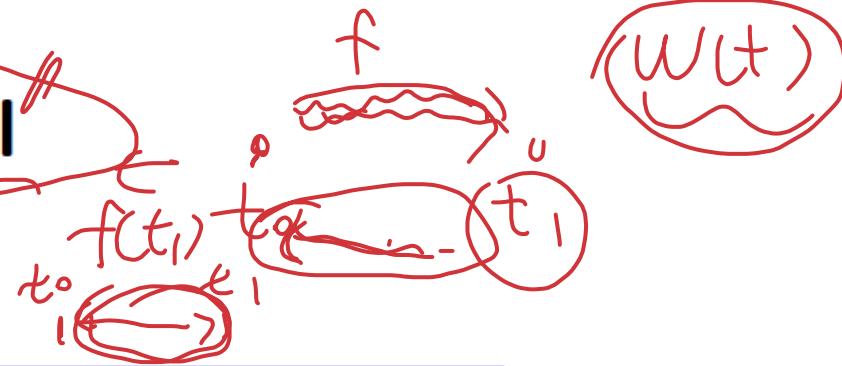
$$\Phi_{s,t}(z) = \Phi_{u,t}(\Phi_{s,u}(z)), \quad s \leq u \leq t, z \in \mathbb{R}^d.$$

The backward flow  $\check{\Psi}_{s,t} := \Phi_{s,t}^{-1}$  satisfies the backward SDE:

$$\begin{aligned} \check{\Psi}_{s,t}(z) &= z - \int_s^t b(\check{\Psi}_{u,t}(z), u) du - \\ &\quad \int_s^t \sigma(\check{\Psi}_{u,t}(z), u) \circ d\check{W}_u, \end{aligned} \quad (3)$$

for all  $z \in \mathbb{R}^d$  and  $s, t \in \mathbb{T}$  such that  $s \leq t$ .

# Stratonovich Integral vs Itô Integral



Stratonovich

$$\int_a^b f(t) \circ dW_t = \lim_{|\Pi| \rightarrow 0} \sum \frac{f(t_i) + f(t_{i+1})}{2} (W_{t_{i+1}} - W_{t_i})$$

$\overbrace{\text{to } t_0 \text{ to } t_1 \dots}$

Itô

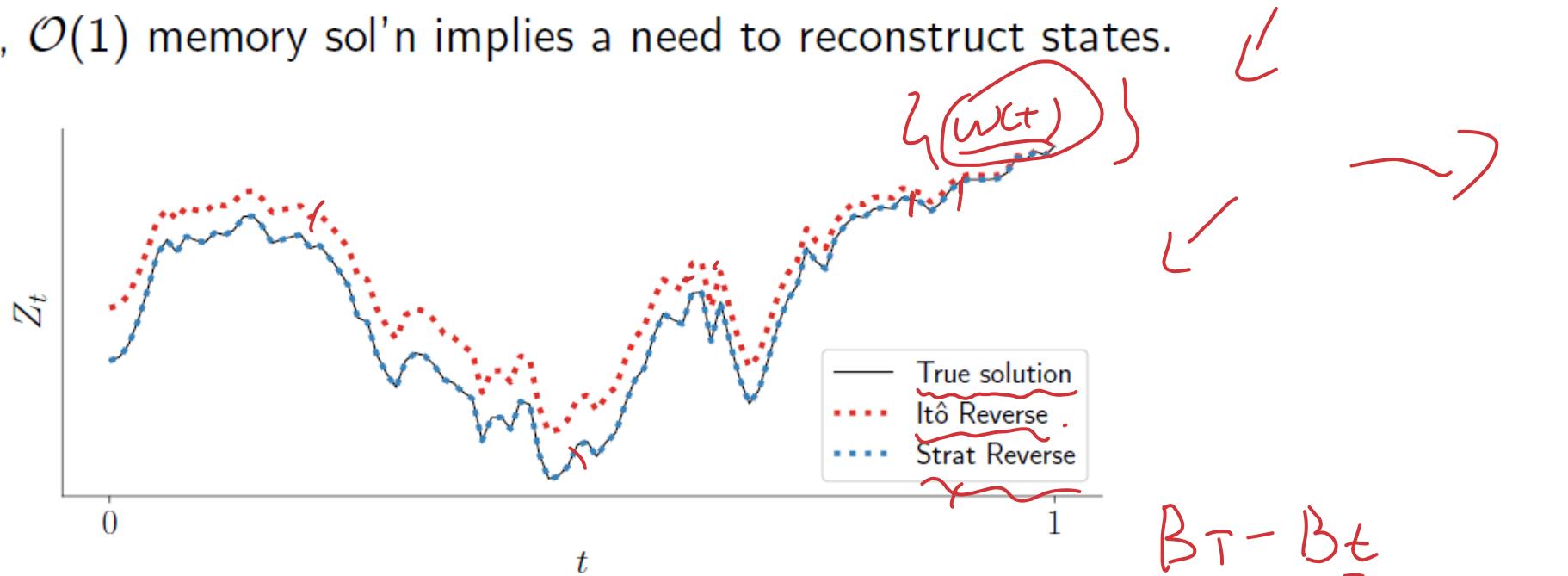
$$\int_a^b f(t) dW_t = \lim_{|\Pi| \rightarrow 0} \sum f(t_i^*) (W_{t_{i+1}} - W_{t_i})$$

- Average of two ends
- $W_t$  not of bounded variation
- **Obeys usual integration by parts rule**
- Convergence is in  $L^2$  and probability

- $t_i^* = t_i$        $f(t_i) f(t_{i+1})$
- $W_t$  not of bounded variation
- Doesn't obey usual integration by parts rule
- Convergence is in  $L^2$  and probability

# Challenge 1: Simulating Backward in Time

Recall,  $\mathcal{O}(1)$  memory sol'n implies a need to reconstruct states.



$$\check{\Psi}_{s,t}(z) = z - \int_s^t b(\check{\Psi}_{r,t}(z), r) dr - \sum_{i=1}^m \int_s^t \sigma_i(\check{\Psi}_{r,t}(z), r) \circ \check{dW}_r^{(i)}.$$

$\check{\Psi}_{s,t}(z)$  (highlighted in red circle)

Stratonovich integral

Time-rev. BM

$B_T - B_t$

$B_t - B_{T-}$

# Extending the Formulation to Adjoints

- With the time-reverse formulation, the task of extending adjoints to SDEs becomes that of rewriting the derivation for ODEs, but with Stratonovich integrals and the time-reverse Wiener process.
- We can still define the adjoint to be the derivative of the loss random variable w.r.t. the state, i.e.  $A_t = \partial \mathcal{L} / \partial Z_t$ , and show it is the sol'n to a time-reverse SDE:

$$\begin{aligned}\check{A}_{s,t}(z) &= \nabla \mathcal{L}(z) + \int_s^t \underbrace{\check{A}_{r,t}(z) \nabla b(\check{\Psi}_{r,t}(z), r)}_{\text{vJp}} dr \\ &\quad + \sum_{i=1}^m \int_s^t \underbrace{\check{A}_{r,t}(z) \nabla \sigma_i(\check{\Psi}_{r,t}(z), r)}_{\text{vJp}} \circ d\check{W}_r^{(i)}.\end{aligned}$$

- Derivative of the expected loss is equal to the expectation of derivative  $A_{0,T}$ , assuming the existence of a dominating R.V.

### Algorithm 1 ODE Adjoint Sensitivity

**Input:** Parameters  $\theta$ , start time  $t_0$ , stop time  $t_1$ , final state  $z_{t_1}$ , loss gradient  $\partial \mathcal{L}/\partial z_{t_1}$ , dynamics  $f(z, t, \theta)$ .

```
def  $\bar{f}([z_t, a_t, \cdot], t, \theta)$ :           ▷ Augmented dynamics  
     $v = f(z_t, -t, \theta)$   
    return  $[-v, a_t \partial v / \partial z, a_t \partial v / \partial \theta]$ 
```

$$\frac{\partial \mathcal{L}}{\partial z} = \alpha(t) \frac{\partial f}{\partial z}$$

```
 $\begin{bmatrix} z_{t_0} \\ \partial \mathcal{L} / \partial z_{t_0} \\ \partial \mathcal{L} / \partial \theta \end{bmatrix} = \text{odeint}\left(\begin{bmatrix} z_{t_1} \\ \partial \mathcal{L} / \partial z_{t_1} \\ 0_p \end{bmatrix}, \bar{f}, -t_1, -t_0\right)$   
return  $\partial \mathcal{L} / \partial z_{t_0}, \partial \mathcal{L} / \partial \theta$ 
```

### Algorithm 2 SDE Adjoint Sensitivity (Ours)

**Input:** Parameters  $\theta$ , start time  $t_0$ , stop time  $t_1$ , final state  $z_{t_1}$ , loss gradient  $\partial \mathcal{L}/\partial z_{t_1}$ , drift  $f(z, t, \theta)$ , diffusion  $\sigma(z, t, \theta)$ , Wiener process sample  $w(t)$ .

```
def  $\bar{f}([z_t, a_t, \cdot], t, \theta)$ :           ▷ Augmented drift  
     $v = f(z_t, -t, \theta)$   
    return  $[-v, a_t \partial v / \partial z, a_t \partial v / \partial \theta]$ 
```

```
def  $\bar{\sigma}([z_t, a_t, \cdot], t, \theta)$ :          ▷ Augmented diffusion  
     $v = \sigma(z_t, -t, \theta)$   
    return  $[-v, a_t \partial v / \partial z, a_t \partial v / \partial \theta]$ 
```

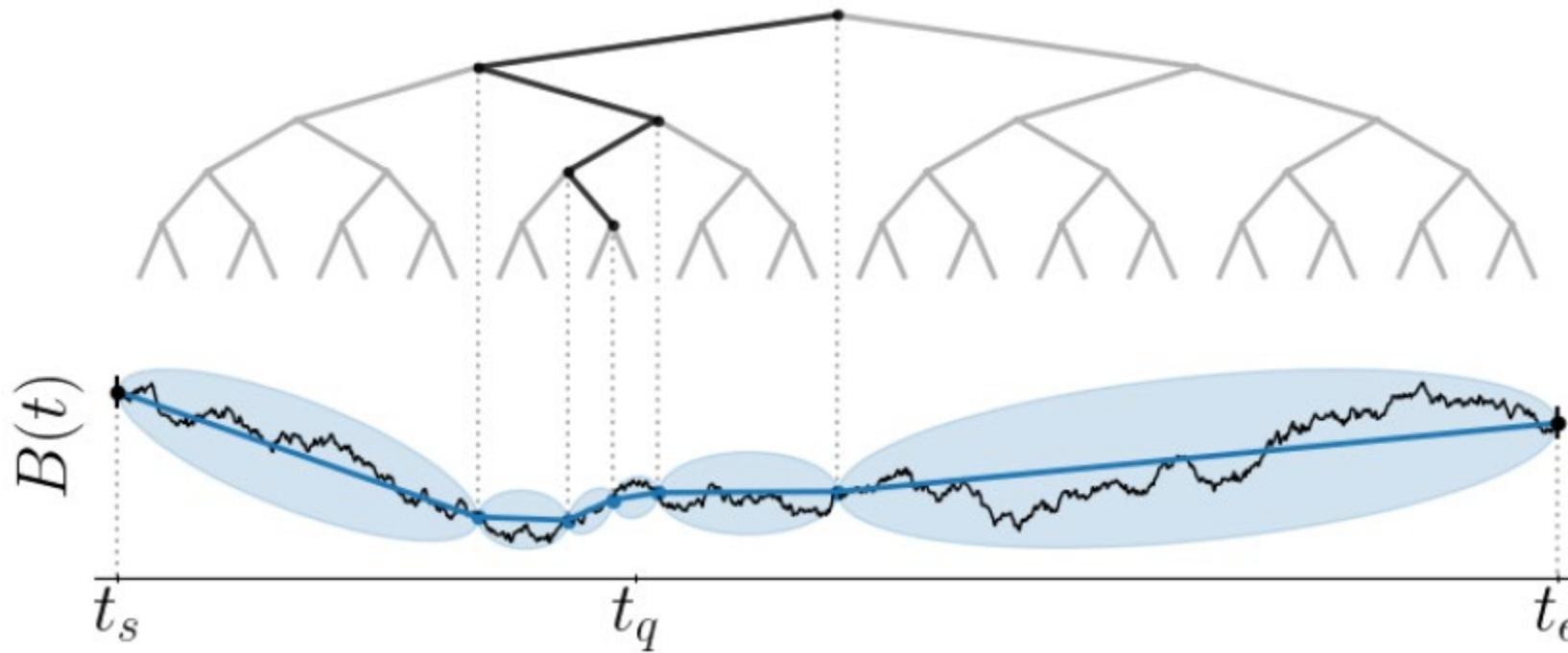
```
def  $\bar{w}(t)$ :                                ▷ Replicated noise  
    return  $[-w(-t), -w(-t), -w(-t)]$ 
```

```
 $\begin{bmatrix} z_{t_0} \\ \partial \mathcal{L} / \partial z_{t_0} \\ \partial \mathcal{L} / \partial \theta \end{bmatrix} = \text{sdeint}\left(\begin{bmatrix} z_{t_1} \\ \partial \mathcal{L} / \partial z_{t_1} \\ 0_p \end{bmatrix}, \bar{f}, \bar{\sigma}, \bar{w}, -t_1, -t_0\right)$   
return  $\partial \mathcal{L} / \partial z_{t_0}, \partial \mathcal{L} / \partial \theta$ 
```

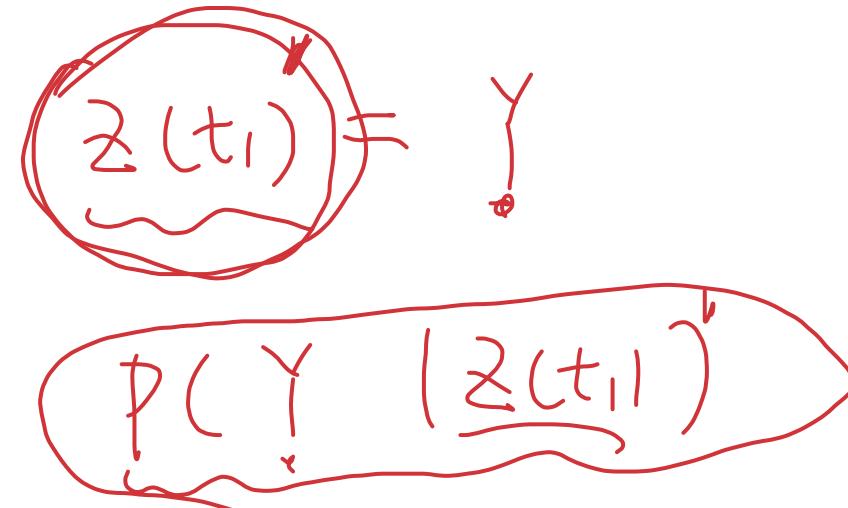
What omit here..

- Some numerical / converge proof..
- Some memory/time saving trick to reproducing noise  
(Brownian tree)

$w(t)$



## Latent SDE



- Fitting SDEs with flexible drift and diffusion coefficients directly to data using maximum likelihood would likely cause a degenerate solution.
- It is reasonable to expect the diffusion to shrink to zero when flexible ODE models can already do a good job.
- One real benefit of using SDEs is when they are treated as latent variables: Place a prior, and infer the posterior given data.

# Latent SDE

$$\begin{aligned} d\tilde{Z}_t &= h_\theta(\tilde{Z}_t, t) dt + \sigma(\tilde{Z}_t, t) dW_t, && \text{(prior)} \\ dZ_t &= h_\phi(Z_t, t) dt + \sigma(Z_t, t) dW_t, && \text{(approx. post.)} \end{aligned}$$

Different function      Same function (no parameter)

$h_\theta$  is circled in red,  $h_\phi$  is circled in blue, and  $\sigma$  is circled in green.

where  $h_\theta, h_\phi$ , and  $\sigma$  are Lipschitz in both arguments, and both processes have the same starting value:  $\tilde{Z}_0 = Z_0 = z_0 \in \mathbb{R}^d$ .

If both processes share the same diffusion function  $\sigma$ , then the KL divergence between them is finite (under additional mild regularity conditions; see Appendix 9.7), and can be estimated by sampling paths from the approximate posterior process. Then, the

evidence lower bound (ELBO) can be written as:

$$\log p(x_1, x_2, \dots, x_N | \theta) \geq \quad (10)$$

$$\mathbb{E}_{Z_t} \left[ \sum_{i=1}^N \log p(x_{t_i} | z_{t_i}) - \int_0^T \frac{1}{2} |u(z_t, t)|^2 dt \right], \checkmark$$

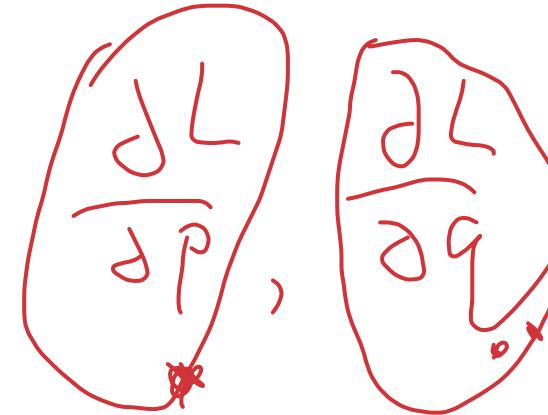
where  $u : \mathbb{R}^d \times [0, T] \rightarrow \mathbb{R}^m$  satisfies

$$\sigma(z, t) u(z, t) = h_\phi(z, t) - h_\theta(z, t);$$

*prior drift*      *post drift.*

and the expectation is taken over the approximate posterior process defined by (approx. post.). The likelihoods of observations  $x_1, \dots, x_N$  at times  $t_1, \dots, t_N$  depend only on latent states  $z_t$  at corresponding times.

# Interpreting the Objective



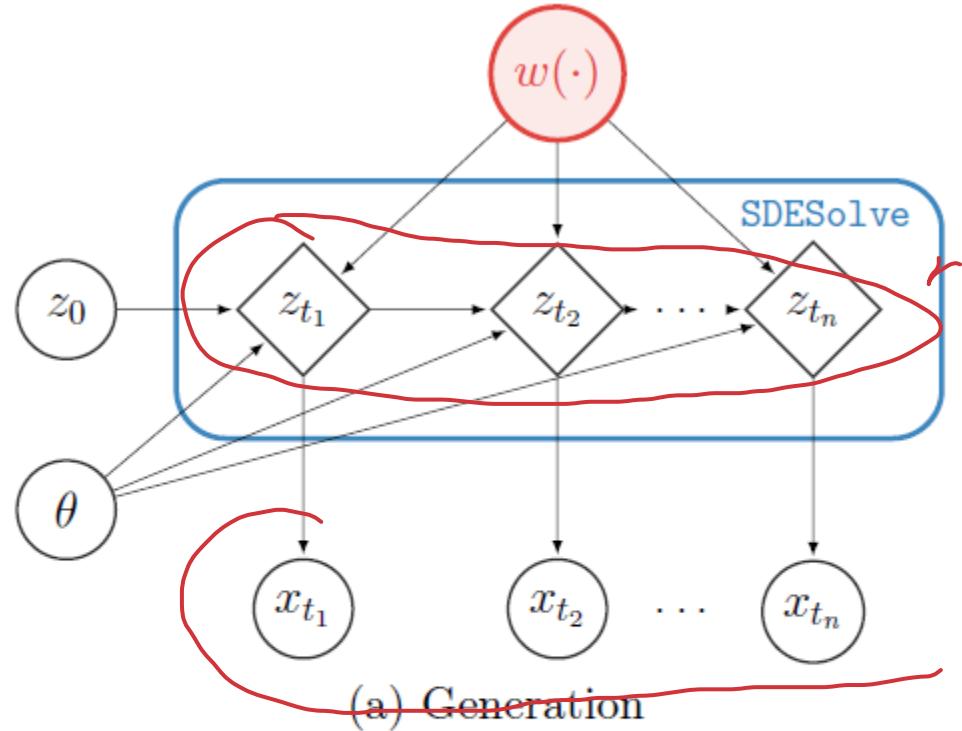
$$\mathcal{L}_{\text{ELBO}} = - \underbrace{\mathbb{E} \left[ \frac{1}{2} \int_0^T |u(q_t, t)|^2 dt \right]}_{\text{KL divergence}} + \underbrace{\mathbb{E} \left[ \sum_{i=1}^N \log p(y_{t_i} | q_{t_i}) \right]}_{\text{negative reconstruction loss}}.$$

As a special case, when the diffusion is a diagonal matrix and non-zero, the loss reduces to

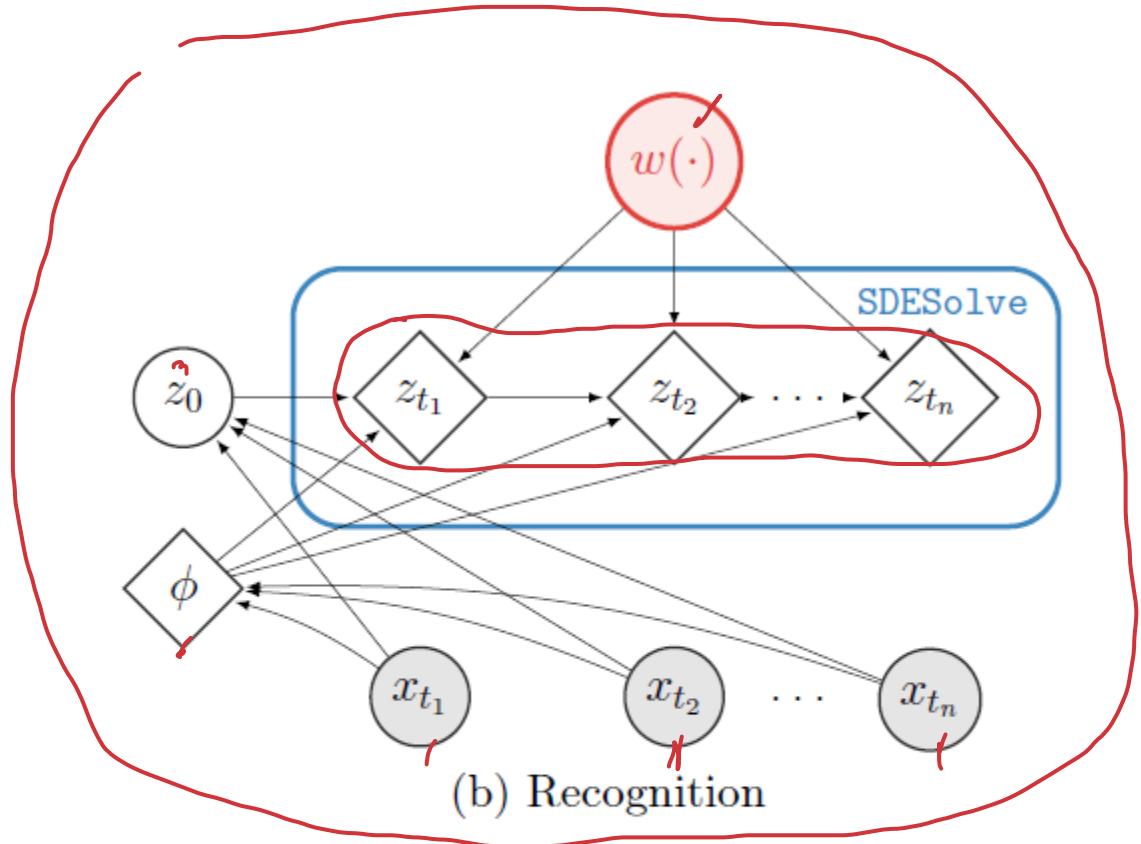
$$\mathcal{L}_{\text{ELBO}} = -\mathbb{E} \left[ \frac{1}{2} \int_0^T |(f_p - f_q)/\sigma|^2 dt \right] + \mathbb{E} \left[ \sum_{i=1}^N \log p(y_{t_i} | q_{t_i}) \right]. \rightarrow$$

↑ prior diff      ↑ post diff

$f_p - f_q$  characterizes the difference;  $\sigma$  scales the penalty.



(a) Generation



(b) Recognition

Figure 4: Graphical models for the generative process (decoder) and recognition network (encoder) of the latent stochastic differential equation model. This model can be viewed as a variational autoencoder with infinite-dimensional noise. Red circles represent entire function draws from Brownian motion. Given the initial state  $z_0$  and a Brownian motion sample path  $w(\cdot)$ , the intermediate states  $z_{t_1}, \dots, z_{t_n}$  are deterministically approximated by a numerical SDE solver.

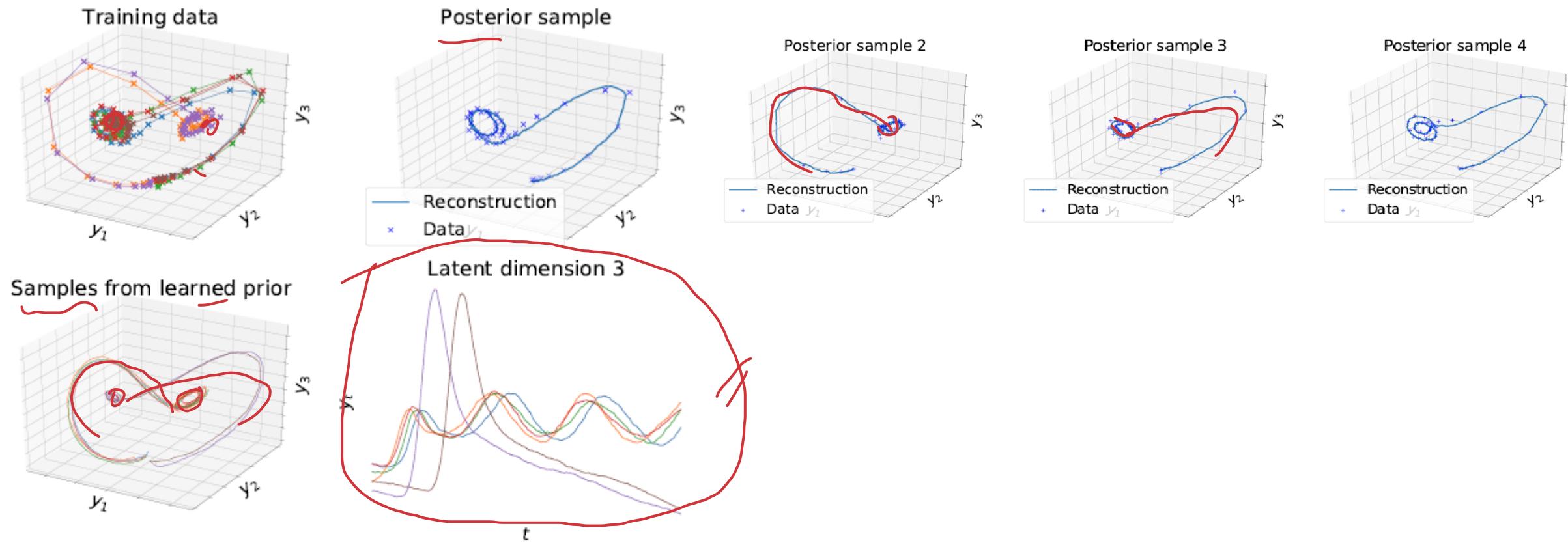


Figure 6: Learned posterior and prior dynamics on data from a stochastic Lorenz attractor. All samples from our model are continuous-time paths, and form a multi-modal, non-Gaussian distribution.

# Applications: Fitting Multiple Trajectories

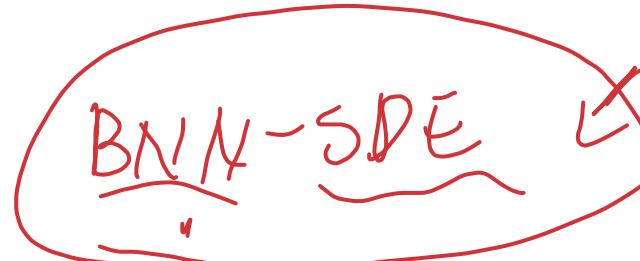
**Table:** Test MSE for a task on a Mocap dataset; 95% confidence interval averaged over 50 samples. †from [10].

Method	Test MSE
npODE [3]	22.96†
NeuralODE [1]	<u>22.49</u> $\pm$ 0.88†
ODE <sup>2</sup> VAE [10]	10.06 $\pm$ 1.4†
ODE <sup>2</sup> VAE-KL [10]	8.09 $\pm$ 1.95†
Latent <u>ODE</u> [9]	5.98 $\pm$ 0.28
Latent <u>SDE</u> (this work)	<b>4.03 <math>\pm</math> 0.20</b>

# Summary

- We presented a generalization of the adjoint sensitivity method to SDEs.
- The procedure relies on vJps, and its memory scaling can interpolate between  $\mathcal{O}(1)$  and  $\mathcal{O}(\text{steps})$  (depending on how much compute time we are willing to trade in).
- We applied adjoints to train latent SDEs.
- Theoretical results also studied a new type of convergence for SDE schemes, i.e. local uniform convergence, for which we focused on the Euler-Maruyama scheme.

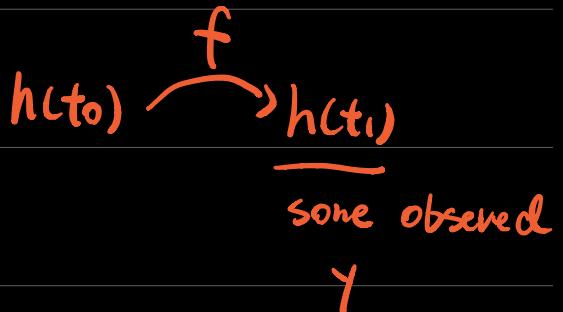
# Next Steps



- Modeling:
  - Comprehensive empirical comparison: GPs vs SDEs.
  - Stationary and multi-timescale models.
  - Bayesian neural networks.
- Numerical analysis:
  - Convert the method to solve high-dimensional PDEs.
  - New theoretical interpretations using rough path theory.
- Other areas of science and engineering:
  - Better models of allele frequency in population genetics.
  - Fast and memory-frugal algorithms for Greeks in finance.
  - Neural science applications?
- Software: A PyTorch package is on its way! Stay tuned!

Neural ODE:

$$\frac{dh(t)}{dt} = f(h, \theta, t) \\ \stackrel{=}{=} \text{NN}$$



↑ if NN given

$$h(t) = h(t_0) + \underbrace{\int_{t_0}^t f(h(t'), \theta(t'), t') dt'}_{\text{solved by numerical method.}} \quad \textcircled{0}$$

Given some loss function  $L$ , target value.

Goal:  $\frac{\partial L}{\partial \theta}$ .

Define:  $\alpha(t) = \frac{\partial L}{\partial h(t)}$  (adjoint state)

after some derivation, we get

$$\frac{d\alpha(t)}{dt} = -\alpha(t) \frac{\partial f(h(t), t, \theta)}{\partial h(t)} \quad \textcircled{1}$$

analytic form / auto-diff.

a new ode system, solve by numerical.

to get  $\alpha(t)$  [ bounty :  $\alpha(t_i) = \frac{\partial L}{\partial h(t_i)}$

We can get  $\alpha(t_0) = \frac{\partial L}{\partial h(t_0)}$  by

inverse inte :

$$\alpha(t_0) = \alpha(t_i) + \int_{t_i}^{t_0} \frac{d\alpha(t)}{dt} dt.$$

Similar, we can define adjoint state

for  $\theta, t$ :

$$\left\{ \begin{array}{l} \alpha_\theta(t) = \frac{\partial L}{\partial \dot{\theta}} \\ \alpha_t = \frac{\partial L}{\partial t} \end{array} \right. \quad ②$$

with some derivation, get:

$$\left\{ \begin{array}{l} \frac{\alpha_\theta(t)}{dt} = -\underbrace{\alpha(t)}_{\text{act}} \frac{\partial f}{\partial \theta} \\ \frac{\alpha_t(t)}{dt} = -\underbrace{\alpha(t)}_{\text{act}} \frac{\partial f}{\partial t}. \end{array} \right. \quad ③$$

when we get  $\alpha(t)$  from ①,

we can use the same numeral solver to

get  $\alpha_\theta(t)$ ,  $\alpha_t(t)$

The general step: (for each  $t$ )

a. solve ② to get  $h(t)$

b. solve ① to get  $\alpha(t)$

c. solve ③ to get  $\frac{\alpha_\theta(t)}{\frac{\partial L}{\partial \theta}}, \frac{\alpha_t(t)}{\frac{\partial L}{\partial t}}$