

Supplementary material

In sections A and B, we give a more thorough introduction to solving CDEs via the log-ODE method.

In section C we discuss the experimental details such as the choice of network structure, computing infrastructure and hyperparameter selection approach.

In section D we give a full breakdown of every experimental result.

A AN INTRODUCTION TO THE LOG-ODE METHOD FOR CONTROLLED DIFFERENTIAL EQUATIONS

The log-ODE method is an effective method for approximating the controlled differential equation:

$$\begin{aligned} dY_t &= f(Y_t) dX_t, \\ Y_0 &= \xi, \end{aligned} \quad (11)$$

where $X : [0, T] \rightarrow \mathbb{R}^d$ has finite length, $\xi \in \mathbb{R}^n$ and $f : \mathbb{R}^n \rightarrow L(\mathbb{R}^d, \mathbb{R}^n)$ is a function with certain smoothness assumptions so that the CDE (11) is well posed. Throughout these appendices, $L(U, V)$ denotes the space of linear maps between the vector spaces U and V . In rough path theory, the function f is referred to as the “vector field” of (11) and usually assumed to have $\text{Lip}(\gamma)$ regularity (see definition 10.2 in Friz & Victoir (2010)). In this section, we assume one of the below conditions on the vector field:

1. f is bounded and has N bounded derivatives.
2. f is linear.

In order to define the log-ODE method, we will first consider the tensor algebra and path signature.

Definition A.1 We say that $T(\mathbb{R}^d) := \mathbb{R} \oplus \mathbb{R}^d \oplus (\mathbb{R}^d)^{\otimes 2} \oplus \dots$ is the tensor algebra of \mathbb{R}^d and $T((\mathbb{R}^d)) := \{\mathbf{a} = (a_0, a_1, \dots) : a_k \in (\mathbb{R}^d)^{\otimes k} \forall k \geq 0\}$ is the set of formal series of tensors of \mathbb{R}^d . Moreover, $T(\mathbb{R}^d)$ and $T((\mathbb{R}^d))$ can be endowed with the operations of addition and multiplication. Given $\mathbf{a} = (a_0, a_1, \dots)$ and $\mathbf{b} = (b_0, b_1, \dots)$, we have

$$\mathbf{a} + \mathbf{b} = (a_0 + b_0, a_1 + b_1, \dots), \quad (12)$$

$$\mathbf{a} \otimes \mathbf{b} = (c_0, c_1, c_2, \dots), \quad (13)$$

where for $n \geq 0$, the n -th term $c_n \in (\mathbb{R}^d)^{\otimes n}$ can be written using the usual tensor product as

$$c_n := \sum_{k=0}^n a_k \otimes b_{n-k}.$$

The operation \otimes given by (13) is often referred to as the “tensor product”.

Definition A.2 The signature of a finite length path $X : [0, T] \rightarrow \mathbb{R}^d$ over the interval $[s, t]$ is defined as the following collection of iterated (Riemann-Stieltjes) integrals:

$$S_{s,t}(X) := \left(1, X_{s,t}^{(1)}, X_{s,t}^{(2)}, X_{s,t}^{(3)}, \dots\right) \in T((\mathbb{R}^d)), \quad (14)$$

where for $n \geq 1$,

$$X_{s,t}^{(n)} := \int \dots \int_{s < u_1 < \dots < u_n < t} dX_{u_1} \otimes \dots \otimes dX_{u_n} \in (\mathbb{R}^d)^{\otimes n}.$$

Similarly, we can define the depth- N (or truncated) signature of the path X on $[s, t]$ as

$$S_{s,t}^N(X) := \left(1, \int_{s < u_1 < t} dX_{u_1}, \dots, \int \dots \int_{s < u_1 < \dots < u_N < t} dX_{u_1} \otimes \dots \otimes dX_{u_N}\right) \in T^N(\mathbb{R}^d), \quad (15)$$

where $T^N(\mathbb{R}^d) := \mathbb{R} \oplus \mathbb{R}^d \oplus (\mathbb{R}^d)^{\otimes 2} \oplus \dots \oplus (\mathbb{R}^d)^{\otimes N}$ denotes the truncated tensor algebra.

The (truncated) signature provides a natural feature set that describes the effects a path X has on systems that can be modelled by (11). That said, defining the log-ODE method actually requires the so-called “log-signature” which efficiently encodes the same integral information as the signature. The log-signature is obtained from the path’s signature by removing certain algebraic redundancies, such as

$$\int_0^t \int_0^s dX_u^i dX_s^j + \int_0^t \int_0^s dX_u^j dX_s^i = X_t^i X_t^j,$$

for $i, j \in \{1, \dots, d\}$, which follows by the integration-by-parts formula. To this end, we will define the logarithm map on the depth- N truncated tensor algebra $T^N(\mathbb{R}^d) := \mathbb{R} \oplus \mathbb{R}^d \oplus \dots \oplus (\mathbb{R}^d)^{\otimes N}$.

Definition A.3 (The logarithm of a formal series) For $\mathbf{a} = (a_0, a_1, \dots) \in T((\mathbb{R}^d))$ with $a_0 > 0$, define $\log(\mathbf{a})$ to be the element of $T((\mathbb{R}^d))$ given by the following series:

$$\log(\mathbf{a}) := \log(a_0) + \sum_{n=1}^{\infty} \frac{(-1)^n}{n} \left(\mathbf{1} - \frac{\mathbf{a}}{a_0} \right)^{\otimes n}, \quad (16)$$

where $\mathbf{1} = (1, 0, \dots)$ is the unit element of $T((\mathbb{R}^d))$ and $\log(a_0)$ is viewed as $\log(a_0)\mathbf{1}$.

Definition A.4 (The logarithm of a truncated series) For $\mathbf{a} = (a_0, a_1, \dots, a_N) \in T((\mathbb{R}^d))$ with $a_0 > 0$, define $\log^N(\mathbf{a})$ to be the element of $T^N(\mathbb{R}^d)$ defined from the logarithm map (16) as

$$\log^N(\mathbf{a}) := P_N(\log(\tilde{\mathbf{a}})), \quad (17)$$

where $\tilde{\mathbf{a}} := (a_0, a_1, \dots, a_N, 0, \dots) \in T((\mathbb{R}^d))$ and P_N denotes the standard projection map from $T((\mathbb{R}^d))$ onto $T^N(\mathbb{R}^d)$.

Definition A.5 The log-signature of a finite length path $X : [0, T] \rightarrow \mathbb{R}^d$ over the interval $[s, t]$ is defined as $\text{LogSig}_{s,t}(X) := \log(S_{s,t}(X))$, where $S_{s,t}(X)$ denotes the path signature of X given by Definition A.2. Likewise, the depth- N (or truncated) log-signature of X is defined for each $N \geq 1$ as $\text{LogSig}_{s,t}^N(X) := \log^N(S_{s,t}^N(X))$.

The final ingredient we use to define the log-ODE method are the derivatives of the vector field f . It is worth noting that these derivatives also naturally appear in the Taylor expansion of (11).

Definition A.6 (Vector field derivatives) We define $f^{\circ k} : \mathbb{R}^n \rightarrow L((\mathbb{R}^d)^{\otimes k}, \mathbb{R}^n)$ recursively by

$$\begin{aligned} f^{\circ(0)}(y) &:= y, \\ f^{\circ(1)}(y) &:= f(y), \\ f^{\circ(k+1)}(y) &:= D(f^{\circ k})(y)f(y), \end{aligned}$$

for $y \in \mathbb{R}^n$, where $D(f^{\circ k})$ denotes the Fréchet derivative of $f^{\circ k}$.

Using these definitions, we can describe two closely related numerical methods for the CDE (11).

Definition A.7 (The Taylor method) Given the CDE (11), we can use the path signature of X to approximate the solution Y on an interval $[s, t]$ via its truncated Taylor expansion. That is, we use

$$\text{Taylor}(Y_s, f, S_{s,t}^N(X)) := \sum_{k=0}^N f^{\circ k}(Y_s) \pi_k(S_{s,t}^N(X)), \quad (18)$$

as an approximation for Y_t where each $\pi_k : T^N(\mathbb{R}^d) \rightarrow (\mathbb{R}^d)^{\otimes k}$ is the projection map onto $(\mathbb{R}^d)^{\otimes k}$.

Definition A.8 (The Log-ODE method) Using the Taylor method (18), we can define the function $\hat{f} : \mathbb{R}^n \rightarrow L(T^N(\mathbb{R}^d), \mathbb{R}^n)$ by $\hat{f}(z) := \text{Taylor}(z, f, \cdot)$. By applying \hat{f} to the truncated log-signature of the path X over an interval $[s, t]$, we can define the following ODE on $[0, 1]$

$$\begin{aligned} \frac{dz}{du} &= \hat{f}(z) \text{LogSig}_{s,t}^N(X), \\ z(0) &= Y_s. \end{aligned} \quad (19)$$

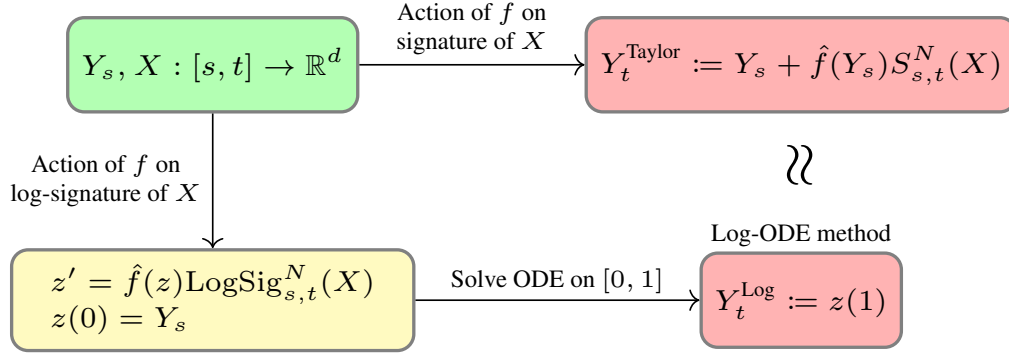


Figure 4: Illustration of the log-ODE and Taylor methods for controlled differential equations.

Then the log-ODE approximation of Y_t (given Y_s and $\text{LogSig}_{s,t}^N(X)$) is defined as

$$\text{LogODE}(Y_s, f, \text{LogSig}_{s,t}^N(X)) := z(1). \quad (20)$$

Remark A.9 Our assumptions of f ensure that $z \mapsto \hat{f}(z)\text{LogSig}_{s,t}^N(X)$ is either globally bounded and Lipschitz continuous or linear. Hence both the Taylor and log-ODE methods are well defined.

Remark A.10 It is well known that the log-signature of a path X lies in a certain free Lie algebra (this is detailed in section 2.2.4 of Lyons et al. (2007)). Furthermore, it is also a theorem that the Lie bracket of two vector fields is itself a vector field which doesn't depend on choices of basis. By expressing $\text{LogSig}_{s,t}^N(X)$ using a basis of the free Lie algebra, it can be shown that only the vector field f and its (iterated) Lie brackets are required to construct the log-ODE vector field $\hat{f}(z)\text{LogSig}_{s,t}^N(X)$. In particular, this leads to our construction of the log-ODE (8) using the Lyndon basis of the free Lie algebra (see Reizenstein (2017) for a precise description of the Lyndon basis). We direct the reader to Lyons (2014) and Boutaib et al. (2014) for further details on this Lie theory.

To illustrate the log-ODE method, we give two examples:

Example A.11 (The “increment-only” log-ODE method) When $N = 1$, the ODE (19) becomes

$$\begin{aligned} \frac{dz}{du} &= f(z)X_{s,t}, \\ z(0) &= Y_s. \end{aligned}$$

Therefore we see that this “increment-only” log-ODE method is equivalent to driving the original CDE (11) by a piecewise linear approximation of the control path X . This is a classical approach for stochastic differential equations (i.e. when $X_t = (t, W_t)$ with W denoting a Brownian motion) and is an example of a Wong-Zakai approximation (see Wong & Zakai (1965) for further details).

Example A.12 (An application for SDE simulation) Consider the following affine SDE,

$$\begin{aligned} dY_t &= a(b - y_t)dt + \sigma y_t \circ dW_t, \\ y(0) &= y_0 \in \mathbb{R}_{\geq 0}, \end{aligned} \quad (21)$$

where $a, b \geq 0$ are the mean reversion parameters, $\sigma \geq 0$ is the volatility and W denotes a standard real-valued Brownian motion. The \circ means that this SDE is understood in the Stratonovich sense. The SDE (21) is known in the literature as Inhomogeneous Geometric Brownian Motion (or IGBM). Using the control path $X = \{(t, W_t)\}_{t \geq 0}$ and setting $N = 3$, the log-ODE (19) becomes

$$\begin{aligned} \frac{dz}{du} &= a(b - z_u)h + \sigma z_u W_{s,t} - ab\sigma A_{s,t} + ab\sigma^2 L_{s,t}^{(1)} + a^2 b\sigma L_{s,t}^{(2)}, \\ z(0) &= Y_s. \end{aligned}$$

where $h := t - s$ denotes the step size and the random variables $A_{s,t}, L_{s,t}^{(1)}, L_{s,t}^{(2)}$ are given by

$$\begin{aligned} A_{s,t} &:= \int_s^t W_{s,r} dr - \frac{1}{2}hW_{s,t}, \\ L_{s,t}^{(1)} &:= \int_s^t \int_s^r W_{s,v} \circ dW_v dr - \frac{1}{2}W_{s,t}A_{s,t} - \frac{1}{6}hW_{s,t}^2, \\ L_{s,t}^{(2)} &:= \int_s^t \int_s^r W_{s,v} dv dr - \frac{1}{2}hA_{s,t} - \frac{1}{6}h^2W_{s,t}. \end{aligned}$$

In Foster et al. (2020), the depth-3 log-signature of $X = \{(t, W_t)\}_{t \geq 0}$ was approximated so that the above log-ODE method became practical and this numerical scheme exhibited state-of-the-art convergence rates. For example, the approximation error produced by 25 steps of the high order log-ODE method was similar to the error of the “increment only” log-ODE method with 1000 steps.

B CONVERGENCE OF THE LOG-ODE METHOD FOR ROUGH DIFFERENTIAL EQUATIONS

In this section, we shall present “rough path” error estimates for the log-ODE method. In addition, we will discuss the case when the vector fields governing the rough differential equation are linear. We begin by stating the main result of Bouteib et al. (2014) which quantifies the approximation error of the log-ODE method in terms of the regularity of the systems vector field f and control path X . Since this section uses a number of technical definitions from rough path theory, we recommend Lyons et al. (2007) as an introduction to the subject.

For $T > 0$, we will use the notation $\Delta_T := \{(s, t) \in [0, T]^2 : s < t\}$ to denote a rescaled 2-simplex.

Theorem B.1 (Lemma 15 in Bouteib et al. (2014)) *Consider the rough differential equation*

$$\begin{aligned} dY_t &= f(Y_t) dX_t, \\ Y_0 &= \xi, \end{aligned} \tag{22}$$

where we make the following assumptions:

- X is a geometric p -rough path in \mathbb{R}^d , that is $X : \Delta_T \rightarrow T^{\lfloor p \rfloor}(\mathbb{R}^d)$ is a continuous path in the tensor algebra $T^{\lfloor p \rfloor}(\mathbb{R}^d) := \mathbb{R} \oplus \mathbb{R}^d \oplus (\mathbb{R}^d)^{\otimes 2} \oplus \dots \oplus (\mathbb{R}^d)^{\otimes \lfloor p \rfloor}$ with increments

$$\begin{aligned} X_{s,t} &= \left(1, X_{s,t}^{(1)}, X_{s,t}^{(2)}, \dots, X_{s,t}^{(\lfloor p \rfloor)}\right), \\ X_{s,t}^{(k)} &:= \pi_k(X_{s,t}), \end{aligned} \tag{23}$$

where $\pi_k : T^{\lfloor p \rfloor}(\mathbb{R}^d) \rightarrow (\mathbb{R}^d)^{\otimes k}$ is the projection map onto $(\mathbb{R}^d)^{\otimes k}$, such that there exists a sequence of continuous finite variation paths $x_n : [0, T] \rightarrow \mathbb{R}^d$ whose truncated signatures converge to X in the p -variation metric:

$$d_p(S^{\lfloor p \rfloor}(x_n), X) \rightarrow 0, \tag{24}$$

as $n \rightarrow \infty$, where the p -variation between two continuous paths Z^1 and Z^2 in $T^{\lfloor p \rfloor}(\mathbb{R}^d)$ is

$$d_p(Z^1, Z^2) := \max_{1 \leq k \leq \lfloor p \rfloor} \sup_{\mathcal{D}} \left(\sum_{t_i \in \mathcal{D}} \left\| \pi_k(Z_{t_i, t_{i+1}}^1) - \pi_k(Z_{t_i, t_{i+1}}^2) \right\|^{\frac{p}{k}} \right)^{\frac{k}{p}}, \tag{25}$$

where the supremum is taken over all partitions \mathcal{D} of $[0, T]$ and the norms $\|\cdot\|$ must satisfy (up to some constant)

$$\|a \otimes b\| \leq \|a\| \|b\|,$$

for $a \in (\mathbb{R}^d)^{\otimes n}$ and $b \in (\mathbb{R}^d)^{\otimes m}$. For example, we can take $\|\cdot\|$ to be the projective or injective tensor norms (see Propositions 2.1 and 3.1 in Ryan (2002)).

- The solution Y and its initial value ξ both take their values in \mathbb{R}^n .
- The collection of vector fields $\{f_1, \dots, f_d\}$ on \mathbb{R}^n are denoted by $f : \mathbb{R}^n \rightarrow L(\mathbb{R}^n, \mathbb{R}^d)$, where $L(\mathbb{R}^n, \mathbb{R}^d)$ is the space of linear maps from \mathbb{R}^n to \mathbb{R}^d . We will assume that f has $\text{Lip}(\gamma)$ regularity with $\gamma > p$. That is, f is bounded with $\lfloor \gamma \rfloor$ bounded derivatives, the last being Hölder continuous with exponent $(\gamma - \lfloor \gamma \rfloor)$. Hence the following norm is finite:

$$\|f\|_{\text{Lip}(\gamma)} := \max_{0 \leq k \leq \lfloor \gamma \rfloor} \|D^k f\|_\infty \vee \|D^{\lfloor \gamma \rfloor} f\|_{(\gamma - \lfloor \gamma \rfloor)\text{-Höl}}, \quad (26)$$

where $D^k f$ is the k -th (Fréchet) derivative of f and $\|\cdot\|_{\alpha\text{-Höl}}$ is the standard α -Hölder norm with $\alpha \in (0, 1)$.

- The RDE (22) is defined in the Lyon's sense. Therefore by the Universal Limit Theorem (see Theorem 5.3 in Lyons et al. (2007)), there exists a unique solution $Y : [0, T] \rightarrow \mathbb{R}^n$.

We define the log-ODE for approximating the solution Y over an interval $[s, t] \subset [0, T]$ as follows:

1. Compute the depth- $\lfloor \gamma \rfloor$ log-signature of the control path X over $[s, t]$. That is, we obtain $\text{LogSig}_{s,t}^{\lfloor \gamma \rfloor}(X) := \log_{\lfloor \gamma \rfloor}(S_{s,t}^{\lfloor \gamma \rfloor}(X)) \in T^{\lfloor \gamma \rfloor}(\mathbb{R}^d)$, where $\log_{\lfloor \gamma \rfloor}(\cdot)$ is defined by projecting the standard tensor logarithm map onto $\{a \in T^{\lfloor \gamma \rfloor}(\mathbb{R}^d) : \pi_0(a) > 0\}$.
2. Construct the following (well-posed) ODE on the interval $[0, 1]$,

$$\begin{aligned} \frac{dz^{s,t}}{du} &= F(z^{s,t}), \\ z_0^{s,t} &= Y_s, \end{aligned} \quad (27)$$

where the vector field $F : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is defined from the log-signature as

$$F(z) := \sum_{k=1}^{\lfloor \gamma \rfloor} f^{\circ k}(z) \pi_k(\text{LogSig}_{s,t}^{\lfloor \gamma \rfloor}(X)). \quad (28)$$

Recall that $f^{\circ k} : \mathbb{R}^n \rightarrow L((\mathbb{R}^d)^{\otimes k}, \mathbb{R}^n)$ was defined previously in Definition A.6.

Then we can approximate Y_t using the $u = 1$ solution of (27). Moreover, there exists a universal constant $C_{p,\gamma}$ depending only on p and γ such that

$$\|Y_t - z_1^{s,t}\| \leq C_{p,\gamma} \|f\|_{\text{Lip}(\gamma)} \|X\|_{p\text{-var};[s,t]}^\gamma, \quad (29)$$

where $\|\cdot\|_{p\text{-var};[s,t]}$ is the p -variation norm defined for paths in $T^{\lfloor p \rfloor}(\mathbb{R}^d)$ by

$$\|X\|_{p\text{-var};[s,t]} := \max_{1 \leq k \leq \lfloor p \rfloor} \sup_{\mathcal{D}} \left(\sum_{t_i \in \mathcal{D}} \|X_{t_i, t_{i+1}}^k\|^{\frac{p}{k}} \right)^{\frac{k}{p}}, \quad (30)$$

with the supremum taken over all partitions \mathcal{D} of $[s, t]$.

Remark B.2 If the vector fields $\{f_1, \dots, f_d\}$ are linear, then it immediately follows that F is linear.

Although the above theorem requires some sophisticated theory, it has a simple conclusion - namely that log-ODEs can approximate controlled differential equations. That said, the estimate (29) does not directly apply when the vector fields $\{f_i\}$ are linear as they would be unbounded. Fortunately, it is well known that linear RDEs are well posed and the growth of their solutions can be estimated.

Theorem B.3 (Theorem 10.57 in Friz & Victoir (2010)) Consider the linear RDE on $[0, T]$

$$\begin{aligned} dY_t &= f(Y_t) dX_t, \\ Y_0 &= \xi, \end{aligned}$$

where X is a geometric p -rough path in \mathbb{R}^d , $\xi \in \mathbb{R}^n$ and the vector fields $\{f_i\}_{1 \leq i \leq d}$ take the form $f_i(y) = A_i y + B$ where $\{A_i\}$ and $\{B_i\}$ are $n \times n$ matrices. Let K denote an upper bound on

$\max_i(\|A_i\| + \|B_i\|)$. Then a unique solution $Y : [0, T] \rightarrow \mathbb{R}^n$ exists. Moreover, it is bounded and there exists a constant C_p depending only on p such that

$$\|Y_t - Y_s\| \leq C_p(1 + \|\xi\|)K\|X\|_{p\text{-var};[s,t]} \exp\left(C_p K^p \|X\|_{p\text{-var};[s,t]}^p\right), \quad (31)$$

for all $0 \leq s \leq t \leq T$.

When the vector fields of the RDE (22) are linear, then the log-ODE (27) also becomes linear. Therefore, the log-ODE solution exists and is explicitly given as the exponential of the matrix F .

Theorem B.4 Consider the same linear RDE on $[0, T]$ as in Theorem B.3,

$$\begin{aligned} dY_t &= f(Y_t) dX_t, \\ Y_0 &= \xi. \end{aligned}$$

Then the log-ODE vector field F given by (28) is linear and the solution of the associated ODE (27) exists and satisfies

$$\|z_u^{s,t}\| \leq \|Y_s\| \exp\left(\sum_{m=1}^{\lfloor \gamma \rfloor} K^m \left\| \pi_m \left(\text{LogSig}_{s,t}^{\lfloor \gamma \rfloor}(X) \right) \right\| \right), \quad (32)$$

for $u \in [0, 1]$ and all $0 \leq s \leq t \leq T$.

Proof B.5 Since F is a linear vector field on \mathbb{R}^n , we can view it as an $n \times n$ matrix and so for $u \in [0, 1]$,

$$z_u^{s,t} = \exp(uF) z_0^{s,t},$$

where \exp denotes the matrix exponential. The result now follows by the standard estimate $\|\exp(F)\| \leq \exp(\|F\|)$.

Remark B.6 Due to the boundedness of linear RDEs (31) and log-ODEs (32), the arguments that established Theorem B.1 will hold in the linear setting as $\|f\|_{\text{Lip}(\gamma)}$ would be finite when defined on the domains that the solutions Y and z lie in.

Given the local error estimate (29) for the log-ODE method, we can now consider the approximation error that is exhibited by a log-ODE numerical solution to the RDE (22). Thankfully, the analysis required to derive such global error estimates was developed by Greg Gyurkó in his PhD thesis. Thus the following result is a straightforward application of Theorem 3.2.1 from Gyurkó (2008).

Theorem B.7 Let X , f and Y satisfy the assumptions given by Theorem B.1 and suppose that $\{0 = t_0 < t_1 < \dots < t_N = T\}$ is a partition of $[0, T]$ with $\max_k \|X\|_{p\text{-var};[t_k, t_{k+1}]}$ sufficiently small. We can construct a numerical solution $\{Y_k^{\log}\}_{0 \leq k \leq N}$ of (22) by setting $Y_0^{\log} := Y_0$ and for each $k \in \{0, 1, \dots, N-1\}$, defining Y_{k+1}^{\log} to be the solution at $u = 1$ of the following ODE:

$$\begin{aligned} \frac{dz^{t_k, t_{k+1}}}{du} &:= F(z^{t_k, t_{k+1}}), \\ z_0^{t_k, t_{k+1}} &:= Y_k^{\log}, \end{aligned} \quad (33)$$

where the vector field F is constructed from the log-signature of X over the interval $[t_k, t_{k+1}]$ according to (28). Then there exists a constant C depending only on p , γ and $\|f\|_{\text{Lip}(\gamma)}$ such that

$$\|Y_{t_k} - Y_k^{\log}\| \leq C \sum_{i=0}^{k-1} \|X\|_{p\text{-var};[t_i, t_{i+1}]}^{\gamma}, \quad (34)$$

for $0 \leq k \leq N$.

Remark B.8 The above error estimate also holds when the vector field f is linear (by Remark B.6)).

Since $\lfloor \gamma \rfloor$ is the truncation depth of the log-signatures used to construct each log-ODE vector field, we see that high convergence rates can be achieved through using more terms in each log-signature. It is also unsurprising that the error estimate (34) increases with the “roughness” of the control path. So just as in our experiments, we see that the performance of the log-ODE method can be improved by choosing an appropriate step size and depth of log-signature.

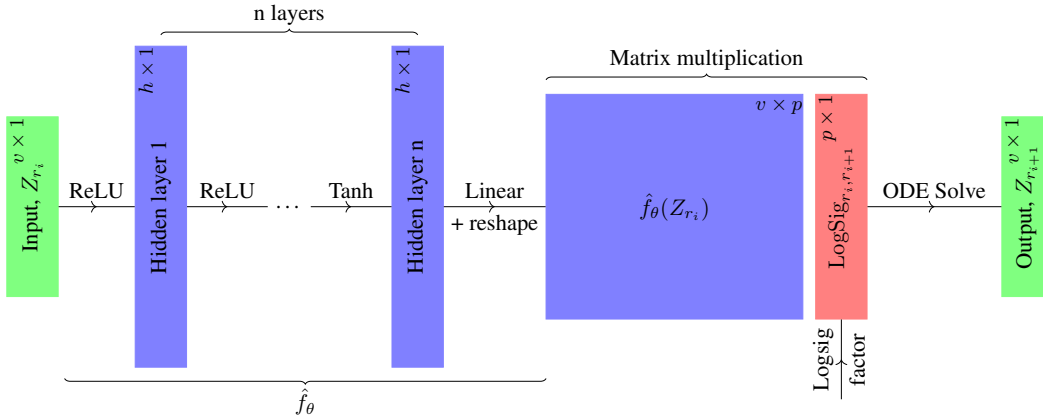


Figure 5: Overview of the hidden state update network structure. We give the dimensions at each layer in the top right hand corner of each box.

C EXPERIMENTAL DETAILS

Code The code to reproduce the experiments is available at [redacted; see supplementary material]

Data splits Each dataset was split into a training, validation, and testing dataset with relative sizes 70%/15%/15%.

Normalisation The training splits of each dataset were normalised to zero mean and unit variance. The statistics from the training set were then used to normalise the validation and testing datasets.

Architecture We give a graphical description of the architecture used for updating the Neural CDE hidden state in figure 5. The input is first run through a multilayer perceptron with n layers of size h , with n, h being hyperparameters. ReLU nonlinearities are used at each layer except the final one, where we instead use a tanh nonlinearity. The goal of this is to help prevent term blow-up over the long sequences.

Note that this is a small inconsistency between this work and the original model proposed in Kidger et al. (2020). Here, we applied the tanh function as the final hidden layer nonlinearity, whilst in the original paper the tanh nonlinearity is applied after the final linear map. Both methods are used to constrain the rate of change of the hidden state; we do not know of a reason to prefer one over the other.

Note that the final linear layer in the multilayer perceptron is reshaped to produce a matrix-valued output, of shape $v \times p$. (As \hat{f}_θ is matrix-valued.) A matrix-vector multiplication with the log-signature then produces the vector field for the ODE solver.

ODE Solver All problems used the ‘rk4’ solver as implemented by `torchdiffeq` (Chen, 2018) version 0.0.1.

Computing infrastructure All EigenWorms experiments were run on a computer equipped with three GeForce RTX 2080 Ti’s. All BIDMC experiments were run on a computer with two GeForce RTX 2080 Ti’s and two Quadro GP100’s.

Optimiser All experiments used the Adam optimiser. The learning rate was initialised at 0.032 divided by batch size. The batch size used was 1024 for EigenWorms and 512 for the BIDMC problems. If the validation loss failed to decrease after 15 epochs the learning rate was reduced by a factor of 10. If the validation loss did not decrease after 60 epochs, training was terminated and the model was rolled back to the point at which it achieved the lowest loss on the validation set.

Validation accuracy	Hidden dim	Num layers	Hidden hidden multiplier	Total params
33.3	16	2	3	5509
43.6	16	2	2	5509
56.4	16	2	1	4453
64.1	16	3	2	8869
38.5	16	3	3	8869
51.3	16	3	1	6517
82.1	16	4	2	12741
35.9	16	4	3	12741
53.8	16	4	1	8581
35.9	32	2	3	21253
74.4	32	2	2	21253
43.6	32	2	1	17093
53.8	32	3	3	34629
87.2	32	3	2	34629
64.1	32	3	1	25317
35.9	32	4	3	50053
71.8	32	4	1	33541
79.5	32	4	2	50053
41.0	64	2	3	83461
64.1	64	2	2	83461
48.7	64	3	3	136837
59.0	64	3	2	136837
51.3	64	2	1	66949
56.4	64	4	2	198405
64.1	64	4	3	198405
64.1	64	3	1	99781
51.3	64	4	1	132613

Table 3: Hyperparamter selection results for the EigenWorms dataset. The blue values denote the selected hyperparameters.

Hyperparameter selection Hyperparameters were selected to optimise the score of the NCDE₁ model on the validation set. For each dataset the search was performed with a step size that meant the total number of hidden state updates was equal to 500, as this represented a good balance between length and speed that allowed us to complete the search in a reasonable time-frame. In particular, this was short enough that we could train using the non-adjoint training method which helped to speed this section up. The hyperparameters that were considered were:

- Hidden dimension: [16, 32, 64] - The dimension of the hidden state Z_t .
- Number of layers: [2, 3, 4] - The number of hidden state layers.
- Hidden hidden multiplier: [1, 2, 3] - Multiplication factor for the hidden hidden state, this being the ‘Hidden layer k ’ in figure 5. The dimension of each of these ‘hidden hidden’ layers will be this value multiplied by ‘Hidden dimension’.

We ran each of these 27 total combinations for every dataset and the parameters that corresponded were used as the parameters when training over the full depth and step grid. The full results from the hyperparameter search are listed in tables (3, 4) with bolded values to show which values were eventually selected.

D EXPERIMENTAL RESULTS

Here we include the full breakdown of all experimental results. Tables 5 and 6 include all results from the EigenWorms and BIDMC datasets respectively.

Validation loss			Hidden dim	Num layers	Hidden hidden multiplier	Total params
RR	HR	SpO2				
1.72	6.10	2.07	16	2	1	2209
1.57	5.58	1.97	16	2	2	3265
1.55	6.10	1.33	16	2	3	3265
1.80	5.16	2.05	16	3	1	3249
1.61	5.22	1.62	16	3	2	5601
1.56	3.34	1.18	16	3	3	5601
1.57	3.86	1.97	16	4	1	4289
1.45	3.54	1.25	16	4	2	8449
1.54	3.93	1.09	16	4	3	8449
1.56	6.81	1.87	32	2	1	8513
1.42	3.11	1.43	32	2	2	12673
1.54	3.60	1.11	32	2	3	12673
1.54	3.52	1.57	32	3	1	12641
1.39	2.96	1.03	32	3	2	21953
1.47	2.95	1.05	32	3	3	21953
1.55	3.00	2.00	32	4	1	16769
1.38	3.20	1.07	32	4	2	33281
1.43	2.58	1.01	32	4	3	33281
1.51	3.21	1.10	64	2	1	33409
1.43	2.22	1.00	64	2	2	49921
1.51	3.34	0.94	64	2	3	49921
1.55	3.24	2.09	64	3	1	49857
1.32	2.53	0.88	64	3	2	86913
1.25	2.57	0.73	64	3	3	86913
1.43	5.78	1.43	64	4	1	66305
1.28	2.26	0.93	64	4	2	132097
1.32	2.46	1.15	64	4	3	132097

Table 4: Hyperparameter selection results for each problem of the BIDMC dataset. The bold values denote the selected hyperparameters for each vitals sign problem. Note that RR and SpO2 had the same parameters selected, hence why only two lines are given in bold.

Model	Step	Test Accuracy	Time (Hrs)	Memory (Mb)
NCDE ₁	1	62.4 \pm 12.1	22.0	176.5
	2	69.2 \pm 4.4	14.6	90.6
	4	66.7 \pm 11.8	5.5	46.6
	6	65.8 \pm 12.9	2.6	31.5
	8	64.1 \pm 13.3	3.1	24.3
	16	64.1 \pm 16.8	1.5	13.4
	32	64.1 \pm 14.3	0.5	8.0
	64	56.4 \pm 6.8	0.4	5.2
	128	48.7 \pm 2.6	0.1	3.9
	256	42.7 \pm 3.0	0.1	3.2
	512	44.4 \pm 5.3	0.0	2.9
	1024	41.9 \pm 14.6	0.0	2.7
	2048	38.5 \pm 5.1	0.0	2.6
NCDE ₂	2	76.1 \pm 13.2	9.8	354.3
	4	83.8 \pm 3.0	2.4	180.0
	6	76.9 \pm 6.8	2.0	82.2
	8	77.8 \pm 5.9	2.1	94.2
	16	78.6 \pm 3.9	1.3	50.2
	32	67.5 \pm 12.1	0.7	28.1
	64	73.5 \pm 7.8	0.4	17.2
	128	76.1 \pm 5.9	0.2	7.8
	256	72.6 \pm 12.1	0.1	8.9
	512	69.2 \pm 11.8	0.0	7.6
	1024	65.0 \pm 7.4	0.0	6.9
	2048	67.5 \pm 3.9	0.0	6.5
NCDE ₃	2	66.7 \pm 4.4	7.4	1766.2
	4	76.9 \pm 9.2	2.8	856.8
	6	70.9 \pm 1.5	1.4	606.1
	8	70.1 \pm 6.5	1.3	460.7
	16	73.5 \pm 3.0	1.4	243.7
	32	75.2 \pm 3.0	0.6	134.7
	64	74.4 \pm 11.8	0.3	81.0
	128	68.4 \pm 8.2	0.1	53.3
	256	60.7 \pm 8.2	0.1	40.2
	512	62.4 \pm 10.4	0.0	33.1
	1024	59.8 \pm 3.9	0.0	29.6
	2048	61.5 \pm 4.4	0.0	27.7

Table 5: Mean and standard deviation of test set accuracy (in %) over three repeats, as well as memory usage and training time, on the EigenWorms dataset for depths 1–3 and a small selection of step sizes. The bold values denote that the model was the top performer for that step size.

Depth	Step	L^2			Time (H)			Memory (Mb)
		RR	HR	SpO ₂	RR	HR	SpO ₂	
NCDE ₁	1	2.79 ± 0.04	9.82 ± 0.34	2.83 ± 0.27	23.8	22.1	28.1	56.5
	2	2.87 ± 0.03	11.69 ± 0.38	3.36 ± 0.2	19.3	9.6	8.8	32.6
	4	2.92 ± 0.08	11.15 ± 0.49	3.69 ± 0.06	5.3	5.7	3.2	20.2
	8	2.8 ± 0.06	10.72 ± 0.24	3.43 ± 0.17	3.0	2.6	4.8	14.3
	16	2.22 ± 0.07	7.98 ± 0.61	2.9 ± 0.11	1.7	1.4	1.8	11.8
	32	2.53 ± 0.23	12.23 ± 0.43	2.68 ± 0.12	1.9	0.9	2.2	9.8
	64	2.63 ± 0.11	12.02 ± 0.09	2.88 ± 0.06	0.2	0.3	0.4	9.1
	128	2.64 ± 0.18	11.98 ± 0.37	2.86 ± 0.04	0.2	0.2	0.3	8.7
	256	2.53 ± 0.04	12.29 ± 0.1	3.08 ± 0.1	0.1	0.1	0.1	8.3
	512	2.53 ± 0.03	12.22 ± 0.11	2.98 ± 0.04	0.1	0.0	0.1	8.4
	1024	2.67 ± 0.12	11.55 ± 0.03	2.91 ± 0.12	0.1	0.1	0.1	8.4
	2048	2.48 ± 0.03	12.03 ± 0.2	3.25 ± 0.01	0.0	0.1	0.0	8.2
NCDE ₂	2	2.91 ± 0.1	11.11 ± 0.23	3.89 ± 0.44	12.7	9.3	8.2	58.3
	4	2.92 ± 0.04	11.14 ± 0.2	4.23 ± 0.57	18.1	5.0	3.4	34.0
	8	2.63 ± 0.12	8.63 ± 0.24	2.88 ± 0.15	2.1	3.4	3.3	21.8
	16	1.8 ± 0.07	5.73 ± 0.45	1.98 ± 0.21	2.2	1.4	2.5	16.0
	32	1.9 ± 0.02	7.9 ± 1.0	1.69 ± 0.2	1.2	1.1	2.0	13.1
	64	1.89 ± 0.04	5.54 ± 0.45	2.04 ± 0.07	0.3	0.3	1.7	11.6
	128	1.86 ± 0.03	6.77 ± 0.42	1.95 ± 0.18	0.3	0.4	0.7	10.9
	256	1.86 ± 0.09	5.64 ± 0.19	2.1 ± 0.19	0.1	0.1	0.5	10.5
	512	1.81 ± 0.02	5.05 ± 0.23	2.17 ± 0.18	0.1	0.2	0.4	10.3
	1024	1.93 ± 0.11	6.0 ± 0.19	2.41 ± 0.07	0.1	0.1	0.2	10.2
	2048	2.03 ± 0.03	7.7 ± 1.46	2.55 ± 0.03	0.1	0.1	0.1	10.2
NCDE ₃	2	2.82 ± 0.08	11.01 ± 0.28	4.1 ± 0.72	8.8	9.4	6.9	125.2
	4	2.97 ± 0.23	10.13 ± 0.62	3.56 ± 0.44	3.2	4.1	2.6	71.6
	8	2.42 ± 0.19	7.67 ± 0.4	2.55 ± 0.13	2.9	3.2	3.1	43.3
	16	1.74 ± 0.05	4.11 ± 0.61	1.4 ± 0.06	1.4	1.4	6.5	29.1
	32	1.67 ± 0.01	4.5 ± 0.7	1.61 ± 0.05	1.3	1.8	7.3	20.5
	64	1.53 ± 0.08	3.05 ± 0.36	1.48 ± 0.14	0.4	1.9	3.3	17.9
	128	1.51 ± 0.08	2.97 ± 0.45	1.37 ± 0.22	0.5	1.7	1.7	17.3
	256	1.51 ± 0.06	3.4 ± 0.74	1.47 ± 0.07	0.3	0.7	0.6	16.6
	512	1.49 ± 0.08	3.46 ± 0.13	1.29 ± 0.15	0.3	0.4	0.4	15.4
	1024	1.83 ± 0.33	5.58 ± 2.5	1.72 ± 0.31	0.2	0.1	0.1	15.7
	2048	2.31 ± 0.27	9.77 ± 1.53	2.45 ± 0.18	0.1	0.1	0.1	15.6

Table 6: Mean and standard deviation of the L^2 losses on the test set for each of the vitals signs prediction tasks (RR, HR, SpO₂) on the BIDMC dataset, across three repeats. Only mean times are shown for space. The memory usage is given as the mean over all three of the tasks as it was approximately the same for any task for a given depth and step. The bold values denote the algorithm with the lowest test set loss for a fixed step size for each task.