

A IDENTIFIABILITY OF STOCHASTIC DIFFERENTIAL EQUATIONS

A.1 DEFINITIONS AND ASSUMPTIONS

In this part, we will introduce some basic definitions and assumptions required for the theory.

First, let us restate assumption 1.

Assumption 1 (Global Lipschitz). *We assume that the drift and diffusion functions in eq. (5) satisfy the global Lipschitz constraints. Namely, we have*

$$|\mathbf{f}_\theta(\mathbf{x}) - \mathbf{f}_\theta(\mathbf{y})| + |\mathbf{g}_\theta(\mathbf{x}) - \mathbf{g}_\theta(\mathbf{y})| \leq C|\mathbf{x} - \mathbf{y}| \quad (6)$$

for some constant C , $\mathbf{x}, \mathbf{y} \in \mathbb{R}^D$ and $|\cdot|$ is the corresponding L_2 norm for vector-valued functions and matrix norm for matrix-valued functions.

This assumption regularizes the Itô diffusion to have a unique strong solution \mathbf{X}_t to eq. (3), which is a standard assumption in the SDE literature. In addition, this diffusion satisfies the Feller continuous property, and its solution is a Feller process (Lemma 8.1.4 in Øksendal & Øksendal (2003)).

Definition 1 (Feller process and semi-group). *A continuous time-homogeneous Markov family \mathbf{X}_t is a Feller process when, for all $\mathbf{x} \in \mathbb{R}^D$, we have $\forall t, \mathbf{y} \rightarrow \mathbf{x} \Rightarrow \mathbf{X}_{y,t} \xrightarrow{d} \mathbf{X}_{x,t}$ and $t \rightarrow 0 \Rightarrow \mathbf{X}_{x,t} \xrightarrow{p} \mathbf{x}$ where $\xrightarrow{d}, \xrightarrow{p}$ means convergence in distribution and in probability, respectively, and $\mathbf{X}_{y,t}$ means the solution with y as the initial condition. A semigroup of linear, positive, conservative contraction operators \mathbb{T}_t is a Feller semigroup if, for every $\mathbf{f} \in C_0, \mathbf{x} \in \mathbb{R}^D$, we have $\mathbb{T}_t \mathbf{f} \in C_0$ and $\lim_{t \rightarrow 0} \mathbb{T}_t \mathbf{f}(\mathbf{x}) = \mathbf{f}(\mathbf{x})$, where C_0 is the space of continuous functions vanishing at infinity.*

Basically, the transition operator of a Feller process is a Feller semigroup. The reason we care about the Feller process is its nice properties related to its infinitesimal generators. In a nutshell, the distributional properties of the Feller process can be uniquely characterised by its generators.

Definition 2 (Infinitesimal generator). *For a Feller process \mathbf{X}_t with a Feller semigroup \mathbb{T}_t , we define the generator A by*

$$A f = \lim_{t \downarrow 0} \frac{\mathbb{T}_t f - f}{t} \quad \text{for any } f \in D(A) \quad (17)$$

where $D(A)$ is the domain of the generator, defined as the function space where the above limit exists.

Next, let us restate assumption 2.

Assumption 2 (Diagonal diffusion). *We assume that the diffusion function \mathbf{g}_θ outputs a non-zero diagonal matrix. That is, it can be simplified to a vector-valued function $\mathbf{g}_\theta(\mathbf{X}_t) : \mathbb{R}^D \rightarrow \mathbb{R}^D$.*

This is a key assumption for structure identifiability. For a general matrix diffusion function, it is easy to come up with unidentifiable examples (see Example 5.5 in (Hansen & Sokol, 2014)). For example, in a driftless process, the distribution of \mathbf{X}_t will depend on $\mathbf{g}_G \mathbf{g}_G^T$, where it can have multiple factorizations that correspond to different graphs.

A.2 STRUCTURE IDENTIFIABILITY FOR OBSERVATIONAL PROCESS

Now, let us re-state theorem 4.1:

Theorem 4.1 (Structure identifiability of the observational process). *Given eq. (16), let us define another process with $\bar{\mathbf{X}}_t, \bar{\mathbf{G}} \neq \mathbf{G}, \bar{\mathbf{f}}_{\bar{\mathbf{G}}}, \bar{\mathbf{g}}_{\bar{\mathbf{G}}}$ and $\bar{\mathbf{W}}_t$. Then, under Assumptions 1-2, and with the same initial condition $\mathbf{X}(0) = \bar{\mathbf{X}}(0) = \mathbf{x}_0$, the solutions \mathbf{X}_t and $\bar{\mathbf{X}}_t$ will have different distributions.*

To prove this theorem, we begin by establishing the corresponding result for the discretized Euler SEMs, and then build the connection to the Itô diffusion through the infinitesimal generator.

Lemma A.1 (Identifiability of Euler SEM). *Assuming assumption 2 is satisfied with nonzero diagonal diffusion functions. For a Euler SEM defined as*

$$\mathbf{X}_{t+1}^\Delta = \mathbf{X}_t^\Delta + \mathbf{f}_G(\mathbf{X}_t^\Delta)\Delta + \mathbf{g}_G(\mathbf{X}_t^\Delta)\boldsymbol{\eta}_t, \quad \boldsymbol{\eta}_t \sim \mathcal{N}(0, \Delta\mathbf{I}), \quad (18)$$

if we have another Euler SEM defined as

$$\bar{\mathbf{X}}_{t+1}^\Delta = \bar{\mathbf{X}}_t^\Delta + \bar{\mathbf{f}}_{\bar{\mathbf{G}}}(\bar{\mathbf{X}}_t^\Delta)\Delta + \bar{\mathbf{g}}_{\bar{\mathbf{G}}}(\bar{\mathbf{X}}_t^\Delta)\bar{\boldsymbol{\eta}}_t, \quad \bar{\boldsymbol{\eta}}_t \sim \mathcal{N}(0, \Delta\mathbf{I}). \quad (19)$$

Then their corresponding transition density $p(\mathbf{X}_{t+1}^\Delta | \mathbf{X}_t^\Delta = \mathbf{a}) = \bar{p}(\bar{\mathbf{X}}_{t+1}^\Delta | \bar{\mathbf{X}}_t^\Delta = \mathbf{a})$ for all $\mathbf{a} \in \mathbb{R}^D$ iff. $\mathbf{G} = \bar{\mathbf{G}}$, $\mathbf{f}_G = \bar{\mathbf{f}}_{\bar{\mathbf{G}}}$ and $|\mathbf{g}_G| = |\bar{\mathbf{g}}_{\bar{\mathbf{G}}}|$.

Proof. If we have $\mathbf{G} = \bar{\mathbf{G}}$, $\mathbf{f} = \bar{\mathbf{f}}$ and $|\mathbf{g}| = |\bar{\mathbf{g}}|$, then it is trivial that their transition densities are the same since they define the same Euler SEM update equations (up to the sign of the diffusion term) with given initial conditions.

On the other hand, we know

$$\begin{aligned} p(\mathbf{X}_{t+1}^\Delta | \mathbf{X}_t^\Delta = \mathbf{a}) &= \mathcal{N}(\mathbf{f}_G(\mathbf{a})\Delta + \mathbf{a}, \mathbf{g}_G^2(\mathbf{a})\Delta) \\ \bar{p}(\bar{\mathbf{X}}_{t+1}^\Delta | \bar{\mathbf{X}}_t^\Delta = \mathbf{a}) &= \mathcal{N}(\bar{\mathbf{f}}_{\bar{\mathbf{G}}}(\mathbf{a})\Delta + \mathbf{a}, \bar{\mathbf{g}}_{\bar{\mathbf{G}}}^2(\mathbf{a})\Delta) \end{aligned}$$

Thus, if two conditional distributions match, we have

$$\mathbf{f}_G(\mathbf{a})\Delta = \bar{\mathbf{f}}_{\bar{\mathbf{G}}}(\mathbf{a})\Delta \quad \mathbf{g}_G^2(\mathbf{a})\Delta = \bar{\mathbf{g}}_{\bar{\mathbf{G}}}^2(\mathbf{a})\Delta \quad (20)$$

Since $\Delta > 0$, we have $\mathbf{f}_G(\mathbf{a}) = \bar{\mathbf{f}}_{\bar{\mathbf{G}}}(\mathbf{a})$, $\mathbf{g}_G^2(\mathbf{a}) = \bar{\mathbf{g}}_{\bar{\mathbf{G}}}^2(\mathbf{a})$ for all $\mathbf{a} \in \mathbb{R}^D$. From the diagonal diffusion assumption, we know $|\mathbf{g}_G(\mathbf{a})| = |\bar{\mathbf{g}}_{\bar{\mathbf{G}}}(\mathbf{a})|$.

Now, assume for contradiction that $\mathbf{G} \neq \bar{\mathbf{G}}$; then there exists $X_{t,i}^\Delta \rightarrow X_{t+1,j}^\Delta$ in \mathbf{G} but not in $\bar{\mathbf{G}}$. Then we have by definition that $\frac{\partial \bar{f}_j(\mathbf{X}_t^\Delta, \bar{\mathbf{G}})}{\partial X_{t,i}^\Delta} = 0$ and $\frac{\partial \bar{g}_j(\mathbf{X}_t^\Delta, \bar{\mathbf{G}})}{\partial X_{t,i}^\Delta} = 0$ for all \mathbf{X}_t^Δ , and also $\frac{\partial f_j(\mathbf{X}_t^\Delta, \mathbf{G})}{\partial X_{t,i}^\Delta} \neq 0$ or $\frac{\partial g_j(\mathbf{X}_t^\Delta, \mathbf{G})}{\partial X_{t,i}^\Delta} \neq 0$ for some \mathbf{X}_t^Δ . In the former case, if $\frac{\partial f_j(\mathbf{X}_t^\Delta, \mathbf{G})}{\partial X_{t,i}^\Delta} \neq 0$ but $\frac{\partial \bar{f}_j(\mathbf{X}_t^\Delta, \bar{\mathbf{G}})}{\partial X_{t,i}^\Delta} = 0$ for some \mathbf{X}_t^Δ , we have a contradiction to $\mathbf{f}_G(\mathbf{a}) = \bar{\mathbf{f}}_{\bar{\mathbf{G}}}(\mathbf{a})$ for $\mathbf{a} \in \mathbb{R}^D$. A similar analysis can be done in the latter case for \mathbf{g}_G , $\bar{\mathbf{g}}_{\bar{\mathbf{G}}}$. Thus, we have $\mathbf{G} = \bar{\mathbf{G}}$, $\mathbf{f}_G = \bar{\mathbf{f}}_{\bar{\mathbf{G}}}$ and $|\mathbf{g}_G| = |\bar{\mathbf{g}}_{\bar{\mathbf{G}}}|$. \square

Next, we will prove a lemma that builds a bridge between the generator of the Itô diffusion and its corresponding Euler SEM.

Lemma A.2 (Generator characterises Euler SEM). *Assume that assumptions 1 and 2. For an Itô diffusion defined as eq. (16), we denote its corresponding variables in Euler SEM with Δ discretization as \mathbf{X}^Δ . Similarly, if we have an alternative Itô diffusion defined with $\bar{\mathbf{f}}_{\bar{\mathbf{G}}}$, $\bar{\mathbf{g}}_{\bar{\mathbf{G}}}$ and $\bar{\mathbf{G}}$, we define the corresponding Euler SEM variables $\bar{\mathbf{X}}^\Delta$. Then, the generators of the Itô diffusions $\mathbf{A} = \bar{\mathbf{A}}$ iff. their Euler SEM variables have the same distribution with given initial conditions.*

Proof. First, assume $\mathbf{A} = \bar{\mathbf{A}}$, then for any $h \in C_0^2$ (twice continuously differentiable functions vanishing at infinity), we can define the generator for Itô diffusion as

$$\mathbf{A} h(\mathbf{x}) = \sum_d f_d(\mathbf{x}, \mathbf{G}) \frac{\partial h(\mathbf{x})}{\partial x_d} + \frac{1}{2} \sum_d g_d^2(\mathbf{x}, \mathbf{G}) \frac{\partial^2 h(\mathbf{x})}{\partial x_d^2} \quad (21)$$

Similarly, we can define $\bar{\mathbf{A}}$. From Lemma A.3 (Hansen & Sokol, 2014), we know if $\mathbf{A} = \bar{\mathbf{A}}$, then $\mathbf{G} = \bar{\mathbf{G}}$, $\mathbf{f}(\cdot, \mathbf{G}) = \bar{\mathbf{f}}(\cdot, \bar{\mathbf{G}})$ and $\mathbf{g}^2(\cdot, \mathbf{G}) = \bar{\mathbf{g}}^2(\cdot, \bar{\mathbf{G}})$ for $\mathbf{x} \in \mathbb{R}^D$. Therefore, by the definition of Euler SEM (eq. (4)), it is trivial that they define the same transition density $p(\mathbf{X}_{t+1}^\Delta | \mathbf{X}_t^\Delta = \mathbf{a}) = \bar{p}(\bar{\mathbf{X}}_{t+1}^\Delta | \bar{\mathbf{X}}_t^\Delta = \mathbf{a})$ for $\mathbf{a} \in \mathbb{R}^D$.

On the other hand, if the two Euler SEMs define the same transition densities, then from Lemma A.1, we have $\mathbf{f}_G = \bar{\mathbf{f}}_{\bar{\mathbf{G}}}$, $|\mathbf{g}_G| = |\bar{\mathbf{g}}_{\bar{\mathbf{G}}}|$ and $\mathbf{G} = \bar{\mathbf{G}}$. Then from eq. (21), $\mathbf{A} = \bar{\mathbf{A}}$. \square

Finally, the following lemma shows why we care about the infinitesimal generator for the Feller process.

Lemma A.3 (Generator uniquely determines Feller semigroup). *Let us define the Feller semigroup transition operator \mathbb{T}_t and $\bar{\mathbb{T}}_t$ associated with generator \mathbf{A} , $\bar{\mathbf{A}}$. Then, $\mathbb{T}_t = \bar{\mathbb{T}}_t$ iff. $\mathbf{A} = \bar{\mathbf{A}}$.*

Proof. We define the resolvent of a Feller process with $\lambda > 0$ as:

$$R_\lambda f = \int_0^\infty \exp(-\lambda t) T_t f dt \quad (22)$$

with $f \in C_0$. This is the Laplace transform of $T_t f$. From Øksendal & Øksendal (2003), we know $R_\lambda = (\lambda \mathbf{I} - \mathbf{A})^{-1}$. Therefore, if $\mathbf{A} = \bar{\mathbf{A}}$, then for $\lambda > 0$, the resolvent $R_\lambda = (\lambda \mathbf{I} - \mathbf{A})^{-1} = (\lambda \mathbf{I} - \bar{\mathbf{A}})^{-1} = \bar{R}_\lambda$. Therefore, for all $h \in C_0$, they define the same Laplace transform of $T_t h$. From the uniqueness of Laplace transform, we have $\bar{T}_t = \bar{T}_t$.

Similarly, if $T_t = \bar{T}_t$, we have $R_\lambda = \bar{R}_\lambda$ from the definition of resolvent. Thus, $\mathbf{A} = \lambda \mathbf{I} - R_\lambda^{-1} = \lambda \mathbf{I} - \bar{R}_\lambda^{-1} = \bar{\mathbf{A}}$. \square

Now, we can prove theorem 4.1.

Proof. Suppose we have two different observation process defined with $\mathbf{G} \neq \bar{\mathbf{G}}$. Then, from Lemma A.1, with any $\Delta > 0$, their Euler transition distribution $\bar{p}(\bar{\mathbf{X}}_{t+1}^\Delta | \bar{\mathbf{X}}_t^\Delta = \mathbf{a}) \neq p(\mathbf{X}_{t+1}^\Delta | \mathbf{X}_t^\Delta = \mathbf{a})$. Thus, from Lemma A.2, these two Itô diffusions have different generators $\mathbf{A} \neq \bar{\mathbf{A}}$. From assumption 1, the solutions of these two Itô diffusions are Feller processes. From Lemma A.3, if $\mathbf{A} \neq \bar{\mathbf{A}}$, their semigroup $T_t \neq \bar{T}_t$, which results in different observation distributions of $\mathbf{X}_t, \bar{\mathbf{X}}_t$. \square

A.3 IDENTIFIABILITY OF LATENT SDE

We begin by re-stating theorem 4.2:

Theorem 4.2 (Structural identifiability with latent formulation). *Consider the distributions p, \bar{p} defined by the latent model in eq. (5) with $(\mathbf{G}, \mathbf{Z}, \mathbf{X}, \mathbf{f}_G, \mathbf{g}_G), (\bar{\mathbf{G}}, \bar{\mathbf{Z}}, \bar{\mathbf{X}}, \bar{\mathbf{f}}_G, \bar{\mathbf{g}}_G)$ respectively, where $\mathbf{G} \neq \bar{\mathbf{G}}$. Further, let t_1, \dots, t_I be the observation times. Then, under Assumptions 1 and 2:*

1. *if $t_{i+1} - t_i = \Delta$ for all $i \in 1, \dots, I - 1$, then $p^\Delta(\mathbf{X}_{t_1}, \dots, \mathbf{X}_{t_I}) \neq \bar{p}^\Delta(\bar{\mathbf{X}}_{t_1}, \dots, \bar{\mathbf{X}}_{t_I})$, where p^Δ is the density generated by the Euler discretized eq. (9) for \mathbf{Z}_t ;*
2. *if we have a fixed time range $[0, T]$, then the path probability $p(\mathbf{X}_{t_1}, \dots, \mathbf{X}_{t_I}) \neq \bar{p}(\bar{\mathbf{X}}_{t_1}, \dots, \bar{\mathbf{X}}_{t_I})$ under the limit of infinite data ($I \rightarrow \infty$).*

We follow the same proof strategy as Hasan et al. (2021); Khemakhem et al. (2020).

Proof. Let's assume $p(\mathbf{X}_{t_1}, \dots, \mathbf{X}_{t_I}) = \bar{p}(\bar{\mathbf{X}}_{t_1}, \dots, \bar{\mathbf{X}}_{t_I})$ even though $\mathbf{G} \neq \bar{\mathbf{G}}$. Then, for any t_{i+1} and t_i , we have $p(\mathbf{X}_{t_{i+1}}, \mathbf{X}_{t_i}) = \bar{p}(\bar{\mathbf{X}}_{t_{i+1}}, \bar{\mathbf{X}}_{t_i})$. Then, we can write

$$\begin{aligned} p(\mathbf{X}_{t_{i+1}}, \mathbf{X}_{t_i}) &= \int p(\mathbf{Z}_{t_{i+1}}, \mathbf{Z}_{t_i}, \mathbf{X}_{t_{i+1}}, \mathbf{X}_{t_i}) d\mathbf{Z}_{t_{i+1}} d\mathbf{Z}_{t_i} \\ &= \int p_z(\mathbf{Z}_{t_{i+1}}, \mathbf{Z}_{t_i}) p_\epsilon(\mathbf{X}_{t_{i+1}} - \mathbf{Z}_{t_{i+1}}) p_\epsilon(\mathbf{X}_{t_i} - \mathbf{Z}_{t_i}) d\mathbf{Z}_{t_{i+1}} d\mathbf{Z}_{t_i} \\ &= [(p_\epsilon \times p_\epsilon) * p_z](\mathbf{X}_{t_{i+1}}, \mathbf{X}_{t_i}) \end{aligned}$$

where p_ϵ is the noise density for the added observational noise ϵ , p_z is the joint density defined by latent Itô diffusion and $*$ is the convolution operator. Thus, by applying the Fourier transform \mathcal{F} , we obtain

$$\mathcal{F}(p_\epsilon \times p_\epsilon)(\omega) \times \mathcal{F}(p_z)(\omega) = \mathcal{F}(p_\epsilon \times p_\epsilon)(\omega) \times \mathcal{F}(\bar{p}_z)(\omega) \quad (23)$$

So $\mathcal{F}(p_z) = \mathcal{F}(\bar{p}_z)$. Then, by inverse Fourier transform, we have $p_z(\mathbf{Z}_{t_{i+1}}, \mathbf{Z}_{t_i}) = \bar{p}_z(\bar{\mathbf{Z}}_{t_{i+1}}, \bar{\mathbf{Z}}_{t_i})$.

If the above distributions are obtained by discretizing the Itô diffusion with a fixed step size Δ , they become the corresponding discretized distribution $p^\Delta(\mathbf{Z}_{t_{i+1}}^\Delta, \mathbf{Z}_{t_i}^\Delta)$ (i.e. defined by Euler SEM). Then the transition density $p^\Delta(\mathbf{Z}_{t_{i+1}}^\Delta | \mathbf{Z}_{t_i}^\Delta) = \bar{p}^\Delta(\bar{\mathbf{Z}}_{t_{i+1}}^\Delta | \bar{\mathbf{Z}}_{t_i}^\Delta)$. From Lemma A.1, we have $\mathbf{G} = \bar{\mathbf{G}}$, resulting in a contradiction. Thus, $p^\Delta(\mathbf{X}_{t_1}, \dots, \mathbf{X}_{t_I}) \neq \bar{p}^\Delta(\bar{\mathbf{X}}_{t_1}, \dots, \bar{\mathbf{X}}_{t_I})$.

If we have a fixed time range $[0, T]$, then, when we have infinite observations $I \rightarrow \infty$, the observation time t follows an independent temporal point process with intensity $\lim_{dt \rightarrow 0} Pr(\text{observe in } [t, t + dt] | \mathcal{H}_t) > 0$ where \mathcal{H}_t is the filtration. Thus, for arbitrary time interval $\Delta > 0$, we have $p(\mathbf{Z}_{t+\Delta}, \mathbf{Z}_t) = \bar{p}(\bar{\mathbf{Z}}_{t+\Delta}, \bar{\mathbf{Z}}_t)$. Since this holds for arbitrarily small $\Delta > 0$,

this equality in densities means they define the same transition density $p(\mathbf{Z}_{t+\Delta}|\mathbf{Z}_t) = \bar{p}(\bar{\mathbf{Z}}_{t+\Delta}|\bar{\mathbf{Z}}_t)$ as $\Delta \rightarrow 0$. By definition of the Feller transition semigroup, we have $\mathbb{T}_t = \bar{\mathbb{T}}_t$. From Lemma A.3, $\mathbf{A} = \bar{\mathbf{A}}$ and $\mathbf{G} = \bar{\mathbf{G}}$ (Lemma A.2, A.1). This leads to contradiction, meaning that $p(\mathbf{X}_{t_1}, \dots, \mathbf{X}_{t_I}) \neq \bar{p}(\mathbf{X}_{t_1}, \dots, \mathbf{X}_{t_I})$ when $I \rightarrow \infty$. \square

A.4 RECOVERY OF THE GROUND TRUTH GRAPH

Before diving into the proof of theorem 4.3, we introduce some necessary assumptions:

Assumption 3 (Correctly specified model). *We say a model is correctly specified w.r.t. the ground truth data generating mechanism iff. there exists a model parameter such that the model coincides with the generating mechanism.*

Assumption 4 (Expressive posterior process). *For a given prior parameter θ , we say the approximate posterior process (eq. (12)) is expressive enough if there exists a measurable function $\mathbf{u}(\mathbf{Z}_t)$ such that (i) $\mathbf{g}_G(\mathbf{Z}_t)\mathbf{u}(\mathbf{Z}_t) = \mathbf{f}_G(\mathbf{Z}_t) - \mathbf{h}_\phi(\mathbf{Z}, t, \mathbf{G})$; (ii) $\mathbf{u}(\mathbf{Z}_t)$ satisfies Novikov's condition and (iii) we define*

$$\mathbf{M}_T = \exp\left(-\frac{1}{2}\int_0^T |\mathbf{u}(\mathbf{Z}_t)|^2 dt - \int_0^T \mathbf{u}(\mathbf{Z}_t)^T d\mathbf{W}_t\right) \quad (24)$$

and for given latent states $\mathbf{Z}_{t_1}, \dots, \mathbf{Z}_{t_I}$ and corresponding observations $\mathbf{X}_{t_1}, \dots, \mathbf{X}_{t_I}$ with $0 \leq t_1 \leq t_2 \leq \dots \leq t_I \leq T$, \mathbf{M}_T can approximate the following arbitrarily well:

$$\mathbf{M}_T \approx \frac{\prod_{i=1}^I p(\mathbf{X}_{t_i}|\mathbf{Z}_{t_i}, \mathbf{G})}{p(\mathbf{X}_{t_1}, \dots, \mathbf{X}_{t_I}|\mathbf{G})} \quad (25)$$

This assumption is to make sure the approximate posterior process is expressive enough to make the variational bound tight. Since we use neural networks to define the drift and diffusion functions, the corresponding approximate posterior is flexible. In fact, Tzen & Raginsky (2019b) showed that the diffusion defined by eq. (12) can be used to obtain samples from any distributions whose Radon-Nikodym derivative w.r.t. standard Gaussian measure can be represented by neural networks. Due to the universal approximation theorem for neural networks (Hornik et al., 1989), the corresponding posterior is indeed flexible.

First, we can re-write the ELBO (eq. (15)) (for a single time series) as the following:

$$\log p(\mathbf{X}_{t_1}, \dots, \mathbf{X}_{t_I}) \geq \mathbb{E}_{\mathbf{G} \sim q_\phi(\mathbf{G})} \left[\mathbb{E}_P \left[\sum_{i=1}^I \log p(\mathbf{X}_{t_i}|\tilde{\mathbf{Z}}_{t_i}) - \frac{1}{2} \int_0^T |\mathbf{u}(\tilde{\mathbf{Z}}_t)|^2 dt \right] \right] - D_{\text{KL}}[q_\phi(\mathbf{G})\|p(\mathbf{G})] \quad (26)$$

where P is the probability measure in the filtered probability space $(\Sigma, \mathcal{F}, \{\mathcal{F}\}_{0 \leq t \leq T}, P)$, and $\tilde{\mathbf{Z}}_t$ is the path sampled from the approximate posterior process (eq. (12)). Let's restate the theorem:

Theorem 4.3 (Consistency of variational formulation). *Suppose Assumptions 1-4 are satisfied for the latent formulation (eq. (5)). Then, for a fixed observation time range $[0, T]$, as the number of observations $I \rightarrow \infty$, when ELBO (eq. (15)) is maximised, $q_\phi(\mathbf{G}) = \delta(\mathbf{G}^*)$, where \mathbf{G}^* is the ground truth graph, and the latent formulation recovers the underlying ground truth mechanism.*

Proof. First, we want to show that the term inside the $\mathbb{E}_{\mathbf{G} \sim q_\phi(\mathbf{G})}[\cdot]$ represents the $\log p(\mathbf{X}_{t_1}, \dots, \mathbf{X}_{t_I}|\mathbf{G})$.

We define a measurable function $\mathbf{u}(\mathbf{Z}_t)$ that satisfies Novikov's condition. From the Girsanov theorem, we can construct another process

$$d\tilde{\mathbf{W}} = \mathbf{u}(\mathbf{Z}_t)dt + d\mathbf{W}_t \quad (27)$$

and another probability measure Q s.t. $\tilde{\mathbf{W}}$ is a Brownian motion under measure Q with

$$\frac{dQ}{dP} = \exp\left(-\frac{1}{2}\int_0^T |\mathbf{u}(\mathbf{Z}_t)|^2 dt - \int_0^T \mathbf{u}(\mathbf{Z}_t)^T d\mathbf{W}_t\right) \quad (28)$$

where P is the probability measure associated with the original Brownian motion \mathbf{W}_t . From Boué & Dupuis (1998); Tzen & Raginsky (2019a), we have the following variational formulation:

$$\log \mathbb{E}_P \left[\prod_{i=1}^I p(\mathbf{X}_{t_i} | \mathbf{Z}_{t_i}, \mathbf{G}) \right] = \sup_{Q \in \mathbb{P}} \left\{ -D_{\text{KL}}[Q \| P] + \mathbb{E}_Q \left[\sum_{i=1}^I \log p(\mathbf{X}_{t_i} | \mathbf{Z}_{t_i}, \mathbf{G}) \right] \right\} \quad (29)$$

where \mathbb{P} represents the set of probability measures for the path \mathbf{Z}_t . Assume measure Q is constructed by \mathbf{u} , we can write down $D_{\text{KL}}[Q \| P]$ by substituting eq. (28):

$$\begin{aligned} D_{\text{KL}}[Q \| P] &= \mathbb{E}_Q \left[\log \frac{dQ}{dP} \right] \\ &= \int \left[-\frac{1}{2} \int_0^T |\mathbf{u}(\mathbf{Z}_t)|^2 dt - \int_0^T \mathbf{u}(\mathbf{Z}_t)^T d\mathbf{W}_t \right] dQ \\ &= \int \left[-\frac{1}{2} \int_0^T |\mathbf{u}(\mathbf{Z}_t)|^2 dt + \int_0^T |\mathbf{u}(\mathbf{Z}_t)|^2 dt \right] dQ \\ &= \mathbb{E}_Q \left[\frac{1}{2} \int_0^T |\mathbf{u}(\mathbf{Z}_t)|^2 dt \right]. \end{aligned}$$

The third equality can be obtained by manipulating eq. (27):

$$\begin{aligned} \mathbf{u}(\mathbf{Z}_t)^T d\tilde{\mathbf{W}}_t &= |\mathbf{u}(\mathbf{Z}_t)|^2 dt + \mathbf{u}(\mathbf{Z}_t)^T d\mathbf{W}_t \\ \Rightarrow \underbrace{\mathbb{E}_Q \left[\int_0^T \mathbf{u}^T d\tilde{\mathbf{W}}_t \right]}_{=0} &= \mathbb{E}_Q \left[\int_0^T |\mathbf{u}|^2 dt \right] + \mathbb{E}_Q \left[\int_0^T \mathbf{u}^T d\mathbf{W}_t \right] \end{aligned}$$

The highlighted term is 0 due to the martingale property under measure Q . Thus, we have

$$\mathbb{E}_Q \left[\int_0^T \mathbf{u}^T d\mathbf{W}_t \right] = -\mathbb{E}_Q \left[\int_0^T |\mathbf{u}|^2 dt \right] \quad (30)$$

Now, let's define

$$\mathbf{u}(\mathbf{Z}_t) = \mathbf{g}_G(\mathbf{Z}_t)^{-1} [\mathbf{f}_\theta(\mathbf{Z}_t, \mathbf{G}) - \mathbf{h}_\phi(\mathbf{Z}_t, t, \mathbf{G})] \quad (31)$$

Note that this is different to the original \mathbf{u} (eq. (14)) by a minus sign. But this does not affect the derivation because we care about \mathbf{u}^2 . By simple manipulation of eq. (27), we have

$$\mathbf{h}_\phi(\mathbf{Z}_t, t, \mathbf{G}) dt + \mathbf{g}_G(\mathbf{Z}_t) d\tilde{\mathbf{W}}_t = \mathbf{f}_G(\mathbf{Z}_t) dt + \mathbf{g}_G(\mathbf{Z}_t) d\mathbf{W}_t \quad (32)$$

This means the prior process (eq. (9)) under probability measure Q is equivalent to the posterior process (eq. (12)) under probability measure P . Next, we can change the probability measure of eq. (29):

$$\begin{aligned} &\sup_{Q \in \mathbb{P}} \left\{ \mathbb{E}_Q \left[\sum_{i=1}^I \log p(\mathbf{X}_{t_i} | \mathbf{Z}_{t_i}, \mathbf{G}) - \frac{1}{2} \int_0^T |\mathbf{u}(\mathbf{Z}_t)|^2 dt \right] \right\} \\ &= \sup_{\mathbf{u}} \left\{ \mathbb{E}_P \left[\sum_{i=1}^I \log p(\mathbf{X}_{t_i} | \tilde{\mathbf{Z}}_{t_i}, \mathbf{G}) - \frac{1}{2} \int_0^T |\mathbf{u}(\tilde{\mathbf{Z}}_t)|^2 dt \right] \right\} \end{aligned}$$

where the second equality is obtained since $\frac{dQ}{dP}$ is fully determined by function \mathbf{u} , and $\tilde{\mathbf{Z}}_t$ is obtained from the posterior process eq. (12). This equation is exactly the term inside $\mathbb{E}_{\mathbf{G} \sim q(\mathbf{G})}[\cdot]$ since $\mathbf{g}_G(\mathbf{Z}_t)^{-2} [\mathbf{f}_G(\mathbf{Z}_t) - \mathbf{h}_\phi(\mathbf{Z}_t, t, \mathbf{G})]^2 = \mathbf{g}_G(\mathbf{Z}_t)^{-2} [\mathbf{h}_\phi(\mathbf{Z}_t, t, \mathbf{G}) - \mathbf{f}_G(\mathbf{Z}_t)]^2$.

From Proposition 2.4.2 in (Dupuis & Ellis, 2011), the supremum is uniquely obtained at

$$\frac{dQ^*}{dP} = \frac{\prod_{i=1}^I p(\mathbf{X}_{t_i} | \mathbf{Z}_{t_i}, \mathbf{G})}{\mathbb{E}_P [\prod_{i=1}^I p(\mathbf{X}_{t_i} | \mathbf{Z}_{t_i}, \mathbf{G})]}.$$

From assumption 4, the measure Q induced by \mathbf{u} can approximate the above arbitrarily well. Thus, the eq. (26) can be written as:

$$\sup_{q(\mathbf{G}), \theta, \phi} \text{ELBO} = \sup_{q(\mathbf{G})} [\log p(\mathbf{X}_{t_1}, \dots, \mathbf{X}_{t_I} | \mathbf{G})] - D_{\text{KL}}[q(\mathbf{G}) \| p(\mathbf{G})]$$

We divide the ELBO by $\frac{1}{I}$, and let $I \rightarrow \infty$, we have

$$\begin{aligned} & \lim_{I \rightarrow \infty} \frac{1}{I} [\log p(\mathbf{X}_{t_1}, \dots, \mathbf{X}_{t_I} | \mathbf{G})] - \frac{1}{I} \text{KL}[q(\mathbf{G}) \| p(\mathbf{G})] \\ &= \lim_{I \rightarrow \infty} \frac{1}{I} [\log p(\mathbf{X}_{t_1}, \dots, \mathbf{X}_{t_I} | \mathbf{G})] \\ &\leq \lim_{I \rightarrow \infty} \frac{1}{I} \log p(\mathbf{X}_{t_1}, \dots, \mathbf{X}_{t_I}; \mathbf{G}^*) \end{aligned}$$

where the first equality is obtained by the fact $D_{\text{KL}}[q(\mathbf{G}) \| p(\mathbf{G})] < \infty$, and the second inequality is due to the property of the ground truth likelihood. From the identifiability theorem 4.2, the equality is uniquely obtained at $q(\mathbf{G}) = \delta(\mathbf{G}^*)$, and the learned system recovers the true generating mechanism under infinite data limits. \square

B MODEL ARCHITECTURE

In this section, we describe the model architecture details used in our experiments for *SCOTCH*.

Prior Drift Function and Diffusion Function As described in Section 3, following Geffner et al. (2022), we use the following design for the prior drift function $\mathbf{f}_{G,d}(\mathbf{Z}_t)$ and diffusion function $\mathbf{g}_{G,d}(\mathbf{Z}_t)$:

$$\mathbf{f}_{G,d}(\mathbf{Z}_t) = \zeta \left(\sum_{i=1}^D G_{i,d} l(\mathbf{Z}_{t,i}, \mathbf{e}_i), \mathbf{e}_d \right) \quad (33)$$

where $\zeta : \mathbb{R}^{D_g \times D_e} \rightarrow \mathbb{R}^D$, $l : \mathbb{R}^{D \times D_e} \rightarrow \mathbb{R}^{D_g}$ are neural networks, and $\mathbf{e}_i \in \mathbb{R}^{D_e}$ is a trainable node embedding for the i^{th} node. The use of node embeddings means that we only need to train two neural networks, regardless of the latent dimensionality D .

We implement both the prior drift and diffusion function using $D_e = D_g = 32$, and as neural networks with two hidden layers of size $\max(2 * D, D_e)$ with residual connections.

Posterior Drift Function In Section 3.1, we described the posterior SDE $d\tilde{\mathbf{Z}}_t^{(n)} = \mathbf{h}_\psi(\tilde{\mathbf{Z}}_t^{(n)}, t; \mathbf{G}, \mathbf{X}^{(n)})dt + \mathbf{g}_G(\tilde{\mathbf{Z}}_t^{(n)})d\mathbf{W}_t$, with posterior drift function $\mathbf{h}_\psi(\tilde{\mathbf{Z}}_t^{(n)}, t; \mathbf{G}, \mathbf{X}^{(n)})$. We now elaborate on how this is implemented.

We design an encoder $\mathbf{K}_\psi(t, \mathbf{G}, \mathbf{X})$, that takes as input the time t , a graph \mathbf{G} and time series $\mathbf{X} = \{\mathbf{X}_{t_1}, \dots, \mathbf{X}_{t_I}\}$, and outputs a *context vector* $\mathbf{c} \in \mathbb{R}^{D_c}$. This encoder consists of a GRU (Cho et al., 2014) that takes as input all future observations (i.e. \mathbf{X}_{t_i} s.t. $t_i > t$) in reverse order; and a single linear layer which takes the input (i) the hidden state of the GRU, and (ii) the flattened graph matrix G , and output the context vector \mathbf{c} . Note that the GRU only takes as input future observations as the future evolution of the latent state is conditionally independent of past observations given the current latent state. We implement the GRU with hidden size 128, and choose $D_c = 64$ for the size of the context vector.

Then, the posterior drift function $\mathbf{h}_\psi(\tilde{\mathbf{Z}}_t^{(n)}, t; \mathbf{G}, \mathbf{X}^{(n)})$ is implemented as a neural network that takes as input $\tilde{\mathbf{Z}}_t^{(n)}$ and the context vector \mathbf{c} computed by the encoder, and outputs a vector of dimension D . This neural network is a MLP with 1 hidden layer of size 128.

Posterior Mean and Covariance In Section 3.1, we also have posterior mean and covariance functions $\boldsymbol{\mu}_\psi(\mathbf{G}, \mathbf{X}^{(n)}) : \{0, 1\}^{D \times D} \times \mathbb{R}^D \rightarrow \mathbb{R}^D$ and $\boldsymbol{\Sigma}_\psi(\mathbf{G}, \mathbf{X}^{(n)}) : \{0, 1\}^{D \times D} \times \mathbb{R}^D \rightarrow \mathbb{R}^{D \times D}$ for the initial state. We reuse the encoder $\mathbf{K}_\psi(t, \mathbf{G}, \mathbf{X})$ with $t = 0$ to encode the entire time series and graph, and then implement $\boldsymbol{\mu}_\psi, \boldsymbol{\Sigma}_\psi$ as a linear transformation of the context vector (i.e. a single linear layer).

Posterior Graph Distribution In Section 3.1, we introduced a variational approximation $q_\phi(\mathbf{G})$ to the true posterior $p(\mathbf{G}|\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(N)})$. To implement this, we use a product of independent Bernoulli distributions for each edge. That is, we have:

$$q_\phi(\mathbf{G}) = \prod_{i,j} \phi_{ij}^{G_{i,j}} (1 - \phi_{ij})^{(1-G_{i,j})} \quad (34)$$

where $\phi_{ij} \in [0, 1]$ are learnable parameters corresponding to the probability of edge $i \rightarrow j$ being present.

Observational Likelihood We choose the observational noise p_{ϵ_d} in the model to follow a standard Laplace distribution with location $\mu = 0$ and scale $b = 0.01$.

C BASELINES

We use the following baselines for all our experiments to evaluate the performance of *SCOTCH*.

- **PCMCI+:** Runge (2018; 2020) proposed a constraint-based causal discovery methods for time series, which leverage the momentary conditional independence test to simultaneously detect the lagged parents and instantaneous effects. This is an improvement over its predecessor called PCMCI, which cannot handle instantaneous effects. In our experiments, we use PCMCI for Netsim and PCMCI+ for the other datasets. We use the opensourced implementation *Tigramite* (<https://github.com/jakobrunge/tigramite>).
- **VARLiNGaM:** Hyvärinen et al. (2010) proposed a linear vector auto-regressive model to learn from time series observations. It is an extension of LiNGaM (Shimizu et al., 2006), where its structural identifiability is guaranteed through the non-Gaussian noise assumption. The major limitation is its linear and discrete nature, which cannot model complex interactions and continuous systems. We also use the opensourced *LiNGaM* package (<https://lingam.readthedocs.io/en/latest/tutorial/var.html>).
- **CUTS:** CUTS (Cheng et al., 2023) is based on Granger causality, and designed for inferring structures from irregularly sampled time series. It treats the irregular samples as a missing data imputation problem. It is capable of imputing missing observations and inferring the graph at the same time. However, it only supports single time series. We use the authors’ opensourced code (<https://github.com/jarrycyx/unn>).
- **Rhino:** Gong et al. (2022) proposed one of the most flexible SEM-based temporal structure learning framework that is capable of modelling (1) lagged parents; (2) history-dependent noise and (3) instantaneous effects. Many SEM-based structure learning approach can be regarded as a special case of Rhino. From the discussion in section 3.2, *SCOTCH* can be regarded as a continuous-time version of Rhino. We use the authors’ opensourced implementation (<https://github.com/microsoft/causica/tree/v0.0.0>).
- **NGM:** NGM (Bellot et al., 2022) proposed to use NeuralODE to learn the mean process of the SDE. Since this is the only baseline we are aware of in terms of structure learning under continuous time, this will be used as our main comparison. We use the authors’ opensourced code (<https://github.com/alexisbellot/Graphical-modelling-continuous-time>).

NGM and CUTS are originally designed for single time series setup and cannot handle multiple time series. For fair comparison, we modify them by concatenating the multiple time series into a single one. That is, given n time series $\{\mathbf{X}^{(n)}\}_{n=1}^N$ with observation times t_1, \dots, t_I , we convert them into a single time series with observation times in $[(n-1) * t_I + t_1, n * t_I]$ for the n^{th} time series. Our assumption is that since their learning routines are batched across time points, and the concatenation points are rarely sampled, this should have small impact to the performance in comparison to the benefit of additional data. Empirically, this approach indeed improves the performance over simply selecting a single time series.

For VARLiNGaM, PCMCI, and Rhino, which cannot handle irregularly sampled data, we use zero-order hold (ZOH) to impute the missing data, which has been found to perform competitively (Cheng et al., 2023) with other imputation methods such as GP regression and GRIN (Cini et al., 2022).

C.1 COMPARISON TO ODE-BASED STRUCTURE LEARNING

In this section, we present an extended version of the example failure case of NGM presented in section 3.2. Bellot et al. (2022) proposed a structure learning method (NGM) for learning from a single time series generated from a SDE. Their approach learns a neural ODE $d\mathbf{Z}(t) = \mathbf{f}_\theta(\mathbf{Z}(t))dt$ that models the mean process of the SDE and extract the graphical structure from the first layer of \mathbf{f}_θ . Given a single observed trajectory $\mathbf{X} = \{\mathbf{X}_{t_i}\}_{i \in [I]}$, they assume that the observed data follows a multivariate Gaussian distribution $(\mathbf{X}_{t_1}, \dots, \mathbf{X}_{t_I}) \sim \mathcal{N}((\mathbf{Z}_{t_1}, \dots, \mathbf{Z}_{t_I}), \Sigma)$ with mean process \mathbf{Z}_t given by the deterministic mean process (ODE), and a diagonal covariance matrix $\Sigma \in \mathbb{R}^{I \times I}$. As such, NGM optimizes the following squared loss:

$$\sum_{i=1}^I \|\mathbf{X}_{t_i} - \mathbf{Z}_{t_i}\|_2^2 \quad (35)$$

Like *SCOTCH*, NGM attempts to model the underlying continuous-time dynamics and can naturally handle irregularly sampled data. However, the Gaussianity assumption only holds when the underlying SDE is linear; that is, SDEs of the form $d\mathbf{X} = (\mathbf{a}(t)\mathbf{X} + \mathbf{b}(t))dt + \mathbf{c}(t)d\mathbf{W}_t$. For general SDEs where the drift and/or diffusion functions are nonlinear functions of the state, the joint distribution can be far from Gaussian, leading to model misspecification, resulting in the incorrect drift function even if the neural network \mathbf{f}_θ has the capacity to express the true drift function.

Another drawback of learning an ODE mean process using the objective in Equation 35 is that it is difficult to generalise to correctly learn from multiple time series, which can be important for recovering the underlying SDEs in practice since a single time series is just a one trajectory sample from the SDE, and thus cannot represent the trajectory multimodality due to stochasticity. In particular, simply computing a batch loss over all time series $\sum_{n=1}^N \sum_{i=1}^I \|\mathbf{X}_{t_i}^{(n)} - \mathbf{Z}_{t_i}\|_2^2$ may fail to recover the underlying dynamics when learning from multiple time series. To demonstrate the above argument, we propose a bi-modal failure case. Consider the following 1D SDE:

$$dX = Xdt + 0.01dW_t \quad (36)$$

where the trajectory will either go upwards or downwards exponentially (bi-modality)

In Figure 2a we show trajectories sampled from this SDE, where the initial state is set to $X_0 = 0$ for all trajectories. The optimal ODE mean process in terms of (batched) squared loss is given by $dZ = 0dt$, whose solution is given by the horizontal axis; in particular, while true graph by definition contains a self-loop, the inferred graph from this ODE has no edges. In Figure 2b we show the ODE mean process \mathbf{f}_θ learned by NGM, together with trajectory samples from the corresponding SDE $dX = \mathbf{f}_\theta(X)dt + 0.01dW_t$. The learned ODE mean process (in black) is close to the horizontal axis (note the scale of the vertical axis), with trajectories that do not match the data. On the other hand, in Figure 2c we see that *SCOTCH* successfully learns the underlying SDE with trajectories closely matching the observed data and demonstrating the bi-modal behavior.

D EXPERIMENTS

D.1 CHOICE OF SDE SOLVER

There are several choices that can affect the accuracy of the SDE solver used for *SCOTCH*. Firstly, discretization step size is an important factor of the solver; a smaller step size generally leads to a more accurate SDE solution, but at the cost of additional time and space complexity. The computational cost (with default Euler discretization) should scale inversely w.r.t. the step size. In the following, we conducted an initial verification run for the Ecol1 dataset with half of the original step size reported below in appendix appendix D.5. Appendix D.1 compares the performance with different step sizes. We can see that $\Delta t = 0.05$ results in similar performance compared to $\Delta t = 0.025$ (while being 2x faster). Therefore, we decide to use the step size $\Delta t = 0.05$. Secondly, we chose to use a pathwise gradient estimator rather than the adjoint method (Li et al., 2020), as we found this was more efficient time-wise and we did not run into space limitations. Although theoretically, they should give the same performance, in practice, the pathwise gradient estimator may have an advantage that computing its gradient does not require solving another SDE, which is subject to the accuracy of the SDE solver. It is also possible to use higher-order numerical solvers such as the Milstein method; however we did not explore this thoroughly in this work.

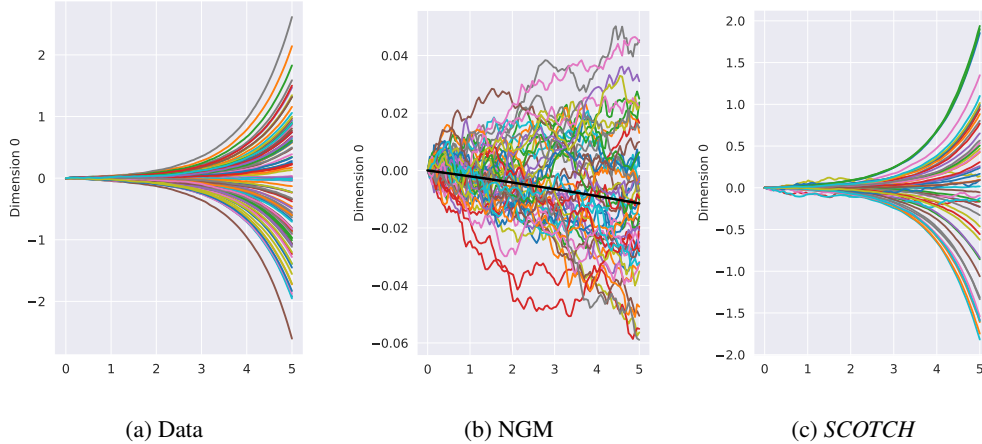


Figure 2: Comparison between NGM and *SCOTCH* for simple SDE (note vertical axis scale)

	AUROC
$\Delta t = 0.025$	0.747±0.005
$\Delta t = 0.05$	0.752±0.008

Table 4: Performance comparisons between different choice of discretization step size Δt for SDE solver.

D.2 COMPARISON TO LATENT SDES

Though appealing at first glance, attempting to directly extract graphical structure from SDEs learned using existing methods, such as that of Li et al. (2020), is very challenging. Firstly, to extract the signature graph, one would have to evaluate the partial derivative of the drift and diffusion networks at every input point in the input domain, which is not practical. Secondly, the learned drift and diffusion functions may have different graphs, and it is unclear how we should combine these. Thirdly, there are no theoretical results to justify this approach (prior to our paper’s theory). For these reasons, prior work does not admit an easy way to extract structure.

In order to construct a simple empirical baseline following this strategy, we follow the setup of Li et al. (2020), and implement each output dimension of the drift and diffusion functions as a separate neural network, i.e.

$$\mathbf{f} = [f_1, \dots, f_D]^T, \mathbf{g} = [g_1, \dots, g_D]^T \tag{37}$$

Using e.g. \mathbf{A}_{g_j} to denote the weight matrix of the first layer of g_j , and $\mathbf{A}_{g_j}^k$ to denote the k^{th} column of that matrix (corresponding to the k^{th} input dimension, then we define:

$$\mathbf{H}_{k,j} = \max(|\mathbf{A}_{f_j}^k|_2, |\mathbf{A}_{g_j}^k|_2) \tag{38}$$

Method	AUROC
PCMCI+	0.530 ± 0.002
NGM	0.611 ± 0.002
CUTS	0.543 ± 0.003
Rhino	0.685 ± 0.003
SCOTCH	0.752 ± 0.008
LSDE	0.496 ± 0.021

Table 5: Performance comparison between methods on DREAM3 Ecoli1 dataset. LSDE refers to latent SDE + extracting first layer weights.

to be our (weighted) estimate of the graph structure. This has the property that whenever $\mathbf{H}_{k,j} = 0$, then $\frac{\partial f_j}{\partial x_k} = 0$ and $\frac{\partial g_j}{\partial x_k} = 0$. This can be extracted from a learned SDE, and we can compute an AUROC using the weights as confidence scores.

Table 5 shows results for this approach (which we call LSDE) in comparison with SCOTCH and other baselines on the DREAM3 Ecoli1 dataset. It can be seen that LSDE performs no better than random guessing at identifying the correct edges.

D.3 SYNTHETIC DATASETS: LORENZ

D.3.1 DATA GENERATION

For the Lorenz dataset, we simulate time-series data according to the following SDE based on the D -dimensional Lorenz-96 system of ODEs:

$$dX_{t,d} = ((X_{t,d+1} - X_{t,d-2})X_{t,d-1} - X_{t,d})dt + F + \sigma dW_{t,i} \quad (39)$$

where $X_{t,-1} := X_{t,D-1}$, $X_{t,0} := X_{t,D}$, and $X_{t,D+1} := X_{t,1}$, with parameters set as $F = 10$ and $\sigma = 0.5$. We generate $N = 100$ 10-dimensional time series, each with length $I = 100$, which are sampled with time interval 1 starting from $t = 0$ (that is, $t_1 = 0, t_2 = 1, \dots, t_{100} = 99$). The initial state $X_{0,i}$ is sampled from a standard Gaussian. To simulate the SDE, we use the Euler-Maruyama scheme with step-size $dt = 0.005$.

For this synthetic dataset, we do not add observation noise to the generated time series.

To produce the irregularly sampled versions of the Lorenz dataset, for each time $t = 0, \dots, 99$, we randomly drop the observed data at that time with probability p , independently at each time t (and for all time series $n = 1, \dots, 100$). We test using $p = 0.3, 0.6$ in our experiments.

D.3.2 HYPERPARAMETERS

SCOTCH We use Adam (Kingma & Ba, 2014) optimizer with learning rate 0.003 and 0.001 for $p = 0.3$ and $p = 0.6$, respectively. We set the $\lambda_s = 500$ and EM discretization step size $\Delta = 1$ for SDE integrator, which coincides with the step size in the data generation process. The time range is set to $[0, 100]$. We enable the residual connections for prior drift and diffusion network. We also adopt a learning rate warm-up schedule, where we linearly increase the learning rate from 0 to the target value within 100 epochs. We do not mini-batch across the time series. We train 5000 epochs for convergence.

NGM We use the same hyperparameter setup as NGM (Bellot et al., 2022) where we set 0.1 for the group lasso regularizer and the learning rate as 0.005. We train NGM for 4000 epochs in total (2000 for the group lasso stage and 2000 for the adaptive group lasso stage).

VARLiNGaM We set the lag to be the same as the ground truth $lag = 1$, and do not prune the inferred adjacency matrix.

PCMCI+ We use *partial correlation* as the underlying conditional independence test. We set the maximum lag at 2, and let the algorithm itself optimise the significance level. We use the threshold 0.07 to determine the graph from the inferred value matrix.

CUTS We use the authors’ suggested hyperparameters (Cheng et al., 2023) for the Lorenz dataset.

Rhino We use hyperparameters with learning rate 0.01, 70 epochs of augmented lagrangian training with 6000 steps each, time lag of 2, sparsity parameter $\lambda_s = 5$, and enable instantaneous effects.

D.3.3 ADDITIONAL RESULTS

Figure 3 shows the curve of other metrics.

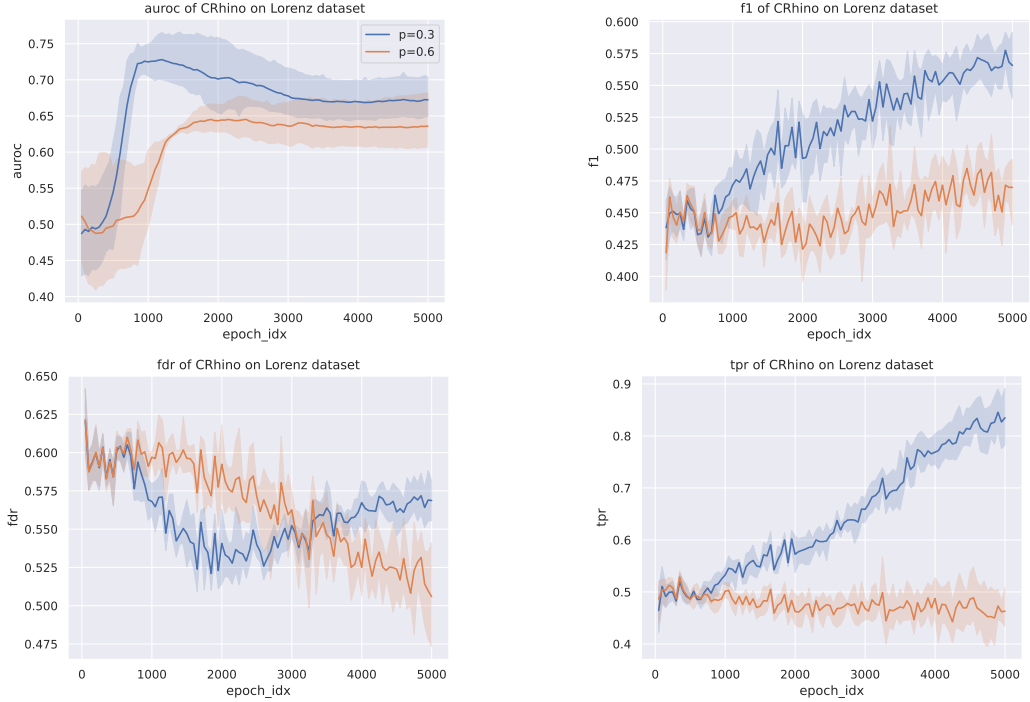


Figure 3: The AUROC (top left), F1 score (top right), false discovery rate (bottom left) and true positive rate (bottom right) curves of *SCOTCH* for Lorenz dataset. The shaded area indicates the 95% confidence intervals. Blue color indicates the dataset with missing probability 0.3 and orange color indicates missing probability 0.6.

D.4 SYNTHETIC DATASETS: GLYCOLYSIS

D.4.1 DATA GENERATION

In this synthetic experiment, we generate data according to the system presented by Daniels & Nemenman (2015), which models a glycolytic oscillator. This is a $D = 7$ dimensional system with the following equations:

$$\begin{aligned}
 dX_{t,1} &= \left(2.5 - \frac{100X_{t,1}X_{t,6}}{1 + (X_{t,6}/0.52)^4} \right) dt + 0.01dW_{t,1} \\
 dX_{t,2} &= \left(\frac{200X_{t,1}X_{t,6}}{1 + (X_{t,6}/0.52)^4} - 6X_{t,2}(1 - X_{t,5}) - 12X_{t,2}X_{t,5} \right) dt + 0.01dW_{t,2} \\
 dX_{t,3} &= (6X_{t,2}(1 - X_{t,5}) - 16X_{t,3}(4 - X_{t,6})) dt + 0.01dW_{t,3} \\
 dX_{t,4} &= (16X_{t,3}(4 - X_{t,6}) - 100X_{t,4}X_{t,5} - 13(X_{t,4} - X_{t,7})) dt + 0.01dW_{t,4} \\
 dX_{t,5} &= (6X_{t,2}(1 - X_{t,5}) - 100X_{t,4}X_{t,5} - 12X_{t,2}X_{t,5}) dt + 0.01dW_{t,5} \\
 dX_{t,6} &= \left(-\frac{200X_{t,1}X_{t,6}}{1 + (X_{t,6}/0.52)^4} + 32X_{t,3}(4 - X_{t,6}) - 1.28X_{t,6} \right) dt + 0.01dW_{t,6} \\
 dX_{t,7} &= (1.3(X_{t,4} - X_{t,7}) - 1.8X_{t,7}) dt + 0.01dW_{t,7}
 \end{aligned}$$

As with the Lorenz dataset, we simulate $N = 100$ time series of length $I = 100$, starting at $t = 0$ and with time interval 1. The initial state is sampled uniformly from the ranges $X_{0,1} \in [0.15, 1.60]$, $X_{0,2} \in [0.19, 2.16]$, $X_{0,3} \in [0.04, 0.20]$, $X_{0,4} \in [0.10, 0.35]$, $X_{0,5} \in [0.08, 0.30]$, $X_{0,6} \in [0.14, 2.67]$, $X_{0,7} \in [0.05, 0.10]$, as indicated in Daniels & Nemenman (2015). To simulate the SDE, we use the Euler-Maruyama scheme with step-size $dt = 0.005$.

For this synthetic dataset, we do not add observation noise to the generated time series.

D.4.2 HYPERPARAMETERS

SCOTCH We use the same hyperparameter as Lorenz experiments. The only differences are that we use learning rate 0.001 and set $\lambda_s = 200$. We train *SCOTCH* for 30000 epochs for convergence.

NGM Since Bellot et al. (2022) did not release the hyperparameters for their glycolysis experiment, we use the default setup in their code. They are the same as the hyperparameters in Lorenz experiments.

VARLiNGaM Same as Lorenz experiment setup.

PCMCI+ Same as Lorenz experiment setup.

CUTS Same as Lorenz experiment setup.

Rhino Same as Lorenz experiment setup.

D.4.3 ADDITIONAL RESULTS

Table 6 shows the performance comparison of *SCOTCH* to NGM with the original glycolysis data, where the data have different variable scales. We can observe that this difference in scale does not affect the AUROC of *SCOTCH* but greatly affects NGM. Since AUROC is threshold free, we can see that *SCOTCH* is more robust in terms of scaling compared to NGM. A possible reason is that the stochastic evolution of the variables in SDE can help stabilise the training when encountering difference in scales, but ODE can easily overshoot due to its deterministic nature.

Figure 4 shows the curves of different metrics. Interestingly, we can see that data normalisation does not improve the AUROC performance (compared to NGM), but does increase the f1 score. This may be because f1 is threshold sensitive and the default threshold of 0.5 might not be optimal. We can see this through the TPR plot, where "Original" has very low value.

Table 6: Performance comparison with original Glycolysis data

	AUROC	TPR \uparrow	FDR \downarrow
<i>SCOTCH</i>	0.7352\pm0.019	0.3623\pm0.007	0.1575\pm0.05
NGM	0.5248 \pm 0.057	0.3478 \pm 0.035	0.4559 \pm 0.094

D.5 DREAM3 DATASET

In this appendix, we will include experiment setups, hyperparameters and additional plots for Dream3 experiment.

D.5.1 HYPERPARAMETERS

SCOTCH We follow similar setup as Lorenz experiment. The differences are that the learning rate is 0.001. The time range is set to $[0, 1.05]$ with EM discretization step size 0.05, which results in exactly 21 observations for each time series. We choose sparsity coefficient $\lambda_s = 200$. For all sub-datasets, we normalize the data to have 0 mean and unit variance for each dimension. We use the above hyperparameters for Ecoli1, Ecoli2 and Yeast1 sub-datasets. For Yeast2, we only change the learning rate to be 0.0005. For Yeast3, we change the $\lambda_s = 50$. We train *SCOTCH* for 30000 epochs until convergence.

NGM For NGM, we follow the same hyperparameter setup as (Cheng et al., 2023), where we set the group lasso regulariser as 0.05, learning rate 0.005. We train NGM with 4000 epochs (2000 each for group lasso and adaptive group lasso stages). For fair comparison, we use the same observation time (i.e. equally spaced time points within $[1, 1.05]$ and step size 0.05).

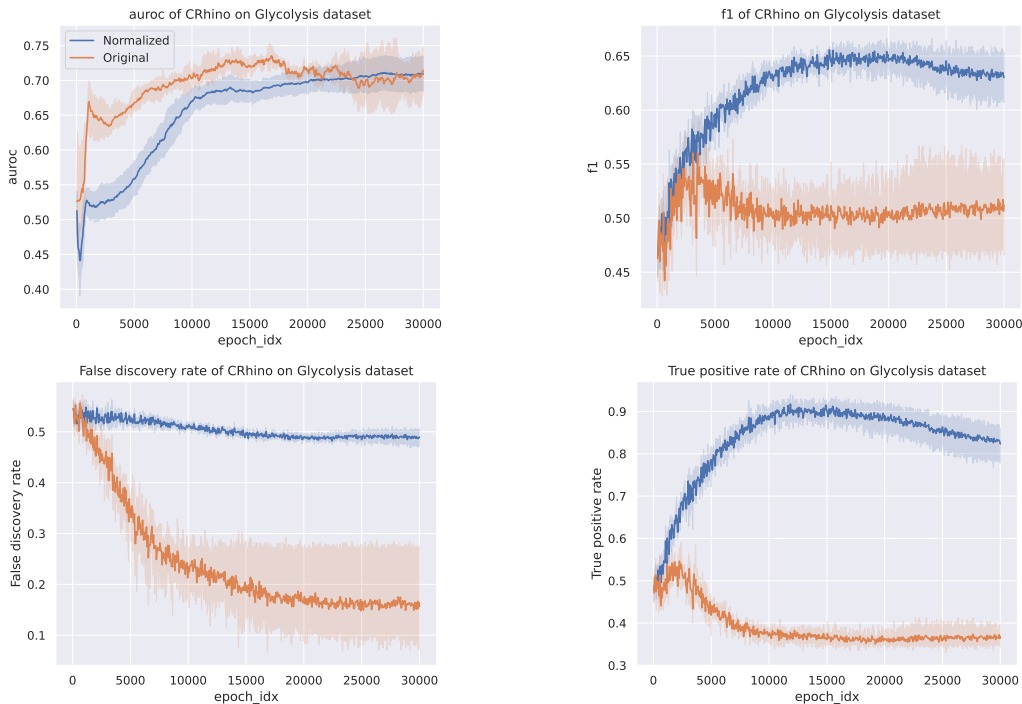


Figure 4: The AUROC (top left), F1 score (top right), false discovery rate (bottom left) and true positive rate (bottom right) curves of *SCOTCH* for Glycolysis dataset. The shaded area indicates the 95% confidence intervals. Blue color indicates the normalized dataset and orange color indicates the original dataset.

PCMCi+ and Rhino As the experiment setup is the same, we directly cite the number from Gong et al. (2022).

CUTS We use the authors’ suggested hyperparameters (Cheng et al., 2023) for the DREAM3 datasets.

D.5.2 ADDITIONAL PLOTS

In this section, we include additional metric curves of *SCOTCH* in fig. 5. Each curve is obtained by averaging over 5 runs and the shaded area indicates the 95% confidence interval. From the value of f1 score, FDR and TPR, we can see DREAM3 is indeed a challenging dataset, where all f1 scores are below 0.5 and FDR only drops to 0.7. From the TPR plot, it is expected to drop at the beginning and then increase during training, which is the case for Ecoli1, Ecoli2 and Yeast1. TPR corresponds well to AUROC and F1 score, since Ecoli1, Ecoli2 and Yeast1 have much better values compared to Yeast2 and Yeast3.

D.6 NETSIM

D.6.1 EXPERIMENT SETUP

For the Netsim dataset, we generate the missing data versions in the same way as the Lorenz dataset (see appendix D.3).

D.6.2 HYPERPARAMETERS

SCOTCH We use similar hyperparameter setup as Dream3 (appendix D.5.1), but we change $\lambda_s = 1000$ and use the raw data without normalisation. We train *SCOTCH* for 10000 epochs.

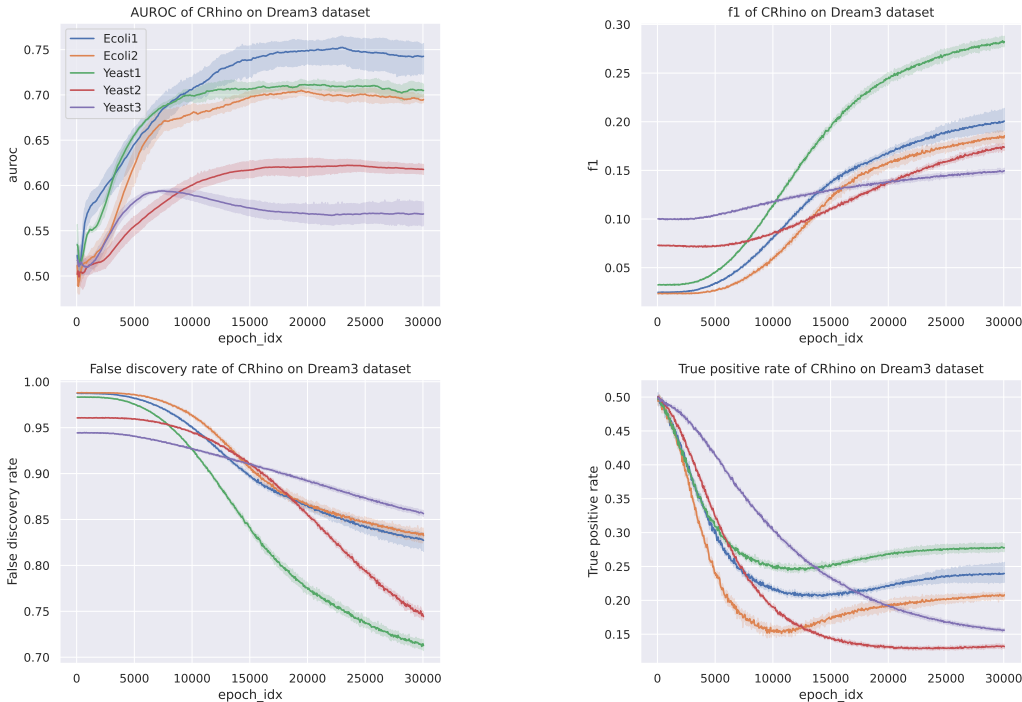


Figure 5: The AUROC (top left), F1 score (top right), false discovery rate (bottom left) and true positive rate (bottom right) curves of *SCOTCH* for each DREAM3 sub-datasets. The shaded area indicates the 95% confidence intervals.

NGM We follow the same setup as DREAM3 experiment, which also coincides with the setup used in Cheng et al. (2023).

PCMCI We follow the same setup as Lorenz and use threshold 0.07 to infer the graph.

CUTS We use the authors’ suggested hyperparameters (Cheng et al., 2023) for the Netsim dataset.

Rhino and Rhino+NoInst We directly cite the number from Gong et al. (2022) for the full dataset, and use the same hyperparameters as Gong et al. (2022) for both $p = 0.1$ and $p = 0.2$ Netsim datasets.

D.6.3 ADDITIONAL PLOTS

We include additional metric curves of *SCOTCH* on Netsim dataset in fig. 6. From the plot, we can see Netsim is a easier dataset compared to DREAM3 since the dimensionality is much smaller. An interesting observation is f1 score does not necessarily correspond well to auROC since f1 score is threshold dependent (by default we use 0.5) but not auROC. To evaluate the robustness of the model, we decide to report AUROC instead of f1 score.

E INTERVENTIONS

Aside from learning the graphical structure between variables, one might also be interested in analysing the effect of applying external changes, or interventions, to the system. Broadly speaking, there are two types of interventions that we can consider in a continuous-time model. The first is to intervene on the dynamics (that is, the drift or diffusion functions), possibly for a set period of time. The second is to directly intervene on the value of (some subset of) variables. The goal is to employ our learned *SCOTCH* model in order to predict the effect of these interventions on the underlying system.

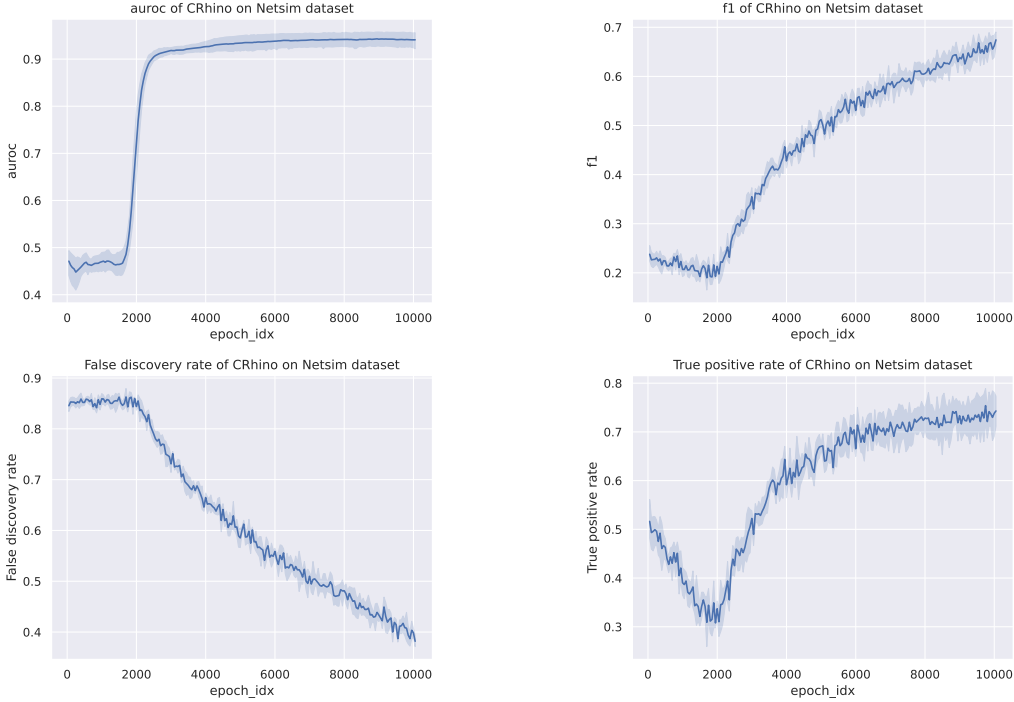


Figure 6: The AUROC (top left), F1 score (top right), false discovery rate (bottom left) and true positive rate (bottom right) curves of *SCOTCH* for Netsim dataset. The shaded area indicates the 95% confidence intervals.

The former is easy to implement as we need only replace (parts of) the learned drift/diffusion function with the intervention. However, the latter is slightly more subtle than it might first appear. (Hansen & Sokol, 2014) proposed to define such an intervention as a function that fixes the value of a particular variable as a function of the other variables. However, it is unclear how we can generalize this to interventions affecting more than one variable. For example, an intervention policy $Z_1 \leftarrow Z_2, Z_2 \leftarrow Z_1 + 1$ creates a feedback loop whose semantics are not easy to resolve. Thus, we propose the following definition:

Definition 3 (State-space Intervention). *Given a D -dimensional SDE, a state-space intervention is an idempotent function $\iota(t, \mathbf{Z}) : \mathbb{R}^{D+1} \rightarrow \mathbb{R}^D$; that is, $\iota(t, \iota(t, \mathbf{Z})) \equiv \iota(t, \mathbf{Z})$. The corresponding intervened stochastic process is defined by:*

$$\tilde{\mathbf{Z}}_t = \iota \left(t, \tilde{\mathbf{Z}}_0 + \sum_{d \in [D]} \int_0^t \mathbf{f}(\tilde{\mathbf{Z}}_s) ds + \sum_{d \in [D]} \int_0^t \mathbf{g}(\tilde{\mathbf{Z}}_s) d\mathbf{W}_s \right) \quad (40)$$

The requirement of idempotence captures the intuition that applying the same intervention twice should result in the same result. Some examples of interventions are given as follows:

- **Identity:** If $\iota(t, \mathbf{Z}) = \mathbf{Z} \forall t \in [T_1, T_2], \mathbf{Z} \in \mathbb{R}^D$, then the system evolves according to the original SDE in this time period, with initial state $\tilde{\mathbf{Z}}_{T_1}$.
- **Ordered Intervention:** Given some ordered subset of the variables, we can consider intervening on each variable in order, as a function of the previous variables in the order. That is, we restrict each dimension ι_i of the intervention output to be of the form

$$\iota_i(t, \mathbf{Z}_t) = \iota_i(t, \mathbf{Z}_{t, < i}) \quad (41)$$

where $\mathbf{Z}_{t, < i} = \{\mathbf{Z}_{t, j} : j < i\}$. It can easily be seen that ι is always idempotent in this case.

- **Projection:** Another example of an idempotent function is a projection. This could simulate a setting where external force is applied to ensure the SDE trajectories satisfy spatial

constraints. Note that a projection cannot necessarily be expressed as an ordered intervention (e.g. consider projection onto a sphere).

In practice, we implement state-space interventions in SDEs learned from *SCOTCH* by modifying the SDE solver (e.g. Euler-Maruyama) such that each step is followed with an intervention assignment $\mathbf{Z}_t \leftarrow \iota(t, \mathbf{Z}_t)$.