

# Framing RNN as a kernel method: A neural ODE approach

Adeline Fermanian, Pierre Marion, Jean-Philippe Vert, Gérard Biau

# ► To cite this version:

Adeline Fermanian, Pierre Marion, Jean-Philippe Vert, Gérard Biau. Framing RNN as a kernel method: A neural ODE approach. Thirty-fifth Conference on Neural Information Processing Systems, Dec 2021, Virtual-only, United States. hal-03943120

# HAL Id: hal-03943120 https://hal.science/hal-03943120

Submitted on 10 Feb 2023  $\,$ 

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers. L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Framing RNN as a kernel method: A neural ODE approach

Adeline Fermanian<sup>1\*</sup>

Pierre Marion<sup>1\*</sup> Jean

Jean-Philippe Vert<sup>2</sup>

Gérard Biau<sup>1</sup>

<sup>1</sup> Sorbonne Université, CNRS, Laboratoire de Probabilités, Statistique et Modélisation, LPSM, F-75005 Paris, France {adeline.fermanian, pierre.marion, gerard.biau}@sorbonne-universite.fr <sup>2</sup> Google Research, Brain team, Paris, France jpvert@google.com

# Abstract

Building on the interpretation of a recurrent neural network (RNN) as a continuoustime neural differential equation, we show, under appropriate conditions, that the solution of a RNN can be viewed as a linear function of a specific feature set of the input sequence, known as the signature. This connection allows us to frame a RNN as a kernel method in a suitable reproducing kernel Hilbert space. As a consequence, we obtain theoretical guarantees on generalization and stability for a large class of recurrent networks. Our results are illustrated on simulated datasets.

# 1 Introduction

Recurrent neural networks (RNN) are among the most successful methods for modeling sequential data. They have achieved state-of-the-art results in difficult problems such as natural language processing (e.g., Mikolov et al., 2010; Collobert et al., 2011) or speech recognition (e.g., Hinton et al., 2012; Graves et al., 2013). This class of neural networks has a natural interpretation in terms of (discretization of) ordinary differential equations (ODE), which casts them in the field of neural ODE (Chen et al., 2018). This observation has led to the development of continuous-depth models for handling irregularly-sampled time-series data, including the ODE-RNN model (Rubanova et al., 2019), GRU-ODE-Bayes (De Brouwer et al., 2019), or neural CDE models (Kidger et al., 2020; Morrill et al., 2020a). In addition, the time-continuous interpretation of RNN allows to leverage the rich theory of differential equations to develop new recurrent architectures (Chang et al., 2019; Herrera et al., 2020; Erichson et al., 2021), which are better at learning long-term dependencies.

On the other hand, the development of kernel methods for deep learning offers theoretical insights on the functions learned by the networks (Cho and Saul, 2009; Belkin et al., 2018; Jacot et al., 2018). Here, the general principle consists in defining a reproducing kernel Hilbert space (RKHS)—that is, a function class  $\mathcal{H}$ —, which is rich enough to describe the architectures of networks. A good example is the construction of Bietti and Mairal (2017, 2019), who exhibit an RKHS for convolutional neural networks. This kernel perspective has several advantages. First, by separating the representation of the data from the learning process, it allows to study invariances of the representations learned by the network. Next, by reducing the learning problem to a linear one in  $\mathcal{H}$ , generalization bounds can be more easily obtained. Finally, the Hilbert structure of  $\mathcal{H}$  provides a natural metric on neural networks, which can be used for example for regularization (Bietti et al., 2019).

35th Conference on Neural Information Processing Systems (NeurIPS 2021).

<sup>\*</sup>Equal contribution

**Contributions.** By taking advantage of the neural ODE paradigm for RNN, we show that RNN are, in the continuous-time limit, linear predictors over a specific space associated with the signature of the input sequence (Levin et al., 2013). The signature transform, first defined by Chen (1958) and central in rough path theory (Lyons et al., 2007; Friz and Victoir, 2010), summarizes sequential inputs by a graded feature set of their iterated integrals. Its natural environment is a tensor space that can be endowed with an RKHS structure (Király and Oberhauser, 2019). We exhibit general conditions under which classical recurrent architectures such as feedforward RNN, Gated Recurrent Units (GRU, Cho et al., 2014), or Long Short-Term Memory networks (LSTM, Hochreiter and Schmidhuber, 1997), can be framed as a kernel method in this RKHS. This enables us to provide generalization bounds for RNN as well as stability guarantees via regularization. The theory is illustrated with some experimental results.

**Related works.** The neural ODE paradigm was first formulated by Chen et al. (2018) for residual neural networks. It was then extended to RNN in several articles, with a focus on handling irregularly sampled data (Rubanova et al., 2019; Kidger et al., 2020) and learning long-term dependencies (Chang et al., 2019). The signature transform has recently received the attention of the machine learning community (Levin et al., 2013; Kidger et al., 2019; Liao et al., 2019; Toth and Oberhauser, 2020; Fermanian, 2021) and, combined with deep neural networks, has achieved state-of-the-art performance for several applications (Yang et al., 2016, 2017; Perez Arribas, 2018; Wang et al., 2019; Morrill et al., 2020b). Király and Oberhauser (2019) use the signature transform to define kernels for sequential data and develop fast computational methods. The connection between continuous-time RNN and signatures has been pointed out by Lim (2021) for a specific model of stochastic RNN. Deriving generalization bounds for RNN is an active research area (Zhang et al., 2018; Akpinar et al., 2019; Tu et al., 2019). By leveraging the theory of differential equations, our approach encompasses a large class of RNN models, ranging from feedforward RNN to LSTM. This is in contrast with most existing generalization bounds, which are architecture-dependent. Close to our point of view is the work of Bietti and Mairal (2017) for convolutional neural networks.

**Mathematical context.** We place ourselves in a supervised learning setting. The input data is a sample of *n* i.i.d. vector-valued sequences  $\{\mathbf{x}^{(1)}, \ldots, \mathbf{x}^{(n)}\}$ , where  $\mathbf{x}^{(i)} = (x_1^{(i)}, \ldots, x_T^{(i)}) \in (\mathbb{R}^d)^T$ ,  $T \ge 1$ . The outputs of the learning problem can be either labels (classification setting) or sequences (sequence-to-sequence setting). Even if we only observe discrete sequences, each  $\mathbf{x}^{(i)}$  is mathematically considered as a regular discretization of a continuous-time process  $X^{(i)} \in BV^c([0,1], \mathbb{R}^d)$ , where  $BV^c([0,1], \mathbb{R}^d)$  is the space of continuous functions from [0,1] to  $\mathbb{R}^d$  of finite total variation. Informally, the total variation of a process corresponds to its length. Formally, for any  $[s,t] \subset [0,1]$ , the total variation of a process  $X \in BV^c([0,1], \mathbb{R}^d)$  on [s,t] is defined by

$$||X||_{TV;[s,t]} = \sup_{(t_0,\dots,t_k)\in D_{s,t}} \sum_{j=1}^k ||X_{t_j} - X_{t_{j-1}}||,$$

where  $D_{s,t}$  denotes the set of all finite partitions of [s, t] and  $\|\cdot\|$  the Euclidean norm. We therefore have that  $x_j^{(i)} = X_{j/T}^{(i)}$ ,  $1 \le j \le T$ , where  $X_t^{(i)} := X^{(i)}(t)$ . We make two assumptions on the processes  $X^{(i)}$ . First, they all begin at zero, and second, their lengths are bounded by  $L \in (0, 1)$ . These assumptions are not too restrictive, since they amount to data translation and normalization, common in practice. Accordingly, we denote by  $\mathscr{X}$  the subset of  $BV^c([0, 1], \mathbb{R}^d)$  defined by

$$\mathscr{X} = \{ X \in BV^c([0,1], \mathbb{R}^d) \mid X_0 = 0 \text{ and } \|X\|_{TV;[0,1]} \le L \}$$

and assume therefore that  $X^{(1)}, \ldots, X^{(n)}$  are i.i.d. according to some  $X \in \mathscr{X}$ . The norm on all spaces  $\mathbb{R}^m, m \ge 1$ , is always the Euclidean one. Observe that assuming that  $X \in \mathscr{X}$  implies that, for any  $t \in [0, 1], \|X_t\| = \|X_t - X_0\| \le \|X\|_{TV;[0,1]} \le L$ .

**Recurrent neural networks.** Classical RNN are defined by a sequence of hidden states  $h_1, \ldots, h_T \in \mathbb{R}^e$ , where, for  $\mathbf{x} = (x_1, \ldots, x_T)$  a generic data sample,

$$h_0 = 0$$
 and  $h_{j+1} = f(h_j, x_{j+1})$  for  $0 \le j \le T - 1$ .

At each time step  $1 \le j \le T$ , the output of the network is  $z_j = \psi(h_j)$ , where  $\psi$  is a linear function. In the present article, we rather consider the following residual version, which is a natural adaptation of classical RNN in the neural ODE framework (see, e.g., Yue et al., 2018):

$$h_0 = 0$$
 and  $h_{j+1} = h_j + \frac{1}{T}f(h_j, x_{j+1})$  for  $0 \le j \le T - 1$ . (1)

The simplest choice for the function f is the feedforward model, say  $f_{\rm RNN}$ , defined by

$$f_{\text{RNN}}(h, x) = \sigma(Uh + Vx + b), \tag{2}$$

where  $\sigma$  is an activation function,  $U \in \mathbb{R}^{e \times e}$  and  $V \in \mathbb{R}^{e \times d}$  are weight matrices, and  $b \in \mathbb{R}^{e}$  is the bias. The function  $f_{\rm RNN}$ , equipped with a smooth activation  $\sigma$  (such as the logistic or hyperbolic tangent functions), will be our leading example throughout the paper. However, the GRU and LSTM models can also be rewritten under the form (1), as shown in Appendix A.1. Thus, model (1) is flexible enough to encompass most recurrent networks used in practice.

**Overview.** Section 2 is devoted to framing RNN as linear functions in a suitable RKHS. We start by embedding iteration (1) into a continuous-time model, which takes the form of a controlled differential equation (CDE). This allows, after introducing the signature transform, to define the appropriate RKHS, and, in turn, to show that model (1) boils down, in the continuous-time limit, to a linear problem on the signature. This framework is used in Section 3 to derive generalization bounds and stability guarantees. We provide some experiments in Section 4 before discussing our results in Section 5. All proofs are postponed to the supplementary material.

#### 2 Framing RNN as a kernel method

**Roadmap.** First, we quantify the difference between the discrete recurrent network (1) and its continuous-time counterpart (Proposition 1). Then, we rewrite the corresponding ODE as a CDE (Proposition 2). Under appropriate conditions, Proposition 4 shows that the solution of this equation is a linear function of the signature of the driving process. Importantly, these assumptions are valid for a feedforward RNN, as stated by Proposition 5. We conclude in Theorem 1.

#### 2.1 From discrete to continuous time

Recall that  $h_0, \ldots, h_T$  denote the hidden states of the RNN (1), and let  $H : [0,1] \to \mathbb{R}^e$  be the solution of the ODE

$$dH_t = f(H_t, X_t)dt, \quad H_0 = h_0.$$
 (3)

By bounding the difference between  $H_{j/T}$  and  $h_j$ , the following proposition shows how to pass from discrete to continuous time, provided f satisfies the following assumption:

 $(A_1)$  The function f is Lipschitz continuous in h and x, with Lipschitz constants  $K_h$  and  $K_x$ . We let  $K_f = \max(K_h, K_r)$ .

**Proposition 1.** Assume that  $(A_1)$  is verified. Then there exists a unique solution H to (3) and, for any  $0 \le j \le T$ ,

$$\|H_{j/T} - h_j\| \le \frac{c_1}{T},$$

 $\|H_{j/T} - h_j\| \leq \frac{1}{T},$ where  $c_1 = K_f e^{K_f} \left( L + \sup_{\|h\| \leq M, \|x\| \leq L} \|f(h, x)\| e^{K_f} \right)$  and  $M = \sup_{\|x\| \leq L} \|f(h_0, x)\| e^{K_f}.$  Moreover, for any  $t \in [0, 1]$ ,  $||H_t|| \le M$ .

Then, following Kidger et al. (2020), we show that the ODE (3) can be rewritten under the form of a CDE. At the cost of increasing the dimension of the hidden state from e to e + d, this allows us to reframe model (3) as a linear model in dX, in the sense that X has been moved 'outside' of f.

**Proposition 2.** Assume that  $(A_1)$  is verified. Let  $H : [0,1] \to \mathbb{R}^e$  be the solution of (3), and let  $\bar{X} : [0,1] \to \mathbb{R}^{d+1}$  be the time-augmented process  $\bar{X}_t = (X_t^{\top}, \frac{1-L}{2}t)^{\top}$ . Then there exists a tensor field  $\mathbf{F}: \mathbb{R}^{\bar{e}} \to \mathbb{R}^{\bar{e} \times \bar{d}}, \bar{e} = e + d, \bar{d} = d + 1$ , such that if  $\bar{H}: [0,1] \to \mathbb{R}^{\bar{e}}$  is the solution of the CDE

$$d\bar{H}_t = \mathbf{F}(\bar{H}_t)d\bar{X}_t, \quad \bar{H}_0 = (H_0^+, X_0^+)^+, \tag{4}$$

then its first e coordinates are equal to H.

Equation (4) can be better understood by the following equivalent integral equation:

$$\bar{H}_t = \bar{H}_0 + \int_0^t \mathbf{F}(\bar{H}_u) d\bar{X}_u,$$

where the integral should be understood as Riemann-Stieljes integral (Friz and Victoir, 2010, Section I.2). Thus, the output of the RNN can be approximated by the solution of the CDE (4), and, according to Proposition 1, the approximation error is  $\mathcal{O}(1/T)$ .

**Example 1.** Consider  $f_{\text{RNN}}$  as in (2). If  $\sigma$  is Lipschitz continuous with constant  $K_{\sigma}$ , then, for any  $h_1, h_2 \in \mathbb{R}^e$ ,  $x_1, x_2 \in \mathbb{R}^d$ ,

$$\|f_{\text{RNN}}(h_1, x_1) - f_{\text{RNN}}(h_2, x_1)\| = \|\sigma(Uh_1 + Vx_1 + b) - \sigma(Uh_2 + Vx_1 + b)\| \le K_{\sigma} \|U\|_{\text{op}} \|h_1 - h_2\|,$$

where  $\|\cdot\|_{op}$  denotes the operator norm—see Appendix A.3. Similarly,  $\|f(h_1, x_1) - f(h_1, x_2)\| \le K_{\sigma} \|V\|_{op} \|x_1 - x_2\|$ . Thus, assumption  $(A_1)$  is satisfied. The tensor field  $\mathbf{F}_{RNN}$  of Proposition 2 corresponding to this network is defined for any  $\bar{h} \in \mathbb{R}^{\bar{e}}$  by

$$\mathbf{F}_{\mathsf{RNN}}(\bar{h}) = \begin{pmatrix} 0_{e \times d} & \frac{2}{1-L}\sigma(W\bar{h}+b) \\ I_{d \times d} & 0_{d \times 1} \end{pmatrix}, \quad \text{where} \quad W = (U \quad V) \in \mathbb{R}^{e \times \bar{e}}.$$
 (5)

# 2.2 The signature

An essential ingredient towards our construction is the signature of a continuous-time process, which we briefly present here. We refer to Chevyrev and Kormilitzin (2016) for a gentle introduction and to Lyons et al. (2007); Levin et al. (2013) for details.

**Tensor Hilbert spaces.** We denote by  $(\mathbb{R}^d)^{\otimes k}$  the *k*th tensor power of  $\mathbb{R}^d$  with itself, which is a Hilbert space of dimension  $d^k$ . The key space to define the signature and, in turn, our RKHS, consists in infinite square-summable sequences of tensors of increasing order:

$$\mathscr{T} = \left\{ a = (a_0, \dots, a_k, \dots) \, \middle| \, a_k \in (\mathbb{R}^d)^{\otimes k}, \, \sum_{k=0}^\infty \|a_k\|_{(\mathbb{R}^d)^{\otimes k}}^2 < \infty \right\}.$$
(6)

Endowed with the scalar product  $\langle a, b \rangle_{\mathscr{T}} := \sum_{k=0}^{\infty} \langle a_k, b_k \rangle_{(\mathbb{R}^d)^{\otimes k}}$ ,  $\mathscr{T}$  is a Hilbert space, as shown in Appendix A.4.

**Definition 1.** Let  $X \in BV^c([0,1], \mathbb{R}^d)$ . For any  $t \in [0,1]$ , the signature of X on [0,t] is defined by  $S_{[0,t]}(X) = (1, \mathbb{X}^1_{[0,t]}, \dots, \mathbb{X}^k_{[0,t]}, \dots)$ , where, for each  $k \ge 1$ ,

$$\mathbb{X}_{[0,t]}^k = k! \int_{0 \le u_1 < \dots < u_k \le t} dX_{u_1} \otimes \dots \otimes dX_{u_k} \in (\mathbb{R}^d)^{\otimes k}.$$

Although this definition is technical, the signature should simply be thought of as a feature map that embeds a bounded variation process into an infinite-dimensional tensor space. The signature has several good properties that make it a relevant tool for machine learning (e.g., Levin et al., 2013; Chevyrev and Kormilitzin, 2016; Fermanian, 2021). In particular, under certain assumptions, S(X) characterizes X up to translations and reparameterizations, and has good approximation properties. We also highlight that fast libraries exist for computing the signature (Reizenstein and Graham, 2020; Kidger and Lyons, 2021).

The expert reader is warned that this definition differs from the usual one by the normalization of  $\mathbb{X}_{[0,t]}^k$ by k!, which is more adapted to our context. In the sequel, for any index  $(i_1, \ldots, i_k) \subset \{1, \ldots, d\}^k$ ,  $S_{[0,t]}^{(i_1,\ldots,i_k)}(X)$  denotes the term associated with the coordinates  $(i_1,\ldots,i_k)$  of  $\mathbb{X}_{[0,t]}^k$ . When the signature is taken on the whole interval [0, 1], we simply write S(X),  $S^{(i_1,\ldots,i_k)}(X)$ , and  $\mathbb{X}^k$ .

**Example 2.** Let X be the d-dimensional linear path defined by  $X_t = (a_1 + b_1 t, \dots, a_d + b_d t)^\top$ ,  $a_i, b_i \in \mathbb{R}$ . Then  $S^{(i_1,\dots,i_k)}(X) = b_{i_1} \cdots b_{i_k}$  and  $\mathbb{X}^k = b^{\otimes k}$ .

The next proposition, which ensures that  $S_{[0,t]}(\bar{X}) \in \mathscr{T}$ , is an important step.

**Proposition 3.** Let  $X \in \mathscr{X}$  and  $\bar{X}_t = (X_t^{\top}, \frac{1-L}{2}t)^{\top}$  as in Proposition 2. Then, for any  $t \in [0, 1]$ ,  $\|S_{[0,t]}(\bar{X})\|_{\mathscr{T}} \leq 2(1-L)^{-1}$ .

**The signature kernel.** By taking advantage of the structure of Hilbert space of  $\mathscr{T}$ , it is natural to introduce the following kernel:

$$\begin{split} K: \mathscr{X} \times \mathscr{X} \to \mathbb{R} \\ (X,Y) \mapsto \langle S(\bar{X}), S(\bar{Y}) \rangle_{\mathscr{T}} \end{split}$$

which is well defined according to Proposition 3. We refer to Király and Oberhauser (2019) for a general presentation of kernel methods with signatures and to Cass et al. (2020) for a kernel trick. The RKHS associated with K is the space of functions

$$\mathscr{H} = \left\{ \xi_{\alpha} : \mathscr{X} \to \mathbb{R} \, | \, \xi_{\alpha}(X) = \langle \alpha, S(X) \rangle_{\mathscr{T}}, \alpha \in \mathscr{T} \right\},\tag{7}$$

with scalar product  $\langle \xi_{\alpha}, \xi_{\beta} \rangle_{\mathscr{H}} = \langle \alpha, \beta \rangle_{\mathscr{T}}$  (see, e.g., Schölkopf and Smola, 2002).

#### 2.3 From the CDE to the signature kernel

An important property of signatures is that the solution of the CDE (4) can be written, under certain assumptions, as a linear function of the signature of the driving process X. This operation can be thought of as a Taylor expansion for CDE. More precisely, let us rewrite (4) as

$$dH_t = \mathbf{F}(H_t)dX_t = \sum_{i=1}^d F^i(H_t)dX_t^i,$$
(8)

where  $X_t = (X_t^1, \ldots, X_t^d)^\top$ ,  $\mathbf{F} : \mathbb{R}^e \to \mathbb{R}^{e \times d}$ , and  $F^i : \mathbb{R}^e \to \mathbb{R}^e$  are the columns of  $\mathbf{F}$ —to avoid heavy notation, we momentarily write e, d, H, and X instead of  $\bar{e}, \bar{d}, \bar{H}$ , and  $\bar{X}$ . Throughout, the bold notation is used to distinguish tensor fields and vector fields. We recall that a vector field  $F : \mathbb{R}^e \to \mathbb{R}^e$  or a tensor field  $\mathbf{F} : \mathbb{R}^e \to \mathbb{R}^{e \times d}$  are said to be smooth if each of their coordinates is  $\mathscr{C}^{\infty}$ .

**Definition 2.** Let  $F, G : \mathbb{R}^e \to \mathbb{R}^e$  be smooth vector fields and denote by  $J(\cdot)$  the Jacobian matrix. Their differential product is the smooth vector field  $F \star G : \mathbb{R}^e \to \mathbb{R}^e$  defined, for any  $h \in \mathbb{R}^e$ , by

$$(F \star G)(h) = \sum_{j=1}^{c} \frac{\partial G}{\partial h_j}(h) F_j(h) = J(G)(h)F(h).$$

In differential geometry,  $F \star G$  is simply denoted by FG. Since the  $\star$  operation is not associative, we take the convention that it is evaluated from right to left, i.e.,  $F^1 \star F^2 \star F^3 := F^1 \star (F^2 \star F^3)$ .

**Taylor expansion.** Let H be the solution of (8), where  $\mathbf{F}$  is assumed to be smooth. We now show that H can be written as a linear function of the signature of X, which is the crucial step to embed the RNN in the RKHS  $\mathcal{H}$ . The step-N Taylor expansion of H (Friz and Victoir, 2008) is defined by

$$H_t^N = H_0 + \sum_{k=1}^N \frac{1}{k!} \sum_{1 \le i_1, \dots, i_k \le d} S_{[0,t]}^{(i_1, \dots, i_k)}(X) F^{i_1} \star \dots \star F^{i_k}(H_0).$$

Throughout, we let

$$\Lambda_k(\mathbf{F}) = \sup_{\|h\| \le M, 1 \le i_1, \dots, i_k \le d} \|F^{i_1} \star \dots \star F^{i_k}(h)\|.$$

**Example 3.** Let  $\mathbf{F} = \mathbf{F}_{\text{RNN}}$  defined by (5) with an identity activation. Then, for any  $\bar{h} \in \mathbb{R}^{\bar{e}}$ ,  $1 \leq i \leq d+1$ ,  $F_{\text{RNN}}^{i}(\bar{h}) = W_i \bar{h} + b_i$ , where  $b_i$  is the (i + d)th vector of the canonical basis of  $\mathbb{R}^{\bar{e}}$ , and

$$W_i = 0_{\bar{e} \times \bar{e}}, \quad W_{d+1} = \begin{pmatrix} \frac{2}{1-L}W\\ 0_{d \times \bar{e}} \end{pmatrix}, \quad and \quad b_{d+1} = \begin{pmatrix} \frac{2}{1-L}b\\ 0_d \end{pmatrix}$$

The vector fields  $F_{\text{RNN}}^i$  are then affine,  $J(F_{\text{RNN}}^i) = W_i$ , and the iterated star products have a simple expression: for any  $1 \le i_1, \ldots, i_k \le d$ ,  $F_{\text{RNN}}^{i_1} \star \cdots \star F_{\text{RNN}}^{i_k}(\bar{h}) = W_{i_k} \cdots W_{i_2}(W_{i_1}\bar{h} + b_{i_1})$ .

The next proposition shows that the step-N Taylor expansion  $H^N$  is a good approximation of H. **Proposition 4.** Assume that the tensor field **F** is smooth. Then, for any  $t \in [0, 1]$ ,

$$||H_t - H_t^N|| \le \frac{d^{N+1}}{(N+1)!} \Lambda_{N+1}(\mathbf{F}).$$
 (9)

Thus, provided that  $\Lambda_N(\mathbf{F})$  is not too large, the right-hand side of (9) converges to zero, hence

$$H_t = H_0 + \sum_{k=1}^{\infty} \frac{1}{k!} \sum_{1 \le i_1, \dots, i_k \le d} S^{(i_1, \dots, i_k)}_{[0,t]}(X) F^{i_1} \star \dots \star F^{i_k}(H_0).$$
(10)

We conclude from the above representation that the solution H of (8) is in fact a linear function of the signature of X. A natural concern is to know whether the upper bound of Proposition 4 vanishes with N for standard architectures. This property is encapsulated in the following more general assumption:

$$(A_2)$$
 The tensor field **F** is smooth and  $\sum_{k=0}^{\infty} \left(\frac{d^k}{k!}\Lambda_k(\mathbf{F})\right)^2 < \infty$ 

Clearly, if  $(A_2)$  is verified, then the right-hand side of (9) converges to 0. The next proposition states formally the conditions under which  $(A_2)$  is verified for  $\mathbf{F}_{\text{RNN}}$ . It is further illustrated in Figure 1, which shows that the convergence is fast with two common activation functions. We let  $\|\sigma\|_{\infty} = \sup_{\|h\| \le M, \|x\| \le L} \|\sigma(Uh + Vx + b)\|$  and  $\|\sigma^{(k)}\|_{\infty} = \sup_{\|h\| \le M, \|x\| \le L} \|\sigma^{(k)}(Uh + Vx + b)\|$ , where  $\sigma^{(k)}$  is the derivative of order k of  $\sigma$ .

**Proposition 5.** Let  $\mathbf{F}_{RNN}$  be defined by (5). If  $\sigma$  is the identity function, then  $(A_2)$  is satisfied. In the general case,  $(A_2)$  holds if  $\sigma$  is smooth and there exists a > 0 such that, for any  $k \ge 0$ ,

$$\|\sigma^{(k)}\|_{\infty} \le a^{k+1}k!$$
 and  $\|W\|_F < \frac{1-L}{8a^2d},$  (11)

where  $\|\cdot\|_F$  is the Frobenius norm. Moreover,  $\Lambda_N(\mathbf{F}_{RNN}) \leq \sqrt{2}a \left(\frac{8a^2 \|W\|_F}{1-L}\right)^{N-1} N!$ .

The proof of Proposition 5, based on the manipulation of higher-order derivatives of tensor fields, is highly non-trivial. We highlight that the conditions on  $\sigma$  are mild and verified for common smooth activations. For example, they are verified for the logistic function (with a = 2) and for the hyperbolic tangent function (with a = 4)—see Appendix A.5. The second inequality of (11) puts a constraint on the norm of the weights, and can be regarded as a radius of convergence for the Taylor expansion.

**Putting everything together.** We now have all the elements at hand to embed the RNN into the RKHS  $\mathscr{H}$ . To fix the idea, we assume in this paragraph that we are in a  $\pm 1$  classification setting. In other words, given an input sequence  $\mathbf{x}$ , we are interested in the final output  $z_T = \psi(h_T) \in \mathbb{R}$ , where  $h_T$  is the solution of (1). The predicted class is  $2 \cdot \mathbf{1}(z_T > 0) - 1$ .

By Propositions 1 and 2,  $z_T$  is approximated by the first e coordinates of the solution of the CDE (4), which outputs a  $\mathbb{R}^{e+d}$ -valued process  $\overline{H}$ . According to Proposition 4,  $\overline{H}$  is a linear function of the signature of the time-augmented process  $\overline{X}$ . Thus, on top of  $\overline{H}$ , it remains to successively apply the projection Proj on the e first coordinates followed by the linear function  $\psi$  to obtain an element of the RKHS  $\mathscr{H}$ . This mechanism is summarized in the following theorem.

**Theorem 1.** Assume that  $(A_1)$  and  $(A_2)$  are verified. Then there exists a function  $\xi_{\alpha} \in \mathcal{H}$  such that

$$|z_T - \xi_{\alpha}(X)| \le \|\psi\|_{\text{op}} \frac{c_1}{T},$$
(12)

where  $\xi_{\alpha}(X) = \langle \alpha, S(\bar{X}) \rangle_{\mathscr{T}}$  and  $\bar{X}_t = (X_t^{\top}, \frac{1-L}{2}t)^{\top}$ . We have  $\alpha = (\alpha_k)_{k=0}^{\infty}$ , where each  $\alpha_k \in (\mathbb{R}^d)^{\otimes k}$  is defined by

$$\alpha_k^{(i_1,\ldots,i_k)} = \frac{1}{k!} \psi \circ \operatorname{Proj} \left( F^{i_1} \star \cdots \star F^{i_k}(\bar{H}_0) \right).$$

Moreover,  $\|\alpha\|_{\mathscr{T}}^2 \leq \|\psi\|_{\text{op}}^2 \sum_{k=0}^{\infty} \left(\frac{d^k}{k!} \Lambda_k(\mathbf{F})\right)^2$ .

We conclude that in the continuous-time limit, the output of the network can be interpreted as a scalar product between the signature of the (time-augmented) process  $\bar{X}$  and an element of  $\mathscr{T}$ . This interpretation is important for at least two reasons: (i) it facilitates the analysis of generalization of RNN by leveraging the theory of kernel methods, and (ii) it provides new insights on regularization strategies to make RNN more robust. These points will be explored in the next section. Finally, we stress that the approach works for a large class of RNN, such as GRU and LSTM. The derivation of conditions ( $A_1$ ) and ( $A_2$ ) beyond the feedforward RNN is left for future work.

# **3** Generalization and regularization

#### 3.1 Generalization bounds

**Learning procedure.** A first consequence of framing a RNN as a kernel method is that it gives natural generalization bounds under mild assumptions. In the learning setup, we are given an i.i.d. sample  $\mathscr{D}_n$  of n random pairs of observations  $(\mathbf{x}^{(i)}, \mathbf{y}^{(i)}) \in (\mathbb{R}^d)^T \times \mathscr{Y}$ , where  $\mathbf{x}^{(i)} = (x_1^{(i)}, \ldots, x_T^{(i)})$ . We distinguish the binary classification problem, where  $\mathscr{Y} = \{-1, 1\}$ , from the sequential prediction problem, where  $\mathscr{Y} = (\mathbb{R}^p)^T$  and  $\mathbf{y}^{(i)} = (y_1^{(i)}, \ldots, y_T^{(i)})$ . The RNN is assumed to be parameterized by  $\theta \in \Theta \subset \mathbb{R}^q$ , where  $\Theta$  is a compact set. To clarify the notation, we use a  $\theta$  subscript whenever a quantity depends on  $\theta$  (e.g.,  $f_\theta$  for f, etc.). In line with Section 2, it is assumed that the tensor field  $\mathbf{F}_{\theta}$  associated with  $f_{\theta}$  satisfies  $(A_1)$  and  $(A_2)$ , keeping in mind that Proposition 5 guarantees that these requirements are fulfilled by a feedforward recurrent network with a smooth activation function.

Let  $g_{\theta} : (\mathbb{R}^d)^T \to \mathscr{Y}$  denote the output of the recurrent network. The parameter  $\theta$  is fitted by empirical risk minimization using a loss function  $\ell : \mathscr{Y} \times \mathscr{Y} \to \mathbb{R}^+$ . The theoretical and empirical risks are respectively defined, for any  $\theta \in \Theta$ , by

$$\mathscr{R}( heta) = \mathbb{E}[\ell(\mathbf{y}, g_{\theta}(\mathbf{x}))] \quad ext{and} \quad \widehat{\mathscr{R}}_n( heta) = rac{1}{n} \sum_{i=1}^n \ellig(\mathbf{y}^{(i)}, g_{ heta}(\mathbf{x}^{(i)})ig),$$

where the expectation  $\mathbb{E}$  is evaluated with respect to the distribution of the generic random pair  $(\mathbf{x}, \mathbf{y})$ . We let  $\hat{\theta}_n \in \operatorname{argmin}_{\theta \in \Theta} \widehat{\mathscr{R}}_n(\theta)$  and aim at upper bounding  $\mathbb{P}(\mathbf{y} \neq g_{\widehat{\theta}_n}(\mathbf{x}))$  in the classification regime (Theorem 2) and  $\mathscr{R}(\widehat{\theta}_n)$  in the sequential regime (Theorem 3). To reach this goal, our strategy is to approximate the RNN by its continuous version and then use the RKHS machinery of Section 2.

**Binary classification.** In this context, the network outputs a real number  $g_{\theta}(\mathbf{x}) = \psi(h_T) \in \mathbb{R}$ and the predicted class is  $2 \cdot \mathbf{1}(g_{\theta}(\mathbf{x}) > 0) - 1$ . The loss  $\ell : \mathbb{R} \times \mathbb{R} \to \mathbb{R}^+$  is assumed to satisfy the assumptions of Bartlett and Mendelson (2002, Theorem 7), that is, for any  $y \in \{-1, 1\}$ ,  $\ell(\mathbf{y}, g_{\theta}(\mathbf{x})) = \phi(\mathbf{y}g_{\theta}(\mathbf{x}))$ , where  $\phi(u) \geq \mathbf{1}(u \leq 0)$ , and  $\phi$  is Lipschitz-continuous with constant  $K_{\ell}$ . For example, the logistic loss satisfies such assumptions. We let  $\xi_{\alpha_{\theta}} \in \mathscr{H}$  be the function of Theorem 1 that approximates the RNN with parameter  $\theta$ . Thus,  $z_T \approx \xi_{\alpha_{\theta}}(\bar{X}) = \langle \alpha_{\theta}, S(\bar{X}) \rangle_{\mathscr{T}}$ , up to a  $\mathscr{O}(1/T)$  term.

**Theorem 2.** Assume that for all  $\theta \in \Theta$ ,  $(A_1)$  and  $(A_2)$  are verified. Assume, in addition, that there exists a constant B > 0 such that for any  $\theta \in \Theta$ ,  $\|\xi_{\alpha_{\theta}}\|_{\mathscr{H}} \leq B$ . Then with probability at least  $1 - \delta$ ,

$$\mathbb{P}\left(\mathbf{y} \neq g_{\widehat{\theta}_n}(\mathbf{x}) | \mathscr{D}_n\right) \le \widehat{\mathscr{R}}_n(\widehat{\theta}_n) + \frac{c_2}{T} + \frac{8BK_\ell}{(1-L)\sqrt{n}} + \frac{2BK_\ell}{1-L}\sqrt{\frac{\log(1/\delta)}{2n}},\tag{13}$$

where  $c_2 = K_\ell \sup_{\theta} \left( \|\psi\|_{\operatorname{op}} K_{f_{\theta}} e^{K_{f_{\theta}}} (L + \|f_{\theta}\|_{\infty} e^{K_{f_{\theta}}}) \right).$ 

Close to our result are the bounds obtained by Zhang et al. (2018), Tu et al. (2019), and Chen et al. (2020). The main difference is that the term in 1/T does not usually appear, since it comes from the Euler discretization error, whereas the speed in  $1/\sqrt{n}$  is the same. For instance, Chen et al. (2020) show that, under some assumptions, the excess risk is of order  $\sqrt{de + e^2}T^{\alpha}K_{\ell}n^{-1/2}$ . We refer to Section 5 for further discussion on the dependency of the different bounds to the parameter T. The take-home message is that the detour by continuous-time neural ODE provides a theoretical framework adapted to RNN, at the modest price of an additional  $\mathcal{O}(1/T)$  term. Moreover, we note that the bound (13) is 'simple' and holds under mild conditions for a large class of RNN. More precisely, for any recurrent network of the form (1), provided  $(A_1)$  and  $(A_2)$  are satisfied, then (13) is valid with constants  $c_2$  and B depending on the architecture. Such constants are given below in the example of a feedforward RNN. We stress that Theorem 2 can be extended without significant effort to the multi-class classification task, with an appropriate choice of loss function.

**Example 4.** Take a feedforward RNN with logistic activation, and  $\Theta = \{(W, b, \psi) \mid ||W||_F \leq K_W < (1-L)/32d, ||b|| \leq K_b, ||\psi||_{op} \leq K_{\psi}\}$ . Then, Proposition 5 states that  $(A_2)$  is satisfied and, with Theorem 1, ensures that

$$\sup_{\theta \in \Theta} \|\xi_{\alpha_{\theta}}\|_{\mathscr{H}} \le \frac{\sqrt{2K_{\psi}(1-L)}}{1-L-32dK_{W}} := B, \quad K_{f_{\theta}} = \max(\|U\|_{\text{op}}, \|V\|_{\text{op}}), \quad and \quad \|f_{\theta}\|_{\infty} = 1.$$

**Sequence-to-sequence learning.** We conclude by showing how to extend both the RKHS embedding of Theorem 1 and the generalization bound of Theorem 2 to the setting of sequence-to-sequence learning. In this case, the output of the network is a sequence

$$g_{\theta}(\mathbf{x}) = (z_1, \dots, z_T) \in (\mathbb{R}^p)^T.$$

An immediate extension of Theorem 1 ensures that there exist p elements  $\alpha_{1,\theta}, \ldots, \alpha_{p,\theta} \in \mathscr{T}$  such that, for any  $1 \leq j \leq T$ ,

$$\left\|z_{j}-\left(\langle\alpha_{1,\theta},S_{[0,j/T]}(\bar{X})\rangle_{\mathscr{T}},\ldots,\langle\alpha_{p,\theta},S_{[0,j/T]}(\bar{X})\rangle_{\mathscr{T}}\right)^{\top}\right\|\leq\|\psi\|_{\mathrm{op}}\frac{c_{1}}{T}.$$
(14)

The properties of the signature guarantee that  $S_{[0,j/T]}(X) = S(\tilde{X}_{[j]})$  where  $\tilde{X}_{[j]}$  is the process equal to  $\bar{X}$  on [0, j/T] and then constant on [j/T, 1]—see Appendix A.6. With this trick, we have, for any  $1 \leq \ell \leq p, \langle \alpha_{\ell,\theta}, S_{[0,j/T]}(\bar{X}) \rangle_{\mathcal{T}} = \langle \alpha_{\ell,\theta}, S(\tilde{X}_{[j]}) \rangle_{\mathcal{T}}$ , so that we are back in  $\mathscr{H}$ . Observe that the only difference with (12) is that we consider vector-valued sequential outputs, which requires to introduce the process  $\tilde{X}_{[j]}$ , but that the rationale is exactly the same.

We let  $\ell : (\mathbb{R}^p)^T \times (\mathbb{R}^p)^T \to \mathbb{R}^+$  be the  $L_2$  distance, that is, for any  $\mathbf{y} = (y_1, \ldots, y_T)$ ,  $\mathbf{y}' = (y'_1, \ldots, y'_T)$ ,  $\ell(\mathbf{y}, \mathbf{y}') = \frac{1}{T} \sum_{j=1}^T ||y_j - y'_j||^2$ . It is assumed that  $\mathbf{y}$  takes its values in a compact subset of  $\mathbb{R}^q$ , i.e., there exists  $K_y > 0$  such that  $||y_j|| \leq K_y$ .

**Theorem 3.** Assume that for all  $\theta \in \Theta$ ,  $(A_1)$  and  $(A_2)$  are verified. Assume, in addition, that there exists a constant B > 0 such that for any  $1 \le \ell \le p$ ,  $\theta \in \Theta$ ,  $\|\xi_{\alpha_{\ell,\theta}}\|_{\mathscr{H}} \le B$ . Then with probability at least  $1 - \delta$ ,

$$\mathscr{R}(\widehat{\theta}_n) \le \widehat{\mathscr{R}}_n(\widehat{\theta}_n) + \frac{c_3}{T} + \frac{4pc_4B(1-L)^{-1}}{\sqrt{n}} + \sqrt{\frac{2c_5\log(1/\delta)}{n}},\tag{15}$$

where  $c_3 = \sup_{\theta} (c_{1,\theta} + \|\psi\|_{op} \|f_{\theta}\|_{\infty}) + 2\sqrt{p}B(1-L)^{-1} + 2K_y$ ,  $c_4 = B(1-L)^{-1} + K_y$ , and  $c_5 = 4pB(1-L)^{-1}c_4 + K_y^2$ .

#### 3.2 Regularization and stability

In addition to providing a sound theoretical framework, framing deep learning in an RKHS provides a natural norm, which can be used for regularization, as shown for example in the context of convolutional neural networks by Bietti et al. (2019). This regularization ensures stability of predictions, which is crucial in particular in a small sample regime or in the presence of adversarial examples (Gao et al., 2018; Ko et al., 2019). In our binary classification setting, for any inputs  $\mathbf{x}, \mathbf{x}' \in (\mathbb{R}^d)^T$ , by the Cauchy-Schwartz inequality, we have

$$||z_T - z'_T|| \le 2||\psi||_{\text{op}}||\frac{c_1}{T} + ||\xi_{\alpha_\theta}(\bar{X}) - \xi_{\alpha_\theta}(\bar{X}')|| \le 2||\psi||_{\text{op}}||\frac{c_1}{T} + ||\xi_{\alpha_\theta}||_{\mathscr{H}}||S(\bar{X}) - S(\bar{X}')||_{\mathscr{T}}.$$

If x and x' are close, so are their associated continuous processes X and X' (which can be approximated for example by taking a piecewise linear interpolation), and so are their signatures. The term  $||S(\bar{X}) - S(\bar{X}')||_{\mathscr{T}}$  is therefore small (Friz and Victoir, 2010, Proposition 7.66). Therefore, when T is large, we see that the magnitude of  $||\xi_{\alpha_{\theta}}||_{\mathscr{H}}$  determines how close the predictions are. A natural training strategy to ensure stable predictions, for the types of networks covered in the present article, is then to penalize the problem by minimizing the loss  $\widehat{\mathscr{R}}_n(\theta) + \lambda ||\xi_{\alpha_{\theta}}||_{\mathscr{H}}^2$ . From a computational point of view, it is possible to compute the norm in  $\mathscr{H}$ , up to a truncation at N of the Taylor expansion, which we know by Proposition 4 to be reasonable. It remains that computing this norm is a non-trivial task, and implementing smart surrogates is an interesting problem for the future. Note however that computing the signature of the data is not necessary for this regularization strategy.

# 4 Numerical illustrations

This section is here for illustration purposes. Our objective is not to achieve competitive performance, but rather to illustrate the theoretical results. We refer to Appendix D for implementation details.



Figure 1: Approximation of the RNN ODE by the step-N Taylor expansion

**Convergence of the Taylor expansion towards the solution of the ODE.** We illustrate Proposition 4 on a toy example. The process X is a 2-dimensional spiral, and we take feedforward RNN with 2 hidden units. Repeating this procedure with  $10^3$  uniform random weight initializations, we observe in Figure 1a that the signature approximation converges exponentially fast in N. As seen in Figure 1b, the rate of convergence depends in particular on the norm of the weight matrices, as predicted by Proposition 5. However, condition (11) seems to be over-restrictive, since convergence happens even for weights with norm larger than the bound (we have  $1/(8a^2d) \simeq 0.01$  here).



Figure 2: Adversarial accuracy as a function of the adversarial perturbation  $\varepsilon$ 

Adversarial robustness. We illustrate the penalization proposed in Section 3.2 on a toy task that consists in classifying the rotation direction of 2-dimensional spirals. We take a feedforward RNN with 32 hidden units and hyperbolic tangent activation. It is trained on 50 examples, with and without penalization, for 200 epochs. Once trained, the RNN is tested on adversarial examples, generated with the projected gradient descent algorithm with Frobenius norm (Madry et al., 2018), which modifies test examples to maximize the error while staying in a ball of radius  $\varepsilon$ . We observe in Figure 2 that adding the penalization seems to make the network more stable.

**Comparison of the trained networks.** The evolution of the Frobenius norm of the weights  $||W||_F$  and the RKHS norm  $||\xi_{\alpha_{\theta}}||_{\mathscr{H}}$  during training is shown in Figure 3. This points out that the penalization, which forces the RNN to keep a small norm in  $\mathscr{H}$ , leads indeed to learning different weights than the non-penalized RNN. The results also suggest that the Frobenius and RKHS norms are decoupled, since both networks have Frobenius norms of similar magnitude but very different RKHS norms. The figures show one random run, but we observe similar qualitative behavior on others.

# 5 Discussion and conclusion

**Role of the discretization procedure.** The starting point of the paper was motivated by the fact that the classical residual RNN formulation coincides with an Euler discretization of the ODE (3).



Figure 3: Evolution of the Frobenius norm of the weights and of the RKHS norm during training

This choice of discretization translates into a 1/T term in Theorems 2 and 3. However, we could have considered higher-order discretization schemes, such as Runge-Kutta schemes, for which the discretization error decreases as  $1/T^p$ . Such schemes correspond to alternative architectures, which were already proposed by Wang and Lin (1998), among others. At the limit, we could also consider directly the continuous model (3), as proposed by Chen et al. (2018), in which case the discretization error term vanishes. Of course, such an option requires to be able to sample the continuous-time data at arbitrary times.

**Long-term stability.** RNN are known to be poor at learning long-term dependencies (Bengio et al., 1993; Hochreiter and Schmidhuber, 1997). This is reflected in the literature by performance bounds increasing in T, which is not the case of our results (13) and (15), seemingly indicating that we fail to capture this phenomenon. This apparent paradox is related to our assumption that the total variation of X is bounded. Indeed, if a time series is observed for a long time, then its total variation may become large. In this case, it is no longer valid to assume that  $||X||_{TV}$  is bounded by L. In other words, in our context, the parameter encapsulating the notion of "long-term" is not T but the regularity of X measured by its total variation. Note that the choice of defining X on [0, 1] and not another interval [0, U] is arbitrary and does not carry any meaning on the problem of learning long-term dependencies. A thorough analysis of these questions is an interesting research direction for future work.

**Radius of convergence.** The assumptions  $||X||_{TV;[0,1]} \le L < 1$  and  $||W||_F \le K_W < (1 - L)/32d$  can be seen as radii of convergence of the Taylor expansion (10). They allow using the Taylor approximation—which is of a local nature—to prove a global result, the RKHS embedding. In return, the condition on the Frobenius norm of the weights puts restrictions on the admissible parameters of the neural network. However, this bound can be improved, in particular by considering more exotic norms, which we did not explicit for clarity purposes.

**Conclusion.** By bringing together the theory of neural ODE, the signature transform, and kernel methods, we have shown that a recurrent network can be framed in the continuous-time limit as a linear function in a well-chosen RKHS. In addition to giving theoretical insights on the function learned by the network and providing generalization guarantees, this framing suggests regularization strategies to obtain more robust RNN. We have only scratched the surface of the potentialities of leveraging this theory to practical applications, which is a subject of its own and will be tackled in future work.

# Acknowledgements

Authors thank T. Lévy for his inputs on the Picard-Lindelöf theorem and N. Doumèche for fruitful discussion. A. Fermanian has been supported by a grant from Région Île-de-France and P. Marion by a stipend from Corps des Mines.

# References

- N.-J. Akpinar, B. Kratzwald, and S. Feuerriegel. Sample complexity bounds for recurrent neural networks with application to combinatorial graph problems. *arXiv:1901.10289*, 2019.
- P. L. Bartlett and S. Mendelson. Rademacher and Gaussian complexities: Risk bounds and structural results. *Journal of Machine Learning Research*, 3:463–482, 2002.
- M. Belkin, S. Ma, and S. Mandal. To understand deep learning we need to understand kernel learning. In J. Dy and A. Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80, pages 541–549. PMLR, 2018.
- Y. Bengio, P. Frasconi, and P. Simard. The problem of learning long-term dependencies in recurrent networks. In 1993 IEEE International Conference on Neural Networks, pages 1183–1188, 1993.
- A. Bietti and J. Mairal. Invariance and stability of deep convolutional representations. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30, pages 6210–6220. Curran Associates, Inc., 2017.
- A. Bietti and J. Mairal. Group invariance, stability to deformations, and complexity of deep convolutional representations. *Journal of Machine Learning Research*, 20:1–49, 2019.
- A. Bietti, G. Mialon, D. Chen, and J. Mairal. A kernel perspective for regularizing deep neural networks. In K. Chaudhuri and R. Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97, pages 664–674. PMLR, 2019.
- T. Cass, T. Lyons, C. Salvi, and W. Yang. Computing the untruncated signature kernel as the solution of a Goursat problem. *arXiv:2006.14794*, 2020.
- B. Chang, M. Chen, E. Haber, and E. H. Chi. AntisymmetricRNN: A dynamical system view on recurrent neural networks. In *International Conference on Learning Representations*, 2019.
- K.-T. Chen. Integration of paths–a faithful representation of paths by non-commutative formal power series. *Transactions of the American Mathematical Society*, 89:395–407, 1958.
- M. Chen, X. Li, and T. Zhao. On generalization bounds of a family of recurrent neural networks. In S. Chiappa and R. Calandra, editors, *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*, volume 108, pages 1233–1243, 2020.
- R. T. Q. Chen, Y. Rubanova, J. Bettencourt, and D. K. Duvenaud. Neural ordinary differential equations. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31, pages 6572–6583. Curran Associates, Inc., 2018.
- I. Chevyrev and A. Kormilitzin. A primer on the signature method in machine learning. *arXiv:1603.03788*, 2016.
- K. Cho, B. van Merriënboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio. Learning phrase representations using RNN encoder-decoder for statistical machine translation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing*, pages 1724–1734. Association for Computational Linguistics, 2014.
- Y. Cho and L. Saul. Kernel methods for deep learning. In Y. Bengio, D. Schuurmans, J. Lafferty, C. Williams, and A. Culotta, editors, *Advances in Neural Information Processing Systems*, volume 22, pages 342–350. Curran Associates, Inc., 2009.
- R. Collobert, J. Weston, L. Bottou, M. Karlen, K. Kavukcuoglu, and P. Kuksa. Natural language processing (almost) from scratch. *Journal of Machine Learning Research*, 12:2493–2537, 2011.
- E. De Brouwer, J. Simm, A. Arany, and Y. Moreau. GRU-ODE-Bayes: Continuous modeling of sporadically-observed time series. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32, pages 7379–7390. Curran Associates, Inc., 2019.

- N. B. Erichson, O. Azencot, A. Queiruga, L. Hodgkinson, and M. W. Mahoney. Lipschitz recurrent neural networks. In *International Conference on Learning Representations*, 2021.
- A. Fermanian. Embedding and learning with signatures. *Computational Statistics & Data Analysis*, 157:107148, 2021.
- P. Friz and N. Victoir. Euler estimates for rough differential equations. *Journal of Differential Equations*, 244:388–412, 2008.
- P. K. Friz and N. B. Victoir. *Multidimensional Stochastic Processes as Rough Paths: Theory and Applications*, volume 120 of *Cambridge Studies in Advanced Mathematics*. Cambridge University Press, Cambridge, 2010.
- J. Gao, J. Lanchantin, M. L. Soffa, and Y. Qi. Black-box generation of adversarial text sequences to evade deep learning classifiers. In 2018 IEEE Security and Privacy Workshops, pages 50–56, 2018.
- A. Graves, A.-r. Mohamed, and G. Hinton. Speech recognition with deep recurrent neural networks. In 2013 IEEE International Conference on Acoustics, Speech and Signal Processing, pages 6645– 6649, 2013.
- C. Herrera, F. Krach, and J. Teichmann. Theoretical guarantees for learning conditional expectation using controlled ODE-RNN. arXiv:2006.04727, 2020.
- G. Hinton, L. Deng, D. Yu, G. E. Dahl, A.-r. Mohamed, N. Jaitly, A. Senior, V. Vanhoucke, P. Nguyen, T. N. Sainath, et al. Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups. *IEEE Signal Processing Magazine*, 29:82–97, 2012.
- S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural Computation*, 9:1735–1780, 1997.
- A. Jacot, F. Gabriel, and C. Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31, pages 8580–8589. Curran Associates, Inc., 2018.
- J. Kelly, J. Bettencourt, M. J. Johnson, and D. K. Duvenaud. Learning differential equations that are easy to solve. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 4370–4380. Curran Associates, Inc., 2020.
- P. Kidger and T. Lyons. Signatory: Differentiable computations of the signature and logsignature transforms, on both CPU and GPU. In *International Conference on Learning Representations*, 2021.
- P. Kidger, P. Bonnier, I. Perez Arribas, C. Salvi, and T. Lyons. Deep signature transforms. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32, pages 3099–3109. Curran Associates, Inc., 2019.
- P. Kidger, J. Morrill, J. Foster, and T. Lyons. Neural controlled differential equations for irregular time series. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 6696–6707. Curran Associates, Inc., 2020.
- D. P. Kingma and J. Ba. Adam: A Method for Stochastic Optimization. In *Proceedings of the 3rd International Conference on Learning Representations (ICLR)*, 2015.
- F. J. Király and H. Oberhauser. Kernels for sequentially ordered data. *Journal of Machine Learning Research*, 20:1–45, 2019.
- Klaus Greff, Aaron Klein, Martin Chovanec, Frank Hutter, and Jürgen Schmidhuber. The Sacred Infrastructure for Computational Research. In Katy Huff, David Lippa, Dillon Niederhut, and M. Pacer, editors, *Proceedings of the 16th Python in Science Conference*, pages 49 – 56, 2017.

- C.-Y. Ko, Z. Lyu, L. Weng, L. Daniel, N. Wong, and D. Lin. POPQORN: Quantifying robustness of recurrent neural networks. In K. Chaudhuri and R. Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97, pages 3468–3477. PMLR, 2019.
- D. Levin, T. Lyons, and H. Ni. Learning from the past, predicting the statistics for the future, learning an evolving system. *arXiv:1309.0260*, 2013.
- S. Liao, T. Lyons, W. Yang, and H. Ni. Learning stochastic differential equations using RNN with log signature features. arXiv:1908.08286, 2019.
- S. H. Lim. Understanding recurrent neural networks using nonequilibrium response theory. *Journal of Machine Learning Research*, 22:1–48, 2021.
- T. Lyons. Rough paths, signatures and the modelling of functions on streams. *arXiv:1405.4537*, 2014.
- T. J. Lyons, M. J. Caruana, and T. Lévy. *Differential Equations Driven by Rough Paths*, volume 1908 of *Lecture Notes in Mathematics*. Springer, Berlin, 2007.
- A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu. Towards deep learning models resistant to adversarial attacks. In *International Conference on Learning Representations*, 2018.
- T. Mikolov, M. Karafiát, L. Burget, J. Černocký, and S. Khudanpur. Recurrent neural network based language model. In *Proceedings of the 11th Annual Conference of the International Speech Communication Association*, volume 2, pages 1045–1048, 2010.
- A. A. Minai and R. D. Williams. On the derivatives of the sigmoid. *Neural Networks*, 6:845–853, 1993.
- J. Morrill, C. Salvi, P. Kidger, J. Foster, and T. Lyons. Neural rough differential equations for long time series. arXiv:2009.08295, 2020a.
- J. H. Morrill, A. Kormilitzin, A. J. Nevado-Holgado, S. Swaminathan, S. D. Howison, and T. J. Lyons. Utilization of the signature method to identify the early onset of sepsis from multivariate physiological time series in critical care monitoring. *Critical Care Medicine*, 48:e976–e981, 2020b.
- A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala. PyTorch: An Imperative Style, High-Performance Deep Learning Library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32, pages 8024–8035. Curran Associates, Inc., 2019.
- I. Perez Arribas. Derivatives pricing using signature payoffs. arXiv:1809.09466, 2018.
- J. F. Reizenstein and B. Graham. Algorithm 1004: The iisignature library: Efficient calculation of iterated-integral signatures and log signatures. *ACM Transactions on Mathematical Software*, 46: article 8, 2020.
- J. Riordan. An Introduction to Combinatorial Analysis. John Wiley & Sons, New York, 1958.
- Y. Rubanova, R. T. Q. Chen, and D. K. Duvenaud. Latent ordinary differential equations for irregularly-sampled time series. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32, pages 5320–5330. Curran Associates, Inc., 2019.
- B. Schölkopf and A. J. Smola. *Learning with kernels: support vector machines, regularization, optimization, and beyond.* MIT press, Cambridge, Massachusetts, 2002.
- C. Toth and H. Oberhauser. Bayesian learning from sequential data using Gaussian processes with signature covariances. In H. Daumé III and A. Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119, pages 9548–9560, 2020.
- Z. Tu, F. He, and D. Tao. Understanding generalization in recurrent neural networks. In *International Conference on Learning Representations*, 2019.

- P. Virtanen, R. Gommers, T. E. Oliphant, M. Haberland, T. Reddy, D. Cournapeau, E. Burovski, P. Peterson, W. Weckesser, J. Bright, S. J. van der Walt, M. Brett, J. Wilson, K. J. Millman, N. Mayorov, A. R. J. Nelson, E. Jones, R. Kern, E. Larson, C. J. Carey, İ. Polat, Y. Feng, E. W. Moore, J. VanderPlas, D. Laxalde, J. Perktold, R. Cimrman, I. Henriksen, E. A. Quintero, C. R. Harris, A. M. Archibald, A. H. Ribeiro, F. Pedregosa, P. van Mulbregt, and SciPy 1.0 Contributors. SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python. *Nature Methods*, 17: 261–272, 2020.
- B. Wang, M. Liakata, H. Ni, T. Lyons, A. J. Nevado-Holgado, and K. Saunders. A path signature approach for speech emotion recognition. In *Proceedings of Interspeech 2019*, pages 1661–1665, 2019.
- Y.-J. Wang and C.-T. Lin. Runge-Kutta neural network for identification of dynamical systems in high accuracy. *IEEE Transactions on Neural Networks*, 9:294–307, 1998.
- W. Yang, L. Jin, and M. Liu. DeepWriterID: An end-to-end online text-independent writer identification system. *IEEE Intelligent Systems*, 31:45–53, 2016.
- W. Yang, T. Lyons, H. Ni, C. Schmid, and L. Jin. Developing the path signature methodology and its application to landmark-based human action recognition. *arXiv:1707.03993*, 2017.
- B. Yue, J. Fu, and J. Liang. Residual recurrent neural networks for learning sequential representations. *Information*, 9:56, 2018.
- J. Zhang, Q. Lei, and I. Dhillon. Stabilizing gradients for deep neural networks via efficient SVD parameterization. In J. Dy and A. Krause, editors, *Proceedings of the 35th International Conference* on Machine Learning, volume 80, pages 5806–5814. PMLR, 2018.

# Framing RNN as a kernel method: A neural ODE approach Supplementary material

# A Mathematical details

#### A.1 Writing the GRU and LSTM in the neural ODE framework

**GRU.** Recall that the equations of a GRU take the following form: for any  $1 \le j \le T$ ,

$$\begin{aligned} r_{j+1} &= \sigma(W_r x_{j+1} + b_r + U_r h_j) \\ z_{j+1} &= \sigma(W_z x_{j+1} + b_z + U_z h_j) \\ n_{j+1} &= \tanh(W_n x_{j+1} + b_n + r_{j+1} * (U_n h_j + c_n)) \\ h_{j+1} &= (1 - z_{j+1}) * h_j + z_{j+1} * n_{j+1}, \end{aligned}$$

where  $\sigma$  is the logistic activation, tanh the hyperbolic tangent, \* the Hadamard product,  $r_j$  the reset gate vector,  $z_j$  the update gate vector,  $W_r$ ,  $U_r$ ,  $W_z$ ,  $U_z$ ,  $W_n$ ,  $U_n$  weight matrices, and  $b_r$ ,  $b_z$ ,  $b_n$ ,  $c_n$ biases. Since  $r_{j+1}$ ,  $z_{j+1}$ , and  $n_{j+1}$  depend only on  $x_{j+1}$  and  $h_j$ , it is clear that these equations can be rewritten in the form

$$h_{j+1} = h_j + f(h_j, x_{j+1})$$

We then obtain equation (1) by normalizing f by 1/T.

**LSTM.** The LSTM networks are defined, for any  $1 \le j \le T$ , by

$$\begin{split} i_{j+1} &= \sigma(W_i x_{j+1} + b_i + U_i h_j) \\ f_{j+1} &= \sigma(W_f x_{j+1} + b_f + U_f h_j) \\ g_{j+1} &= \tanh(W_g x_{j+1} + b_g + U_g h_j) \\ o_{j+1} &= \sigma(W_o x_{j+1} + b_o + U_o h_j) \\ c_{j+1} &= f_{j+1} * c_j + i_{j+1} * g_{j+1} \\ h_{j+1} &= o_{j+1} * \tanh(c_{j+1}), \end{split}$$

where  $\sigma$  is the logistic activation, tanh the hyperbolic tangent, \* the Hadamard product,  $i_j$  the input gate,  $f_j$  the forget gate,  $g_j$  the cell gate,  $o_j$  the output gate,  $c_j$  the cell state,  $W_i$ ,  $U_i$ ,  $W_f$ ,  $U_f$ ,  $W_g$ ,  $U_g$ ,  $W_o$ ,  $U_o$  weight matrices, and  $b_i$ ,  $b_f$ ,  $b_g$ ,  $b_o$  biases. Since  $i_{j+1}$ ,  $f_{j+1}$ ,  $g_{j+1}$ ,  $o_{j+1}$  depend only on  $x_{j+1}$  and  $h_j$ , these equations can be rewritten in the form

$$h_{j+1} = f_1(h_j, x_{j+1}, c_{j+1})$$
  
$$c_{j+1} = f_2(h_j, x_{j+1}, c_j).$$

Let  $\tilde{h}_j = (h_j^{\top}, c_j^{\top})^{\top}$  be the hidden state defined by stacking the hidden and cell state. Then, clearly,  $\tilde{h}$  follows an equation of the form

$$\tilde{h}_{j+1} = f(\tilde{h}_j, x_{j+1}).$$

We obtain (1) by subtracting  $h_j$  and normalizing by 1/T.

# A.2 Picard-Lindelöf theorem

Consider a CDE of the form (8). We recall the Picard-Lindelöf theorem as given by Lyons et al. (2007, Theorem 1.3), and provide a proof for the sake of completeness.

**Theorem 4** (Picard-Lindelöf theorem). Assume that  $X \in BV^c([0,1], \mathbb{R}^d)$  and that  $\mathbf{F}$  is Lipschitzcontinuous with constant  $K_{\mathbf{F}}$ . Then, for any  $H_0 \in \mathbb{R}^e$ , the differential equation (8) admits a unique solution  $H : [0,1] \to \mathbb{R}^e$ .

*Proof.* Let  $\mathscr{C}([s,t]), \mathbb{R}^e)$  be the set of continuous functions from [s,t] to  $\mathbb{R}^e$ . For any  $[s,t] \subset [0,1]$ ,  $\zeta \in \mathbb{R}^e$ , let  $\Psi$  be the function

$$\Psi: \mathscr{C}([s,t]), \mathbb{R}^e) \to \mathscr{C}([s,t], \mathbb{R}^e)$$
$$Y \mapsto \left(v \mapsto \zeta + \int_s^v \mathbf{F}(Y_u) dX_u\right).$$

For any  $Y, Y' \in \mathscr{C}([s,t]), \mathbb{R}^e), v \in [s,t],$ 

$$\begin{aligned} \|\Psi(Y)_v - \Psi(Y')_v\| &\leq \int_s^v \left\| \left( \mathbf{F}(Y_u) - \mathbf{F}(Y'_u) \right) dX_u \right\| \\ &\leq \int_s^v \|\mathbf{F}(Y_u) - \mathbf{F}(Y'_u)\|_{\text{op}} \|dX_u\| \\ &\leq \int_s^v K_{\mathbf{F}} \|Y_u - Y'_u\| \|dX_u\| \\ &\leq K_{\mathbf{F}} \|Y - Y'\|_{\infty} \int_s^v \|dX_u\| \\ &\leq K_{\mathbf{F}} \|Y - Y'\|_{\infty} \|X\|_{TV;[s,t]}. \end{aligned}$$

This shows that the function  $\Psi$  is Lipschitz-continuous on  $\mathscr{C}([s,t]), \mathbb{R}^e)$  endowed with the supremum norm, with Lipschitz constant  $K_{\mathbf{F}} \|X\|_{TV;[s,t]}$ . Clearly, the function  $t \mapsto \|X\|_{TV;[0,t]}$  is nondecreasing and uniformly continuous on the compact interval [0,1]. Therefore, for any  $\varepsilon > 0$ , there exists  $\delta > 0$  such that

$$|t-s| < \delta \Rightarrow \left| \|X\|_{TV;[0,t]} - \|X\|_{TV;[0,s]} \right| < \varepsilon.$$

Take  $\varepsilon = 1/K_{\mathbf{F}}$ . Then on any interval [s, t] of length smaller than  $\delta$ , one has  $||X||_{TV;[s,t]} = ||X||_{TV;[0,t]} - ||X||_{TV;[0,s]} < 1/K_{\mathbf{F}}$ , so that the function  $\Psi$  is a contraction. By the Banach fixed-point theorem, for any initial value  $\zeta$ ,  $\Psi$  has a unique fixed point. Hence, there exists a solution to (8) on any interval of length  $\delta$  with any initial condition. To obtain a solution on [0, 1] it is sufficient to concatenate these solutions.

A corollary of this theorem is a Picard-Lindelöf theorem for initial value problems of the form

$$dH_t = f(H_t, X_t)dt, \quad H_0 = \zeta, \tag{16}$$

where  $f : \mathbb{R}^e \times \mathbb{R}^d \to \mathbb{R}^e, \zeta \in \mathbb{R}^e$ .

**Corollary 1.** Assume that f is Lipschitz continuous in its first variable. Then, for any  $\zeta \in \mathbb{R}^e$ , the initial value problem (16) admits a unique solution.

*Proof.* Let  $f_X : (h,t) \mapsto f(h,X_t)$ . Then the solution of (16) is solution of the differential equation  $dH_t = f_X(H_t,t)dt.$ 

Let d = 1,  $\bar{e} = e + 1$ , and **F** be the vector field defined by

$$\mathbf{F}: h \mapsto \begin{pmatrix} f_X(h^{1:e}, h^{e+1}) \\ 1 \end{pmatrix},$$

where  $h^{1:e}$  denotes the projection of h on its first e coordinates. Then, since  $f_X$  is Lipschitz, so is the vector field **F**. Theorem 4 therefore applies to the differential equation

$$dH_t = \mathbf{F}(H_t)dt, \quad H_0 = (\zeta^{\top}, 0)^{\top}.$$

Projecting this differential equation on the last coordinate gives  $dH_t^{e+1} = dt$ , that is,  $H_t^{e+1} = t$ . Projecting on the first *e* coordinates exactly provides equation (16), which therefore has a unique solution, equal to  $H^{1:e}$ .

#### A.3 Operator norm

**Definition 3.** Let  $(E, \|\cdot\|_E)$  and  $(F, \|\cdot\|_F)$  be two normed vector spaces and let  $f \in \mathscr{L}(E, F)$ , where  $\mathscr{L}(E, F)$  is the space of linear functions from E to F. The operator norm of f is defined by

$$||f||_{\text{op}} = \sup_{u \in E, ||u||_E = 1} ||f(u)||_F.$$

Equipped with this norm,  $\mathscr{L}(E, F)$  is a normed vector space.

This definition is valid when f is represented by a matrix.

#### A.4 Tensor Hilbert space

Let us first briefly recall some elements on tensor spaces. If  $e_1, \ldots, e_d$  is the canonical basis of  $\mathbb{R}^d$ , then  $(e_{i_1} \otimes \cdots \otimes e_{i_k})_{1 \leq i_1, \ldots, i_k \leq d}$  is a basis of  $(\mathbb{R}^d)^{\otimes k}$ . Any element  $a \in (\mathbb{R}^d)^{\otimes k}$  can therefore be written as

$$a = \sum_{1 \le i_1, \dots, i_k \le d} a^{(i_1, \dots, i_k)} e_{i_1} \otimes \dots \otimes e_{i_k},$$

where  $a^{(i_1,...,i_k)} \in \mathbb{R}$ . The tensor space  $(\mathbb{R}^d)^{\otimes k}$  is a Hilbert space of dimension  $d^k$ , with scalar product

$$\langle a,b\rangle_{(\mathbb{R}^d)^{\otimes k}} = \sum_{1 \le i_1,\dots,i_k \le d} a^{(i_1,\dots,i_k)} b^{(i_1,\dots,i_k)}$$

and associated norm  $\|\cdot\|_{(\mathbb{R}^d)^{\otimes k}}$ .

We now consider the space  $\mathscr{T}$  defined by (6). The sum, multiplication by a scalar, and scalar product on  $\mathscr{T}$  are defined as follows: for any  $a = (a_0, \ldots, a_k, \ldots) \in \mathscr{T}, b = (b_0, \ldots, b_k, \ldots) \in \mathscr{T}, \lambda \in \mathbb{R},$ 

$$a + \lambda b = (a_0 + \lambda b_0, \dots, a_k + \lambda b_k, \dots)$$
 and  $\langle a, b \rangle_{\mathscr{T}} = \sum_{k=0}^{\infty} \langle a_k, b_k \rangle_{(\mathbb{R}^d)^{\otimes k}},$ 

with the convention  $(\mathbb{R}^d)^{\otimes 0} = \mathbb{R}$ .

**Proposition 6.**  $(\mathcal{T}, +, \cdot, \langle \cdot, \cdot \rangle_{\mathcal{T}})$  is a Hilbert space.

*Proof.* By the Cauchy-Schwartz inequality,  $\langle \cdot, \cdot \rangle_{\mathscr{T}}$  is well-defined: for any  $a, b \in \mathscr{T}$ ,

$$\begin{aligned} |\langle a,b\rangle_{\mathscr{T}}| &\leq \sum_{k=0}^{\infty} |\langle a_k,b_k\rangle_{(\mathbb{R}^d)^{\otimes k}}| \leq \sum_{k=0}^{\infty} \|a_k\|_{(\mathbb{R}^d)^{\otimes k}} \|b_k\|_{(\mathbb{R}^d)^{\otimes k}} \\ &\leq \Big(\sum_{k=0}^{\infty} \|a_k\|_{(\mathbb{R}^d)^{\otimes k}}^2\Big)^{1/2} \Big(\sum_{k=0}^{\infty} \|b_k\|_{(\mathbb{R}^d)^{\otimes k}}^2\Big)^{1/2} < \infty. \end{aligned}$$

Moreover,  $\mathscr{T}$  is a vector space: for any  $a, b \in \mathscr{T}, \lambda \in \mathbb{R}$ , since

 $a + \lambda b = (a_0 + \lambda b_0, \dots, a_k + \lambda b_k, \dots),$ 

and

$$\begin{split} \sum_{k=0}^{\infty} \|a_k + \lambda b_k\|_{(\mathbb{R}^d)^{\otimes k}}^2 &= \sum_{k=0}^{\infty} \|a_k\|_{(\mathbb{R}^d)^{\otimes k}}^2 + \lambda^2 \sum_{k=0}^{\infty} \|b_k\|_{(\mathbb{R}^d)^{\otimes k}}^2 \\ &+ 2\lambda \sum_{k=0}^{\infty} \langle a_k, b_k \rangle_{(\mathbb{R}^d)^{\otimes k}} \\ &\leq \sum_{k=0}^{\infty} \|a_k\|_{(\mathbb{R}^d)^{\otimes k}}^2 + \lambda^2 \sum_{k=0}^{\infty} \|b_k\|_{(\mathbb{R}^d)^{\otimes k}}^2 + 2\lambda \langle a, b \rangle_{\mathscr{T}} < \infty, \end{split}$$

we see that  $a + \lambda b \in \mathscr{T}$ . The operation  $\langle \cdot, \cdot \rangle_{\mathscr{T}}$  is also bilinear, symmetric, and positive definite:

$$\langle a,a\rangle_{\mathscr{T}} = 0 \Leftrightarrow \sum_{k=0}^{\infty} \|a_k\|^2_{(\mathbb{R}^d)^{\otimes k}} = 0 \Leftrightarrow \forall k \in \mathbb{N}, \|a_k\|^2_{(\mathbb{R}^d)^{\otimes k}} = 0 \Leftrightarrow \forall k \in \mathbb{N}, a_k = 0 \Leftrightarrow a = 0.$$

Therefore  $\langle \cdot, \cdot \rangle_{\mathscr{T}}$  is an inner product on  $\mathscr{T}$ . Finally, let  $(a^{(n)})_{n \in \mathbb{N}}$  be a Cauchy sequence in  $\mathscr{T}$ . Then, for any  $n, m \geq 0$ ,

$$\|a^{(n)} - a^{(m)}\|_{\mathscr{T}}^2 = \sum_{k=0}^{\infty} \|a_k^{(n)} - a_k^{(m)}\|_{(\mathbb{R}^d)^{\otimes k}}^2,$$

so for any  $k \in \mathbb{N}$ , the sequence  $(a_k^{(n)})_{n \in \mathbb{N}}$  is Cauchy in  $(\mathbb{R}^d)^{\otimes k}$ . Since  $(\mathbb{R}^d)^{\otimes k}$  is a Hilbert space,  $(a_k^{(n)})_{n \in \mathbb{N}}$  converges to a limit  $a_k^{(\infty)} \in (\mathbb{R}^d)^{\otimes k}$ . Let  $a^{(\infty)} = (a_0^{(\infty)}, \dots, a_k^{(\infty)}, \dots)$ . To finish the

proof, we need to show that  $a^{(\infty)} \in \mathscr{T}$  and that  $a^{(n)}$  converges to  $a^{(\infty)}$  in  $\mathscr{T}$ . First, note that there exists a constant B > 0 such that for any  $n \in \mathbb{N}$ ,

$$\|a^{(n)}\|_{\mathscr{T}} \le B.$$

To see this, observe that for  $\varepsilon > 0$ , there exists  $N \in \mathbb{N}$  such that for any  $n \ge N$ ,  $||a^{(n)} - a^{(N)}||_{\mathscr{T}} < \varepsilon$ , and so  $||a^{(n)}||_{\mathscr{T}} \le \varepsilon + ||a^{(N)}||_{\mathscr{T}}$ . Take  $B = \max(||a^{(1)}||_{\mathscr{T}}, \ldots, ||a^{(N)}||_{\mathscr{T}}, \varepsilon + ||a^{(N)}||_{\mathscr{T}})$ . Then, for any  $K \in \mathbb{N}$ ,

$$\sum_{k=0}^{K} \|a_{k}^{(n)}\|_{(\mathbb{R}^{d})^{\otimes k}}^{2} \leq \|a^{(n)}\|_{\mathscr{T}} \leq B.$$

Letting  $K \to \infty$ , we obtain that  $\|a^{(\infty)}\|_{\mathscr{T}} \leq B$ , and therefore  $a^{(\infty)} \in \mathscr{T}$ . Finally, let  $\varepsilon > 0$  and let  $N \in \mathbb{N}$  be such that for any  $n, m \geq N$ ,  $\|a^{(n)} - a^{(m)}\|_{\mathscr{T}} < \varepsilon$ . Clearly, for any  $K \in \mathbb{N}$ ,

$$\sum_{k=0}^{K} \|a_k^{(n)} - a_k^{(m)}\|_{(\mathbb{R}^d)^{\otimes k}}^2 < \varepsilon^2$$

Letting  $m \to \infty$  leads to

$$\sum_{k=1}^{K} \|a_k^{(n)} - a_k^{(\infty)}\|_{(\mathbb{R}^d)^{\otimes k}}^2 < \varepsilon^2,$$

and letting  $K \to \infty$  gives

$$\|a^{(n)} - a^{(\infty)}\|_{\mathscr{T}} < \varepsilon,$$

which completes the proof.

# A.5 Bounding the derivatives of the logistic and hyperbolic tangent activations

**Lemma 1.** Let  $\sigma$  be the logistic function defined, for any  $x \in \mathbb{R}$ , by  $\sigma(x) = 1/(1+e^{-x})$ . Then, for any  $n \ge 0$ ,  $\|\sigma^{(n)}\|_{\infty} \le 2^{n-1}n!$ .

*Proof.* For any  $x \in \mathbb{R}$ , one has (Minai and Williams, 1993, Theorem 2)

$$\sigma^{(n)}(x) = \sum_{k=1}^{n+1} (-1)^{k-1} (k-1)! \left\{ \binom{n+1}{k} \sigma(x)^k \right\},$$

where  $\binom{n}{k}$  stands for the Stirling number of the second kind (see, e.g., Riordan, 1958). Let

$$u_n = \sum_{k=1}^{n+1} (k-1)! \binom{n+1}{k}$$

for  $n \ge 1$  and  $u_0 = 1$ . Since  $0 \le \sigma(x) \le 1$ , it is clear that  $|\sigma^{(n)}(x)| \le u_n$ . Using the fact that the Stirling numbers satisfy the recurrence relation

$$\binom{n+1}{k} = k \binom{n}{k} + \binom{n}{k-1},$$

valid for all  $0 \le k \le n$ , we have

$$u_n = \sum_{k=1}^n (k-1)! \left( k \begin{Bmatrix} n \\ k \end{Bmatrix} + \begin{Bmatrix} n \\ k-1 \end{Bmatrix} \right) + n! = \sum_{k=1}^n k! \begin{Bmatrix} n \\ k \end{Bmatrix} + \sum_{k=0}^{n-1} k! \begin{Bmatrix} n \\ k \end{Bmatrix} + n! = 2 \sum_{k=1}^n k! \begin{Bmatrix} n \\ k \end{Bmatrix}$$
  
(since  $\begin{Bmatrix} n \\ 0 \end{Bmatrix} = 0$ )  
 $\leq 2n \sum_{k=1}^n (k-1)! \begin{Bmatrix} n \\ k \end{Bmatrix} = 2nu_{n-1}.$ 

Thus, by induction,  $u_n \leq 2^{n-1}n!$ , from which the claim follows.

**Lemma 2.** Let tanh be the hyperbolic tangent function. Then, for any  $n \ge 0$ ,

$$\| \operatorname{tanh}^{(n)} \|_{\infty} \leq 4^n n!$$

*Proof.* Let  $\sigma$  be the logistic function. Straightforward calculations yield the equality, valid for any  $x \in \mathbb{R}$ ,

$$\tanh(x) = 2\sigma(2x) - 1.$$

But, for any  $n \ge 1$ ,

$$\tanh^{(n)}(x) = 2^{n+1}\sigma^{(n)}(2x),$$

and thus, by Lemma 1,

$$\| \tanh^{(n)} \|_{\infty} \le 2^{n+1} \| \sigma^{(n)} \|_{\infty} \le 4^n n!$$

The inequality is also true for n = 0 since  $\|\tanh\|_{\infty} \leq 1$ .

#### A.6 Chen's formula

First, note that it is straightforward to extend the definition of the signature to any interval  $[s, t] \subset [0, 1]$ . The next proposition, known as Chen's formula (Lyons et al., 2007, Theorem 2.9), tells us that the signature can be computed iteratively as tensor products of signatures on subintervals.

**Proposition 7.** Let  $X \in BV^c([s,t], \mathbb{R}^d)$  and  $u \in (s,t)$ . Then

$$S_{[s,t]}(X) = S_{[s,u]}(X) \otimes S_{[u,t]}(X).$$

Next, it is clear that the signature of a constant path is equal to  $\mathbf{1} = (1, 0, \dots, 0, \dots)$  which is the null element in  $\mathscr{T}$ . Indeed, let  $Y \in BV^c([s, t], \mathbb{R}^d)$  be a constant path. Then, for any  $k \ge 1$ ,

$$\mathbb{Y}_{[s,t]}^k = k! \int \cdots \int _{s \le u_1 < \cdots < u_k \le t} dY_{u_1} \otimes \cdots \otimes dY_{u_k} = k! \int \cdots \int _{s \le u_1 < \cdots < u_k \le t} 0 \otimes \cdots \otimes 0 = 0.$$

Now let  $X \in BV^c([0,1], \mathbb{R}^d)$  and consider the path  $\tilde{X}_{[j]}$  equal to the time-augmented path  $\bar{X}$  on [0, j/T] and then constant on [j/T, 1]—see Figure 4. We have by Proposition 7

$$S_{[0,1]}(\tilde{X}_{[j]}) = S_{[0,j/T]}(\tilde{X}_{[j]}) \otimes S_{[j/T,1]}(\tilde{X}_{[j]}) = S_{[0,j/T]}(\bar{X}) \otimes \mathbf{1} = S_{[0,j/T]}(\bar{X}).$$



Figure 4: Example of a path  $X \in BV^c([0,1],\mathbb{R})$  (left) and its corresponding paths  $\tilde{X}_{[j]}$ , plotted against time, for different values of  $j \in \{1, \ldots, T\}$  (right)

# **B Proofs**

## **B.1** Proof of Proposition 1

According to Assumption  $(A_1)$ , for any  $h_1, h_2 \in \mathbb{R}^e, x_1, x_2 \in \mathbb{R}^d$ , one has

$$||f(h_1, x_1) - f(h_2, x_1)|| \le K_f ||h_1 - h_2||$$
 and  $||f(h_1, x_1) - f(h_1, x_2)|| \le K_f ||x_1 - x_2||$ .

Under assumption  $(A_1)$ , by Corollary 1, the initial value problem (3) admits a unique solution H. Let us first show that for any  $t \in [0, 1]$ ,  $H_t$  is bounded independently of X. For any  $t \in [0, 1]$ ,

$$\begin{aligned} \|H_t - H_0\| &= \left\| \int_0^t f(H_u, X_u) du \right\| \le \int_0^t \|f(H_u, X_u)\| du \\ &= \int_0^t \|f(H_u, X_u) - f(H_0, X_u) + f(H_0, X_u)\| du \\ &\le \int_0^t \|f(H_u, X_u) - f(H_0, X_u)\| + \int_0^t \|f(H_0, X_u)\| du \\ &\le K_f \int_0^t \|H_u - H_0\| du + t \sup_{\|x\| \le L} \|f(H_0, x)\|. \end{aligned}$$

Applying Grönwall's inequality to the function  $t \mapsto ||H_t - H_0||$  yields

$$||H_t - H_0|| \le t \sup_{||x|| \le L} ||f(H_0, x)|| \exp\left(\int_0^t K_f du\right) \le \sup_{||x|| \le L} ||f(H_0, x)|| e^{K_f} := M.$$

Given that  $H_0 = h_0 = 0$ , we conclude that  $||H_t|| \le M$ . Next, let

xi, lei

$$||f||_{\infty} = \sup_{||x|| \le L, ||h|| \le M} f(h, x)$$

By similar arguments, for any  $[s, t] \subset [0, 1]$ , Grönwall's inequality applied to the function  $t \mapsto ||H_t - H_s||$  yields

$$||H_t - H_s|| \le (t - s)||f||_{\infty} e^{K_f}.$$

Therefore, for any partition  $(t_0, \ldots, t_k)$  of [s, t],

$$\sum_{i=1}^{k} \|H_{t_i} - H_{t_{i-1}}\| \le \|f\|_{\infty} e^{K_f} \sum_{i=1}^{k} (t_i - t_{i-1}) \le \|f\|_{\infty} e^{K_f} (t-s),$$

and, taking the supremum over all partitions of [s,t],  $||H||_{TV;[s,t]} \leq ||f||_{\infty} e^{K_f}(t-s)$ . In other words, H is of bounded variation on any interval  $[s,t] \subset [0,1]$ . Let  $(t_0,\ldots,t_T)$  denote the regular partition of [0,1] with  $t_j = j/T$ . For any  $1 \leq j \leq T$ , we have

$$||H_{t_j} - h_j|| = ||H_{t_{j-1}} + \int_{t_{j-1}}^{t_j} f(H_u, X_u) du - h_{j-1} - \frac{1}{T} f(h_{j-1}, x_j)||$$
  
$$\leq ||H_{t_{j-1}} - h_{j-1}|| + \int_{t_{j-1}}^{t_j} ||f(H_u, X_u) - f(h_{j-1}, x_j)|| du.$$

Writing

$$\begin{aligned} \left\| f(H_u, X_u) - f(h_{j-1}, x_j) \right\| &= \left\| f(H_u, X_u) - f(H_u, x_j) + f(H_u, x_j) - f(h_{j-1}, x_j) \right\| \\ &\leq \left\| f(H_u, X_u) - f(H_u, x_j) \right\| + \left\| f(H_u, x_j) - f(h_{j-1}, x_j) \right\| \\ &\leq K_f \left\| X_u - x_j \right\| + K_f \left\| H_u - h_{j-1} \right\|, \end{aligned}$$

we obtain

$$\begin{aligned} \|H_{t_j} - h_j\| &\leq \|H_{t_{j-1}} - h_{j-1}\| + K_f \int_{t_{j-1}}^{t_j} \|H_u - h_{j-1}\| du + K_f \int_{t_{j-1}}^{t_j} \|X_u - x_j\| du \\ &\leq \|H_{t_{j-1}} - h_{j-1}\| + K_f \int_{t_{j-1}}^{t_j} \left( \|H_u - H_{t_{j-1}}\| + \|H_{t_{j-1}} - h_{j-1}\| \right) du \\ &\quad + \frac{K_f}{T} \|X\|_{TV;[t_{j-1},t_j]} \\ &\leq \left(1 + \frac{K_f}{T}\right) \|H_{t_{j-1}} - h_{j-1}\| + \frac{K_f}{T} \left( \|H\|_{TV;[t_{j-1},t_j]} + \|X\|_{TV;[t_{j-1},t_j]} \right). \end{aligned}$$

By induction, we are led to

$$\begin{aligned} \|H_{t_j} - h_j\| &\leq \frac{K_f}{T} \sum_{k=0}^{j-1} \left( 1 + \frac{K_f}{T} \right)^k \left( \|H\|_{TV;[t_k, t_{k+1}]} + \|X\|_{TV;[t_k, t_{k+1}]} \right) \\ &\leq \frac{K_f}{T} \left( 1 + \frac{K_f}{T} \right)^T \left( \|X\|_{TV;[0,1]} + \|H\|_{TV;[0,1]} \right) \\ &\leq \frac{K_f e^{K_f}}{T} \left( L + \|f\|_{\infty} e^{K_f} \right), \end{aligned}$$

which concludes the proof.

# **B.2** Proof of Proposition 2

Let  $\bar{h} \in \mathbb{R}^{\bar{e}}$  and let  $\bar{h}^{i:j} = (\bar{h}^i, \dots, \bar{h}^j)$  be its projection on a subset of coordinates. It is sufficient to take  $\mathbf{F}$  defined by

$$\mathbf{F}(\bar{h}) = \begin{pmatrix} 0_{e \times d} & \frac{2}{1-L}f(\bar{h}^{1:e}, \bar{h}^{e+1:e+d}) \\ I_{d \times d} & 0_{d \times 1} \end{pmatrix},$$

where  $I_{d \times d}$  denotes the identity matrix and  $0_{\times}$  the matrix full of zeros. The function  $\overline{H}$  is then solution of

$$d\bar{H}_t = \begin{pmatrix} 0_{e \times d} & \frac{2}{1-L}f(\bar{H}_t^{1:e}, \bar{H}_t^{e+1:e+d})\\ I_{d \times d} & 0_{d \times 1} \end{pmatrix} \begin{pmatrix} dX_t\\ \frac{1-L}{2}dt \end{pmatrix}$$

Note that under assumption  $(A_1)$ , the tensor field **F** satisfies the assumptions of the Picard-Lindelöf theorem (Theorem 4) so that  $\overline{H}$  is well-defined. The projection of this equation on the last d coordinates gives

$$d\bar{H}_t^{e+1:e+d} = dX_t, \quad \bar{H}_0^{e+1:e+d} = X_0$$

and therefore  $\bar{H}_t^{e+1:e+d} = X_t$ . The projection on the first *e* coordinates gives

$$d\bar{H}_t^{1:e} = \frac{2}{1-L} f(\bar{H}_t^{1:e}, X_t) \frac{1-L}{2} dt = f(\bar{H}_t^{1:e}, X_t) dt, \quad \bar{H}_0^{1:e} = h_0,$$

which is exactly (3).

#### **B.3** Proof of Proposition 3

According to Lyons (2014, Lemma 5.1), one has

$$\|\bar{\mathbb{X}}_{[0,t]}^k\|_{(\mathbb{R}^d)^{\otimes k}} \le \|\bar{X}\|_{TV;[0,t]}^k.$$

Let  $(t_0, \ldots, t_k)$  be a partition of [0, t]. Then

$$\begin{split} \sum_{j=1}^{k} \|\bar{X}_{t_{j}} - \bar{X}_{t_{j-1}}\| &= \sum_{j=1}^{k} \sqrt{\|X_{t_{j}} - X_{t_{j-1}}\|^{2} + \left(\frac{1-L}{2}\right)^{2} (t_{j} - t_{j-1})^{2}} \\ &\leq \sum_{j=1}^{k} \|X_{t_{j}} - X_{t_{j-1}}\| + \frac{1-L}{2} \sum_{j=1}^{k} (t_{j} - t_{j-1}) \\ &= \sum_{j=1}^{k} \|X_{t_{j}} - X_{t_{j-1}}\| + \frac{1-L}{2} t. \end{split}$$

Taking the supremum over any partition of [0, t] we obtain

$$\|\bar{X}\|_{TV;[0,t]} \le \|X\|_{TV;[0,t]} + \frac{1-L}{2}t \le L + \frac{1-L}{2} = \frac{1+L}{2} < 1,$$

and thus  $\|\bar{\mathbb{X}}^k_{[0,t]}\|_{(\mathbb{R}^d)^{\otimes k}} \leq \left(\frac{1+L}{2}\right)^k$ . It is then clear that

$$\|S_{[0,t]}(\bar{X})\|_{\mathscr{T}} = \left(\sum_{k=0}^{\infty} \|\bar{\mathbb{X}}_{[0,t]}^{k}\|_{(\mathbb{R}^{d})^{\otimes k}}^{2}\right)^{1/2} \le \sum_{k=0}^{\infty} \|\bar{\mathbb{X}}_{[0,t]}^{k}\|_{(\mathbb{R}^{d})^{\otimes k}} \le \sum_{k=0}^{\infty} \left(\frac{1+L}{2}\right)^{k} = 2(1-L)^{-1}.$$

# **B.4** Proof of Proposition 4

We first recall the fundamental theorem of calculus for line integrals (also known as gradient theorem). **Theorem 5.** Let  $g : \mathbb{R}^e \to \mathbb{R}$  be a continuously differentiable function, and let  $\gamma : [a, b] \to \mathbb{R}^e$  be a smooth curve in  $\mathbb{R}^e$ . Then

$$\int_{a}^{b} \nabla g(\gamma_t) d\gamma_t = g(\gamma_b) - g(\gamma_a),$$

where  $\nabla g$  denotes the gradient of g.

The identity above immediately generalizes to a function  $g : \mathbb{R}^e \to \mathbb{R}^e$ :

$$\int_{a}^{b} J(g)(\gamma_t) d\gamma_t = g(\gamma_b) - g(\gamma_a),$$

where  $J(g) \in \mathbb{R}^{e \times e}$  is the Jacobian matrix of g. Let us apply Theorem 5 to the vector field  $F^i$  between 0 and t, with  $\gamma = H$ . We have

$$F^{i}(H_{t}) - F^{i}(H_{0}) = \int_{0}^{t} J(F^{i})(H_{u}) dH_{u} = \int_{0}^{t} J(F^{i})(H_{u}) \sum_{j=1}^{d} F^{j}(H_{u}) dX_{u}$$
$$= \sum_{j=1}^{d} \int_{0}^{t} J(F^{i})(H_{u}) F^{j}(H_{u}) dX_{u} = \sum_{j=1}^{d} \int_{0}^{t} F^{j} \star F^{i}(H_{u}) dX_{u}.$$

Iterating this procedure (N-1) times for the vector fields  $F^1, \ldots, F^d$  yields

$$\begin{split} H_t &= H_0 + \sum_{i=1}^d \int_0^t F^i(H_u) dX_u^i \\ &= H_0 + \sum_{i=1}^d \int_0^t F^i(H_0) dX_u^i + \sum_{i=1}^d \int_0^t \sum_{j=1}^d \int_0^u F^j \star F^i(H_v) dX_v^j dX_u^i \\ &= H_0 + \sum_{i=1}^d F^i(H_0) S_{[0,t]}^{(i)}(X) + \sum_{1 \le i,j \le d} \int_{0 \le v \le u \le t} F^j \star F^i(H_v) dX_v^j dX_u^i \\ &= \cdots \\ &= H_0 + \sum_{k=1}^N \sum_{1 \le i_1, \dots, i_k \le d} F^{i_1} \star \dots \star F^{i_k}(H_0) \frac{1}{k!} S_{[0,t]}^{(i_1,\dots,i_k)}(X) \\ &+ \sum_{1 \le i_1,\dots, i_{N+1} \le d} \int_{\Delta_{N+1;[0,t]}} F^{i_1} \star \dots \star F^{i_{N+1}}(H_{u_1}) dX_{u_1}^{i_1} \cdots dX_{u_{N+1}}^{i_{N+1}}, \end{split}$$

where  $\Delta_{N;[0,t]} := \{(u_1, \cdots, u_N) \in [0,t]^N \mid 0 \le u_1 < \cdots < u_N \le t\}$  is the simplex in  $[0,t]^N$ . The first (N+1) terms equal  $H_t^N$ . Hence,

$$\begin{split} \|H_{t} - H_{t}^{N}\| &= \Big\| \sum_{1 \leq i_{1}, \dots, i_{N+1} \leq d} \int_{\Delta_{N+1;[0,t]}} F^{i_{1}} \star \dots \star F^{i_{N+1}}(H_{u_{1}}) dX_{u_{1}}^{i_{1}} \dots dX_{u_{N+1}}^{i_{N+1}} \Big\| \\ &\leq \sum_{1 \leq i_{1}, \dots, i_{N+1} \leq d} \int_{\Delta_{N+1;[0,t]}} \|F^{i_{1}} \star \dots \star F^{i_{N+1}}(H_{u_{1}})\| |dX_{u_{1}}^{i_{1}}| \dots |dX_{u_{N+1}}^{i_{N+1}}| \\ &\leq \sum_{1 \leq i_{1}, \dots, i_{N+1} \leq d} \int_{\Delta_{N+1;[0,t]}} \sup_{1 \leq i_{1}, \dots, i_{N+1} \leq d, \|h\| \leq M} \|F^{i_{1}} \star \dots \star F^{i_{N+1}}(h)\| |dX_{u_{1}}^{i_{1}}| \dots |dX_{u_{N+1}}^{i_{N+1}}| \\ &\leq \Lambda_{N+1}(\mathbf{F}) \sum_{1 \leq i_{1}, \dots, i_{N+1} \leq d} \int_{\Delta_{N+1;[0,t]}} |dX_{u_{1}}^{i_{1}}| \dots |dX_{u_{N+1}}^{i_{N+1}}|. \end{split}$$

Thus,

$$\begin{aligned} \|H_t - H_t^N\| &\leq \Lambda_{N+1}(\mathbf{F}) \sum_{1 \leq i_1, \dots, i_{N+1} \leq d} \int_{\Delta_{N+1;[0,t]}} |dX_{u_1}^{i_1}| \cdots |dX_{u_{N+1}}^{i_{N+1}}| \\ &\leq \Lambda_{N+1}(\mathbf{F}) \sum_{1 \leq i_1, \dots, i_{N+1} \leq d} \int_{\Delta_{N+1;[0,t]}} \|dX_{u_1}\| \cdots \|dX_{u_{N+1}}\| \\ &= \Lambda_{N+1}(\mathbf{F}) \frac{d^{N+1}}{(N+1)!} \int_{[0,t]^{N+1}} \|dX_{u_1}\| \cdots \|dX_{u_{N+1}}\| \\ &= \Lambda_{N+1}(\mathbf{F}) \frac{d^{N+1}}{(N+1)!} \Big( \int_0^t \|dX_u\| \Big)^{N+1} \\ &= \Lambda_{N+1}(\mathbf{F}) \frac{d^{N+1}}{(N+1)!} \|X\|_{TV;[0,t]}^{N+1} \leq \Lambda_{N+1}(\mathbf{F}) \frac{d^{N+1}}{(N+1)!}. \end{aligned}$$

#### **B.5 Proof of Proposition 5**

For simplicity of notation, since the context is clear, we now use the notation  $\|\cdot\|$  instead of  $\|\cdot\|_{(\mathbb{R}^e)^{\otimes k}}$ . According to Proposition 1, the solution  $\bar{H}$  of (4) verifies  $\|\bar{H}_t\| \leq M + L := \bar{M}$ . We therefore place ourselves in the ball  $\mathscr{B}_{\bar{M}}$ . Recall that for any  $1 \leq i_1, \ldots, i_N \leq d$ ,  $\bar{h} \in \mathscr{B}_{\bar{M}}$ ,

$$F^{i_1} \star \dots \star F^{i_N}(\bar{h}) = J(F^{i_2} \star \dots \star F^{i_N})(\bar{h})F^{i_1}(\bar{h}).$$
(17)

**Linear case.** We start with the proof of the linear case before moving on to the general case. When  $\sigma$  is chosen to be the identity function, each  $F_{\text{RNN}}^i$  is an affine vector field, in the sense that  $F_{\text{RNN}}^i(\bar{h}) = W_i \bar{h} + b_i$ , where  $W_i = 0_{\bar{e} \times \bar{e}}$ ,  $b_i$  is the i + dth vector of the canonical basis of  $\mathbb{R}^{e+d}$ , and

$$W_{d+1} = \begin{pmatrix} \frac{2}{1-L}W\\ 0_{d\times\bar{e}} \end{pmatrix} \quad \text{and} \quad b_{d+1} = \begin{pmatrix} \frac{2}{1-L}b\\ 0_d \end{pmatrix}.$$

Since  $J(F_{\text{RNN}}^i) = W_i$ , we have, for any  $\bar{h} \in \mathbb{R}^{e+d}$  and any  $1 \le i_1, \ldots, i_k \le d$ ,

$$F_{\text{RNN}}^{i_1} \star \cdots \star F_{\text{RNN}}^{i_k}(\bar{h}) = W_{i_k} \cdots W_{i_2}(W_{i_1}\bar{h} + b_{i_1})$$

Thus, for any  $\bar{h} \in \mathscr{B}_{\bar{M}}$ ,

$$\|F_{\text{RNN}}^{i_1} \star \cdots \star F_{\text{RNN}}^{i_k}(\bar{h})\| \le \|W_{i_k}\|_{\text{op}} \cdots \|W_{i_2}\|_{\text{op}}(\|W_{i_1}\|_{\text{op}}\bar{M} + \|b_{i_1}\|).$$
  
For  $i \ne d+1$ ,  $\|W_{i_1}\|_{\text{op}} = 0$ , and so

$$\Lambda_k(\mathbf{F}_{\mathsf{RNN}}) \le C \|W_{d+1}\|_{\mathsf{op}}^{k-1},$$

with  $C = ||W_{d+1}||_{op} \overline{M} + \max(1, 2(1-L)^{-1}||b||)$ . Therefore,

$$\sum_{k=1}^{\infty} \frac{d^k}{k!} \Lambda_k(\mathbf{F}_{\text{RNN}}) \le Cd \sum_{k=0}^{\infty} \frac{1}{k!} \left( 2d(1-L)^{-1} \|W\|_{\text{op}} \right)^{k-1} < \infty$$

**General case.** In the general case, the proof is two-fold. First, we upper bound (17) by a function of the norms of higher-order Jacobians of  $F^{i_1}, \ldots, F^{i_N}$ . We then apply this bound to the specific case  $\mathbf{F} = \mathbf{F}_{\text{RNN}}$ . We refer to Appendix C for details on higher-order derivatives in tensor spaces. Let  $F : \mathbb{R}^e \to \mathbb{R}^e$  be a smooth vector field. If  $F(h) = (F_1(h), \ldots, F_e(h))^{\top}$ , each of its coordinates  $F_i$  is a function from  $\mathbb{R}^e$  to  $\mathbb{R}$ ,  $\mathscr{C}^{\infty}$  with respect to all its input variables. We define the derivative of order k of F as the tensor field

$$J^{k}(F) : \mathbb{R}^{e} \to (\mathbb{R}^{e})^{\otimes k+1}$$
$$h \mapsto J^{k}(F)(h),$$

where

$$J^{k}(F)(h) = \sum_{1 \leq j, i_{1}, \dots, i_{k} \leq e} \frac{\partial^{k} F_{j}(h)}{\partial h_{i_{1}} \dots \partial h_{i_{k}}} e_{j} \otimes e_{i_{1}} \otimes \dots \otimes e_{i_{k}}.$$

We take the convention  $J^0(F) = F$ , and note that  $J(F) = J^1(F)$  is the Jacobian matrix, and that  $J^k(J^{k'}(F)) = J^{k+k'}(F)$ .

**Lemma 3.** Let  $A^1, \ldots, A^k : \mathbb{R}^e \to \mathbb{R}^e$  be smooth vector fields. Then, for any  $h \in \mathbb{R}^e$ 

$$\left\|A^{k} \star \dots \star A^{1}(h)\right\| \leq \sum_{n_{1} + \dots + n_{k} = k-1} C(k; n_{1}, \dots, n_{k}) \|J^{n_{1}}(A^{1})(h)\| \dots \|J^{n_{k}}(A^{k})(h)\|,$$

where  $C(k; n_1, \ldots, n_k)$  is defined by the following recurrence on k: C(1; 0) = 1 and for any  $n_1, \ldots, n_{k+1} \ge 0$ ,

$$C(k+1; n_1, \dots, n_{k+1}) = \sum_{\ell=1}^{k} C(k; n_1, \dots, n_{\ell} - 1, \dots, n_k)$$
 if  $n_{k+1} = 0$ , (18)  

$$C(k+1; n_1, \dots, n_{k+1}) = 0$$
 otherwise.

*Proof.* We refer to Appendix C for the definitions of the tensor dot product  $\odot$  and tensor permutations, as well as for computation rules involving these operations. We show in fact by induction a stronger result, namely that there exist tensor permutations  $\pi_p$  such that

$$A^{k} \star \dots \star A^{1}(h) = \sum_{n_{1} + \dots + n_{k} = k-1} \sum_{1 \le p \le C(k; n_{1}, \dots, n_{k})} \pi_{p} \left[ J^{n_{1}}(A^{1})(h) \odot \dots \odot J^{n_{k}}(A^{k})(h) \right].$$
(19)

Note that we do not make explicit the permutations nor the axes of the tensor dot operations since we are only interested in bounding the norm of the iterated star products. Also, for simplicity, we denote all permutations by  $\pi$ , even though they may change from line to line.

We proceed by induction on k. For k = 1, the formula is clear. Assume that the formula is true at order k. Then

$$J(A^{n} \star \dots \star A^{r}) = \sum_{n_{1} + \dots + n_{k} = k-1} \sum_{1 \le p \le C(k;n_{1},\dots,n_{k})} J\Big[\pi_{p}[J^{n_{1}}(A^{1}) \odot \dots \odot J^{n_{k}}(A^{k})]\Big]$$
  
= 
$$\sum_{n_{1} + \dots + n_{k} = k-1} \sum_{1 \le p \le C(k;n_{1},\dots,n_{k})} \pi_{p}\Big[J[J^{n_{1}}(A^{1}) \odot \dots \odot J^{n_{k}}(A^{k})]\Big]$$
  
= 
$$\sum_{n_{1} + \dots + n_{k} = k-1} \sum_{1 \le p \le C(k;n_{1},\dots,n_{k})} \sum_{\ell=1}^{k} \pi_{p} \circ \pi_{\ell}\Big[J^{n_{1}}(A^{1}) \odot \dots \odot J^{n_{\ell}}(A^{\ell}) \odot \dots \odot J^{n_{k}}(A^{k})\Big].$$

In the inner sum, we introduce the change of variable  $p_i = n_i$  for  $i \neq \ell$  and  $p_\ell = n_\ell + 1$ . This yields  $J(A^k \star \cdots \star A^1)$ 

$$= \sum_{p_1 + \dots + p_k = k} \sum_{\ell=1}^k \sum_{1 \le p \le C(k; p_1, \dots, p_\ell - 1, \dots, p_k)} \pi_p \circ \pi_\ell \Big[ J^{n_1}(A^1) \odot \cdots \odot J^{n_\ell + 1}(A^\ell) \odot \cdots \odot J^{n_k}(A^k) \Big]$$
$$= \sum_{p_1 + \dots + p_{k+1} = k} \sum_{1 \le q \le C(k+1; p_1, \dots, p_{k+1})} \pi_q \Big[ J^{n_1}(A^1) \odot \cdots \odot J^{p_k}(A^k) \Big],$$

where in the last sum the only non-zero term is for  $p_{k+1} = 0$ . To conclude the induction, it remains to note that

$$A^{k+1} \star \cdots \star A^1 = J(A^k \star \cdots \star A^1) \odot A^{k+1} = J(A^k \star \cdots \star A^1) \odot J^0(A^{k+1}).$$

Hence,

$$A^{k+1} \star \dots \star A^{1} = \sum_{p_{1} + \dots + p_{k+1} = k} \sum_{1 \le q \le C(k+1;p_{1},\dots,p_{k+1})} \pi_{q} \left[ J^{n_{1}}(A^{1}) \odot \dots \odot J^{p_{k}}(A^{k}) \right] \odot J^{p_{k+1}}(A^{k+1}) \\ = \sum_{p_{1} + \dots + p_{k+1} = k} \sum_{1 \le q \le C(k+1;p_{1},\dots,p_{k+1})} \pi_{q} \left[ J^{n_{1}}(A^{1}) \odot \dots \odot J^{p_{k}}(A^{k}) \odot J^{p_{k+1}}(A^{k+1}) \right].$$

The result is then a consequence of (19) and of Lemma 6.

We now restrict ourselves to the case  $\mathbf{F} = \mathbf{F}_{\text{RNN}}$  as defined by (5) and give an upper bound on the higher-order derivatives of the tensor fields  $F^{i_1}, \ldots, F^{i_N}$ .

**Lemma 4.** For any  $i \in \{1, \ldots, d+1\}$ ,  $\bar{h} \in \mathscr{B}_{\bar{M}}$ , for any  $k \ge 0$ ,

$$\|J^k(F^i_{\text{RNN}})(\bar{h})\| \le \left(\frac{2}{1-L}\|W\|_F\right)^k \|\sigma^{(k)}\|_{\infty}.$$

*Proof.* For any  $1 \le i \le d$ ,  $F_{\text{RNN}}^i(\bar{h})$  is constant, so  $J^k(F_{\text{RNN}}^1) = \cdots = J^k(F_{\text{RNN}}^d) = 0$ . For i = d+1, we have, for any  $1 \le j \le e$ ,

$$\frac{\partial^k F^{d+1}_{\mathsf{RNN},j}(\bar{h})}{\partial \bar{h}_{i_1} \dots \partial \bar{h}_{i_k}} = \left(\frac{2}{1-L}\right)^k W_{ji_1} \cdots W_{ji_k} \sigma^{(k)}(W_{j}.\bar{h}+b),$$

where  $W_{j}$  denotes the *j*th row of W and for  $e + 1 \le j \le \bar{e}$ ,  $F_{j}^{d+1} = 0$ . Therefore,

$$\begin{split} \|J^{k}(F_{\text{RNN}}^{d+1})(\bar{h})\|^{2} &\leq \left(\frac{2}{1-L}\right)^{2k} \sum_{1 \leq j, i_{1}, \dots, i_{k} \leq e} |W_{ji_{1}} \cdots W_{ji_{k}} \sigma^{(k)}(W_{j}.\bar{h}+b)|^{2} \\ &= \left(\frac{2}{1-L}\right)^{2k} \|\sigma^{(k)}\|_{\infty}^{2} \sum_{j} \left(\sum_{i} |W_{ji}|^{2}\right)^{k} \\ &\leq \left(\frac{2}{1-L}\right)^{2k} \|\sigma^{(k)}\|_{\infty}^{2} \|W\|_{F}^{2k}. \end{split}$$

We are now in a position to conclude the proof using condition (11). By Lemma 3 and 4, for any  $1 \le i_1, \ldots, i_N \le d+1$ ,

$$\begin{split} \left\| F_{\mathsf{RNN}}^{i_{1}} \star \cdots \star F_{\mathsf{RNN}}^{i_{N}}(\bar{h}) \right\| \\ &\leq \sum_{n_{1}+\dots+n_{N}=N-1} C(N;n_{N},\dots,n_{1}) \| J^{n_{N}}(F_{\mathsf{RNN}}^{i_{N}})(\bar{h}) \| \cdots \| J^{n_{1}}(F_{\mathsf{RNN}}^{i_{1}})(\bar{h}) \| \\ &\leq \left( \frac{2}{1-L} \| W \|_{F} \right)^{N-1} \sum_{n_{1}+\dots+n_{N}=N-1} C(N;n_{N},\dots,n_{1}) a^{n_{1}+1} n_{1}! \cdots a^{n_{N}+1} n_{N}! \\ &\leq a \left( \frac{2}{1-L} a^{2} \| W \|_{F} \right)^{N-1} \sum_{n_{1}+\dots+n_{N}=N-1} C(N;n_{N},\dots,n_{1}) n_{1}! \cdots n_{N}! \,. \end{split}$$

Assume for the moment that  $C(N; n_N, \ldots, n_1)$  is smaller than the multinomial coefficient  $\binom{N}{n_N, \ldots, n_1}$ . Then, using the fact that there are  $\binom{n+k-1}{k-1}$  weak compositions of n in k parts and Stirling's approximation, we have

$$\begin{split} \Lambda_{N}(\mathbf{F}) &\leq a \Big( \frac{2}{1-L} a^{2} \|W\|_{F} \Big)^{N-1} N! \times \operatorname{Card} \big( \{ n_{1} + \dots + n_{N} = N-1 \} \big) \\ &\leq a \Big( \frac{2}{1-L} a^{2} \|W\|_{F} \Big)^{N-1} N! \binom{2N-2}{N-1} \\ &\leq \frac{a}{2} \Big( \frac{2}{1-L} a^{2} \|W\|_{F} \Big)^{N-1} N! \binom{2N}{N} \\ &\leq a \frac{\sqrt{2}e}{\pi} \Big( \frac{8}{1-L} a^{2} \|W\|_{F} \Big)^{N-1} \frac{N!}{\sqrt{N}}. \end{split}$$

Hence, provided  $||W||_F < (1-L)/8a^2d$ ,

$$\sum_{k=1}^{\infty} \frac{d^k}{k!} \Lambda_k(\mathbf{F}) \le a d \frac{\sqrt{2}e}{\pi} \sum_{k=1}^{\infty} \left(\frac{8da^2 \|W\|_F}{1-L}\right)^{k-1} \frac{1}{\sqrt{k}} < \infty,$$

and  $(A_2)$  is verified.

To conclude the proof, it remains to prove the following lemma.

**Lemma 5.** For any  $k \ge 1$  and  $n_1, \ldots, n_k \ge 0$ ,  $C(k; n_1, \ldots, n_k) \le {\binom{k-1}{n_1, \ldots, n_k}}$ .

*Proof.* The proof is done by induction, by comparing the recurrence formula (18) with the following recurrence formula for multinomial coefficients:

$$\binom{k}{n_1, \dots, n_{k+1}} = \sum_{\ell=1}^{k+1} \binom{k-1}{n_1, \dots, n_\ell - 1, \dots, n_{k+1}}$$

More precisely, for k = 1,  $C(1;0) = 1 \le {0 \choose 0} = 1$  and  $C(1;1) = 0 \le {0 \choose 1} = 0$ . Assume that the formula is true at order k. Then, at order k + 1, there are two cases. If  $n_{k+1} \ne 0$ ,  $C(k+1;n_1,\ldots,n_{k+1}) = 0$ , and the result is clear. On the other hand, if  $n_{k+1} = 0$ ,

$$C(k+1; n_1, \dots, n_k, 0) = \sum_{\ell=1}^k C(k; n_1, \dots, n_\ell - 1, \dots, n_k)$$
  
$$\leq \sum_{\ell=1}^k \binom{k-1}{n_1, \dots, n_\ell - 1, \dots, n_k}$$
  
$$\leq \sum_{\ell=1}^{k+1} \binom{k-1}{n_1, \dots, n_\ell - 1, \dots, n_{k+1}}$$
  
$$\leq \binom{k}{n_1, \dots, n_{k+1}}.$$

#### **B.6** Proof of Theorem 1

First, Propositions 1 and 2 state that if  $\overline{H}$  is the solution of (4) and Proj denotes the projection on the first *e* coordinates, then

$$\left|z_{T}-\psi\left(\operatorname{Proj}(\bar{H}_{1})\right)\right|=\left|\psi(h_{T})-\psi\left(\operatorname{Proj}(\bar{H}_{1}]\right)\right)\right|\leq \left\|\psi\right\|_{\operatorname{op}}\left\|h_{T}-\operatorname{Proj}(\bar{H}_{1})\right\|\leq \left\|\psi\right\|_{\operatorname{op}}\frac{c_{1}}{T}$$

For any  $1 \le k \le N$ , we let  $\mathscr{D}^k(\bar{H}_0) : (\mathbb{R}^d)^{\otimes k} \to \mathbb{R}^e$  be the linear function defined by

$$\mathscr{D}^{k}(\bar{H}_{0})(e_{i_{1}}\otimes\cdots\otimes e_{i_{k}})=F^{i_{1}}\star\cdots\star F^{i_{k}}(\bar{H}_{0}),$$
(20)

where  $e_1, \ldots, e_d$  denotes the canonical basis of  $\mathbb{R}^{\bar{d}}$ . Then, under assumptions  $(A_1)$  and  $(A_2)$ , if  $\bar{\mathbb{X}}^k$  denotes the signature of order k of the path  $\bar{X}_t = (X_t^{\top}, \frac{1-L}{2}t)^{\top}$ , according to Propositions 4 and 5,

$$\bar{H}_1 = \bar{H}_0 + \sum_{k=1}^{\infty} \frac{1}{k!} \sum_{1 \le i_1, \dots, i_k \le d} S^{(i_1, \dots, i_k)}_{[0,t]}(X) F^{i_1} \star \dots \star F^{i_k}(\bar{H}_0) = \sum_{k=1}^{\infty} \frac{1}{k!} \mathscr{D}^k(\bar{H}_0)(\mathbb{X}^k_{[0,t]}),$$

and

$$\psi \circ \operatorname{Proj}(\bar{H}_1) = \psi \circ \operatorname{Proj}\Big(\sum_{k=0}^{\infty} \frac{1}{k!} \mathscr{D}^k(\bar{H}_0)(\bar{\mathbb{X}}^k)\Big) = \sum_{k=0}^{\infty} \frac{1}{k!} \psi \circ \operatorname{Proj}\big(\mathscr{D}^k(\bar{H}_0)(\bar{\mathbb{X}}^k)\big),$$

by linearity of  $\psi$  and Proj. Since the maps  $\mathscr{D}^k(\bar{H}_0) : (\mathbb{R}^d)^{\otimes k} \to \mathbb{R}^e$  are linear, the above equality takes the form

$$\psi \circ \operatorname{Proj}(\bar{H}_1) = \sum_{k=0}^{\infty} \langle \alpha^k, \bar{\mathbb{X}}^k \rangle_{(\mathbb{R}^d)^{\otimes k}},$$
(21)

where  $\alpha^k \in (\mathbb{R}^d)^{\otimes k}$  is the coefficient of the linear map  $\frac{1}{k!}\psi \circ \operatorname{Proj} \circ \mathscr{D}^k(\bar{H}_0)$  in the canonical basis. Let  $\alpha = (\alpha^0, \ldots, \alpha^k, \ldots)$ . Under assumption  $(A_2)$ ,

$$\begin{split} \sum_{k=0}^{\infty} \|\alpha^{k}\|_{(\mathbb{R}^{d})^{\otimes k}}^{2} &\leq \sum_{k=0}^{\infty} \sum_{1 \leq i_{1}, \dots, i_{k} \leq d} \left(\frac{1}{k!}\right)^{2} \|\psi\|_{\mathrm{op}}^{2} \|F^{i_{1}} \star \dots \star F^{i_{k}}(\bar{H}_{0})\|^{2} \\ &\leq \|\psi\|_{\mathrm{op}}^{2} \sum_{k=0}^{\infty} \sum_{1 \leq i_{1}, \dots, i_{k} \leq d} \left(\frac{1}{k!}\right)^{2} \Lambda_{k}(\mathbf{F})^{2} \\ &\leq \|\psi\|_{\mathrm{op}}^{2} \sum_{k=0}^{\infty} \left(\frac{d^{k}}{k!} \Lambda_{k}(\mathbf{F})\right)^{2} < \infty. \end{split}$$

This shows that  $\alpha \in \mathscr{T}$ , and therefore, using (21), we conclude

$$||z_T - \langle \alpha, S(\bar{X}) \rangle_{\mathscr{T}}|| \le ||\psi||_{\text{op}} \frac{c_1}{T}.$$

# B.7 Proof of Theorem 2

Let

$$\mathscr{G} = \left\{ g_{\theta} : (\mathbb{R}^d)^T \to \mathbb{R} \, | \, g_{\theta}(\mathbf{x}) = z_T, \theta \in \Theta \right\}$$

be the function class of (discrete) RNN and

$$\mathscr{S} = \Big\{ \xi_{\alpha_{\theta}} : \mathscr{X} \to \mathbb{R} \, | \, \xi_{\alpha_{\theta}}(X) = \langle \alpha_{\theta}, S(\bar{X}) \rangle_{\mathscr{T}}, \theta \in \Theta \Big\},\$$

be the class of their RKHS embeddings, where  $\alpha_{\theta}$  is defined by (21). For any  $\theta \in \Theta$ , we let

$$\mathscr{R}_{\mathscr{G}}(\theta) = \mathbb{E}[\ell(\mathbf{y}, g_{\theta}(\mathbf{x}))], \quad \text{ and } \quad \mathscr{R}_{\mathscr{S}}(\theta) = \mathbb{E}[\ell(\mathbf{y}, \xi_{\alpha_{\theta}}(\bar{X}))],$$

and denote by  $\widehat{\mathscr{R}}_{n,\mathscr{G}}$  and  $\widehat{\mathscr{R}}_{n,\mathscr{G}}$  the corresponding empirical risks. We also let  $\theta_{\mathscr{G}}^*$ ,  $\theta_{\mathscr{S}}^*$ ,  $\widehat{\theta}_{n,\mathscr{G}}$ , and  $\widehat{\theta}_{n,\mathscr{G}}$  be the corresponding minimizers. We have

$$\begin{split} \mathbb{P} \big( \mathbf{y} \neq g_{\widehat{\theta}_{n,\mathscr{G}}}(\mathbf{x}) \big) - \widehat{\mathscr{R}}_{n,\mathscr{G}}(\widehat{\theta}_{n,\mathscr{G}}) &\leq \mathbb{E} \big[ \ell(\mathbf{y}, g_{\widehat{\theta}_{n,\mathscr{G}}}(\mathbf{x})) \big] - \widehat{\mathscr{R}}_{n,\mathscr{G}}(\widehat{\theta}_{n,\mathscr{G}}) \\ &= \mathscr{R}_{\mathscr{G}}(\widehat{\theta}_{n,\mathscr{G}}) - \widehat{\mathscr{R}}_{n,\mathscr{G}}(\widehat{\theta}_{n,\mathscr{G}}) \\ &= \mathscr{R}_{\mathscr{G}}(\widehat{\theta}_{n,\mathscr{G}}) - \mathscr{R}_{\mathscr{F}}(\widehat{\theta}_{n,\mathscr{G}}) + \mathscr{R}_{\mathscr{F}}(\widehat{\theta}_{n,\mathscr{G}}) - \widehat{\mathscr{R}}_{n,\mathscr{F}}(\widehat{\theta}_{n,\mathscr{G}}) \\ &+ \widehat{\mathscr{R}}_{n,\mathscr{F}}(\widehat{\theta}_{n,\mathscr{G}}) - \widehat{\mathscr{R}}_{n,\mathscr{G}}(\widehat{\theta}_{n,\mathscr{G}}) \\ &\leq \sup_{\theta} |\mathscr{R}_{\mathscr{G}}(\theta) - \mathscr{R}_{\mathscr{F}}(\theta)| + \sup_{\theta} |\mathscr{R}_{\mathscr{F}}(\theta) - \widehat{\mathscr{R}}_{n,\mathscr{F}}(\theta)| \\ &+ \sup_{\theta} |\widehat{\mathscr{R}}_{n,\mathscr{G}}(\theta) - \widehat{\mathscr{R}}_{n,\mathscr{F}}(\theta)|. \end{split}$$

Using Theorem 1, we have

$$\begin{split} \sup_{\theta} |\mathscr{R}_{\mathscr{G}}(\theta) - \mathscr{R}_{\mathscr{S}}(\theta)| &= \sup_{\theta} |\mathbb{E} \left[ \ell(\mathbf{y}, g_{\theta}(\mathbf{x})) - \ell(\mathbf{y}, \xi_{\alpha_{\theta}}(\bar{X})) \right] | \\ &\leq \sup_{\theta} \mathbb{E} \left[ |\phi(\mathbf{y}g_{\theta}(\mathbf{x})) - \phi(\mathbf{y}\xi_{\alpha_{\theta}}(\bar{X}))| \right] \\ &\leq \sup_{\theta} \mathbb{E} \left[ K_{\ell} |\mathbf{y}| \times |g_{\theta}(\mathbf{x}) - \xi_{\alpha_{\theta}}(\bar{X})| \right] \\ &\leq K_{\ell} \sup_{\theta} (\|\psi\|_{\text{op}} c_{1,\theta}) \frac{1}{T} := \frac{c_{2}}{2T}, \end{split}$$

where  $c_{1,\theta} = K_{f_{\theta}} e^{K_{f_{\theta}}} (L + ||f_{\theta}||_{\infty} e^{K_{f_{\theta}}})$  (the infinity norm  $||f_{\theta}||_{\infty}$  is taken on the balls  $\mathscr{B}_L$  and  $\mathscr{B}_M$ ). One proves with similar arguments that

$$\sup_{\theta} |\widehat{\mathscr{R}}_{n,\mathscr{G}}(\theta) - \widehat{\mathscr{R}}_{n,\mathscr{S}}(\theta)| \le \frac{c_2}{2T}.$$

Under the assumption of the theorem, there exists a ball  $\mathscr{B} \subset \mathscr{H}$  of radius B such that  $\mathscr{S} \subset \mathscr{B}$ . This yields

$$\sup_{\theta} |\mathscr{R}_{\mathscr{S}}(\theta) - \widehat{\mathscr{R}}_{n,\mathscr{S}}(\theta)| \leq \sup_{\alpha \in \mathscr{T}, \|\alpha\|_{\mathscr{T}} \leq B} |\mathscr{R}_{\mathscr{B}}(\alpha) - \widehat{\mathscr{R}}_{n,\mathscr{B}}(\alpha)|,$$

where

$$\mathscr{R}_{\mathscr{B}}(\alpha) = \mathbb{E}[\ell(Y,\xi_{\alpha}(\bar{X}))] \quad \text{and} \quad \widehat{\mathscr{R}}_{n,\mathscr{B}}(\alpha) = \frac{1}{n} \sum_{i=1}^{n} \ell(Y^{(i)},\xi_{\alpha}(\bar{X}^{(i)}))$$

We now have reached a familiar situation where the supremum is over a ball in an RKHS. A slight extension of Bartlett and Mendelson (2002, Theorem 8) yields that with probability at least  $1 - \delta$ ,

$$\sup_{\alpha \in \mathscr{T}, \|\alpha\|_{\mathscr{T}} \le B} |\mathscr{R}_{\mathscr{B}}(\alpha) - \widehat{\mathscr{R}}_{n,\mathscr{B}}(\alpha)| \le 4K_{\ell} \mathbb{E} \operatorname{Rad}_{n}(\mathscr{B}) + 2BK_{\ell}(1-L)^{-1} \sqrt{\frac{\log(1/\delta)}{2n}},$$

where  $\operatorname{Rad}_n(\mathscr{B})$  denotes the Rademacher complexity of  $\mathscr{B}$ . Observe that we have used the fact that the loss is bounded by  $2BK_\ell(1-L)^{-1}$  since, for any  $\xi_\alpha \in \mathscr{B}$ , by the Cauchy-Schwartz inequality,

$$\ell(\mathbf{y},\xi_{\alpha}(X)) = \phi(\mathbf{y}\langle\alpha,S(X)\rangle_{\mathscr{T}}) \leq K_{\ell}|\mathbf{y}\langle\alpha,S(X)\rangle_{\mathscr{T}}| \leq K_{\ell}\|\alpha\|_{\mathscr{T}}\|S(X)\|_{\mathscr{T}}$$
$$\leq 2K_{\ell}B(1-L)^{-1}.$$

Finally, the proof follows by noting that Rademacher complexity of  $\mathcal{B}$  is bounded by

$$\operatorname{Rad}_{n}(\mathscr{B}) \leq \frac{B}{n} \sqrt{\sum_{i=1}^{n} K(X^{(i)}, X^{(i)})} = \frac{B}{n} \sqrt{\sum_{i=1}^{n} \|S(\bar{X}^{(i)})\|_{\mathscr{T}}^{2}} \leq \frac{2B(1-L)^{-1}}{\sqrt{n}}$$

#### **B.8** Proof of Theorem 3

Let

$$\mathscr{G} = \left\{ g_{\theta} : (\mathbb{R}^d)^T \to (\mathbb{R}^p)^T \, | \, g_{\theta}(\mathbf{x}) = (z_1, \dots, z_T), \theta \in \Theta \right\}$$

be the function class of discrete RNN in a sequential setting. Let

$$\mathscr{S} = \left\{ \Gamma_{\theta} : \mathscr{X} \to (\mathbb{R}^p)^T \,|\, \Gamma_{\theta}(X) = \left( \Xi_{\theta}(\tilde{X}_{[1]}), \dots, \Xi_{\theta}(\tilde{X}_{[T]}) \right) \right\},\,$$

be the class of their RKHS embeddings, where  $\tilde{X}_{[j]}$  is the path equal to X on [0, j/T] and then constant on [j/T, 1] (see Figure 4). For any  $X \in \mathscr{X}$ ,

$$\Xi_{\theta}(a) = \begin{pmatrix} \langle \alpha_{1,\theta}, S(\bar{X}) \rangle_{\mathscr{T}} \\ \vdots \\ \langle \alpha_{p,\theta}, S(\bar{X}) \rangle_{\mathscr{T}} \end{pmatrix} = \begin{pmatrix} \xi_{\alpha_{1,\theta}}(X) \\ \vdots \\ \xi_{\alpha_{p,\theta}}(X) \end{pmatrix} \in \mathbb{R}^{p},$$

where  $(\alpha_{1,\theta},\ldots,\alpha_{p,\theta})^{\top} \in (\mathscr{T})^p$  are the coefficients of the linear maps  $\frac{1}{k!}\psi \circ \operatorname{Proj} \circ \mathscr{D}^k(\bar{H}_0)$ :  $(\mathbb{R}^d)^{\otimes k} \to \mathbb{R}^p, k \geq 0$ , in the canonical basis, where  $\mathscr{D}^k$  is defined by (20).

We start the proof as in Theorem 2, until we obtain

$$\begin{aligned} \mathscr{R}_{\mathscr{G}}(\widehat{\theta}_{n,\mathscr{G}}) - \widehat{\mathscr{R}}_{n,\mathscr{G}}(\widehat{\theta}_{n,\mathscr{G}}) &\leq \sup_{\theta} |\mathscr{R}_{\mathscr{G}}(\theta) - \mathscr{R}_{\mathscr{S}}(\theta)| + \sup_{\theta} |\mathscr{R}_{\mathscr{S}}(\theta) - \widehat{\mathscr{R}}_{n,\mathscr{S}}(\theta)| \\ &+ \sup_{\theta} |\widehat{\mathscr{R}}_{n,\mathscr{G}}(\theta) - \widehat{\mathscr{R}}_{n,\mathscr{S}}(\theta)|. \end{aligned}$$

By definition of the loss, for any  $\theta \in \Theta$ ,

$$\begin{aligned} |\mathscr{R}_{\mathscr{G}}(\theta) - \mathscr{R}_{\mathscr{S}}(\theta)| &= \left| \mathbb{E} \left[ \ell \left( \mathbf{y}, g_{\theta}(\mathbf{x}) \right) - \ell \left( \mathbf{y}, \Gamma_{\theta}(X) \right) \right] \right| \\ &\leq \mathbb{E} \left[ \left| \frac{1}{T} \sum_{j=1}^{T} \left( \|y_j - z_j\|^2 - \|y_j - \Xi_{\theta}(\tilde{X}_{[j]})\|^2 \right) \right| \right] \\ &\leq \mathbb{E} \left[ \frac{1}{T} \sum_{j=1}^{T} \left| \left\langle z_j + \Xi_{\theta}(\tilde{X}_{[j]}) - 2y_j, z_j - \Xi_{\theta}(\tilde{X}_{[j]}) \right\rangle \right| \right] \\ &\leq \mathbb{E} \left[ \frac{1}{T} \sum_{j=1}^{T} \|z_j + \Xi_{\theta}(\tilde{X}_{[j]}) - 2y_j\| \times \|z_j - \Xi_{\theta}(\tilde{X}_{[j]})\| \right] \end{aligned}$$

(by the Cauchy-Schwartz inequality).

According to inequality (14), one has

$$\|z_j - \Xi_{\theta}(\tilde{X}_{[j]})\| \le \|\psi\|_{\operatorname{op}} \frac{c_{1,\theta}}{T},$$

where  $c_{1,\theta} = K_{f_{\theta}} e^{K_{f_{\theta}}} \left( L + \|f_{\theta}\|_{\infty} e^{K_{f_{\theta}}} \right)$ . Moreover,

$$\left\|\Xi_{\theta}(\tilde{X}_{[j]})\right\|^{2} = \sum_{\ell=1}^{p} \left|\langle \alpha_{\ell,\theta}, S(\tilde{X}_{[j]})\rangle_{\mathscr{T}}\right|^{2} \le \sum_{\ell=1}^{p} \|\alpha_{\ell,\theta}\|_{\mathscr{T}}^{2} \|S(\tilde{X}_{[j]})\|_{\mathscr{T}}^{2} \le pB^{2} \left(2(1-L)^{-1}\right)^{2},$$

since  $\|S(\tilde{X}_{[j]})\|_{\mathscr{T}} = \|S_{[0,j/T]}(\bar{X})\|_{\mathscr{T}} \le \|S(\bar{X})\|_{\mathscr{T}}$ . This yields

$$\begin{aligned} \|z_j + \Xi_{\theta}(\dot{X}_{[j]}) - 2y_j\| &\leq \|z_j\| + \|\Xi_{\theta}(\dot{X}_{[j]})\| + 2\|y_j\| \\ &\leq \|\psi\|_{\text{op}} \|f_{\theta}\|_{\infty} + 2\sqrt{p}B(1-L)^{-1} + 2K_y \end{aligned}$$

Finally,

$$\sup_{\theta} |\mathscr{R}_{\mathscr{G}}(\theta) - \mathscr{R}_{\mathscr{S}}(\theta)| \le \frac{c_3}{2T},$$

where  $c_3 = \sup_{\theta} (c_{1,\theta} + \|\psi\|_{\text{op}} \|f_{\theta}\|_{\infty}) + 2\sqrt{p}B(1-L)^{-1} + 2K_y$ . One proves with similar arguments that

$$\sup_{\theta} |\widehat{\mathscr{R}}_{n,\mathscr{G}}(\theta) - \widehat{\mathscr{R}}_{n,\mathscr{S}}(\theta)| \le \frac{c_3}{2T}.$$

We now turn to the term  $\sup_{\theta} |\mathscr{R}_{\mathscr{S}}(\theta) - \widehat{\mathscr{R}}_{n,\mathscr{S}}(\theta)|.$  We have

$$\begin{aligned} \mathscr{R}_{\mathscr{S}}(\theta) &- \widehat{\mathscr{R}}_{n,\mathscr{S}}(\theta) \\ &= \mathbb{E}[\ell(\mathbf{y}, \Gamma_{\theta}(X))] - \frac{1}{n} \sum_{i=1}^{n} \ell(\mathbf{y}^{(i)}, \Gamma_{\theta}(X^{(i)})) \\ &= \frac{1}{T} \sum_{j=1}^{T} \left( \mathbb{E}[\|y_{j} - \Xi_{\theta}(\tilde{X}_{[j]})]\|^{2} - \frac{1}{n} \sum_{i=1}^{n} \|y_{j}^{(i)} - \Xi_{\theta}(\tilde{X}_{[j]}^{(i)})\|^{2} \right) \end{aligned}$$

Therefore,

$$\sup_{\theta} |\mathscr{R}_{\mathscr{S}}(\theta) - \widehat{\mathscr{R}}_{n,\mathscr{S}}(\theta)| \le \frac{1}{T} \sum_{j=1}^{T} \sup_{\theta} \Big| \mathbb{E}[\|y_j - \Xi_{\theta}(\tilde{X}_{[j]})]\|^2 - \frac{1}{n} \sum_{i=1}^{n} \|y_j^{(i)} - \Xi_{\theta}(\tilde{X}_{[j]}^{(i)})\|^2 \Big|.$$

Note that for a fixed j, the pairs  $(\tilde{X}_{[j]}^{(i)}, y_j^{(i)})$  are i.i.d. Under the assumptions of the theorem, there exists a ball  $\mathscr{B} \subset \mathscr{H}$  such that for any  $1 \leq \ell \leq p, \theta \in \Theta, \xi_{\alpha_{\ell,\theta}} \in \mathscr{B}$ . We denote by  $\mathscr{B}_p$  the sum of p such spaces, that is,

$$\mathscr{B}_p = \left\{ f_\alpha : \mathscr{X} \to \mathbb{R}^p \,|\, f_\alpha(X) = (f_{\alpha_1}(X), \dots, f_{\alpha_p}(X))^\top, f_{\alpha_\ell} \in \mathscr{B} \right\}.$$

Clearly,  $\Xi_{\theta} \in \mathscr{B}_p$ , and it follows that

$$\begin{split} \sup_{\theta} & \left| \mathbb{E}[\|y_j - \Xi_{\theta}(\tilde{X}_{[j]})]\|^2 - \frac{1}{n} \sum_{i=1}^n \|y_j^{(i)} - \Xi_{\theta}(\tilde{X}_{[j]}^{(i)})\|^2 \right| \\ & \leq \sup_{f_{\alpha} \in \mathscr{B}_p} \left| \mathbb{E}\big[\|y_j - f_{\alpha}(\tilde{X}_{[j]})\|^2\big] - \frac{1}{n} \sum_{i=1}^n \|y_j^{(i)} - f_{\alpha}(\tilde{X}_{[j]}^{(i)})\|^2 \big] \end{split}$$

We have once again reached a familiar situation, which can be dealt with by an easy extension of Bartlett and Mendelson (2002, Theorem 12). For any  $f_{\alpha} \in \mathscr{B}_p$ , let  $\tilde{\phi} \circ f_{\alpha} : \mathscr{X} \times \mathbb{R}^p : (X, y) \mapsto$  $\|y - f_{\alpha}(X)\|^2 - \|y\|^2$ . Then,  $\tilde{\phi} \circ f_{\alpha}$  is upper bounded by

$$\begin{split} |\tilde{\phi} \circ f_{\alpha}(X,y)| &= \left| \|y - f_{\alpha}(X)\|^{2} - \|y\|^{2} \right| \leq \|f_{\alpha}(X)\| \left( \|f_{\alpha}(X)\| + 2\|y\| \right) \\ &\leq 2\sqrt{p}B(1-L)^{-1}(2\sqrt{p}B(1-L)^{-1} + 2K_{y}) \\ &\leq 4pB(1-L)^{-1}(B(1-L)^{-1} + K_{y}). \end{split}$$

Let  $c_4 = B(1-L)^{-1} + K_y$  and  $c_5 = 4pB(1-L)^{-1}c_4 + K_y^2$ . Then with probability at least  $1 - \delta$ ,

$$\sup_{f_{\alpha}\in\mathscr{B}_{p}}\left|\mathbb{E}\left[\|y_{j}-f_{\alpha}(\tilde{X}_{[j]})\|\right]-\frac{1}{n}\sum_{i=1}^{n}\|y_{j}^{(i)}-f_{\alpha}(\tilde{X}_{[j]}^{(i)})\|\right|\leq \operatorname{Rad}_{n}(\tilde{\phi}\circ\mathscr{B}_{p})+\sqrt{\frac{2c_{5}\log(1/\delta)}{n}},$$

where  $\tilde{\phi} \circ \mathscr{B}_p = \{(X, y) \mapsto \tilde{\phi} \circ f_{\alpha}(X, y) | f_{\alpha} \in \mathscr{B}_p\}$ . Elementary computations on Rademacher complexities yield

$$\operatorname{Rad}_n(\tilde{\phi} \circ \mathscr{B}_p) \leq 2pc_4 \operatorname{Rad}_n(\mathscr{B}) \leq \frac{4pc_4 B(1-L)^{-1}}{\sqrt{n}},$$

which concludes the proof.

# C Differentiation with higher-order tensors

## C.1 Definition

We define the generalization of matrix product between square tensors of order k and  $\ell$ .

**Definition 4.** Let  $a \in (\mathbb{R}^e)^{\otimes k}$ ,  $b \in (\mathbb{R}^e)^{\otimes \ell}$ ,  $p \in \{1, \ldots, k\}$ ,  $q \in \{1, \ldots, \ell\}$ . Then the tensor dot product along (p, q), denoted by  $a \odot_{p,q} b \in (\mathbb{R}^e)^{\otimes (k+\ell-2)}$ , is defined by

$$(a \odot_{p,q} b)_{(i_1,\dots,i_{k-1},j_1,\dots,j_{\ell-1})} = \sum_{j=1}^c a_{(i_1,\dots,i_{p-1},j,i_p,\dots,i_{k-1})} b_{(j_1,\dots,j_{q-1},j,j_q,\dots,j_{\ell-1})}.$$

This operation just consists in computing  $a \otimes b$ , and then summing the *p*th coordinate of *a* with the *q*th coordinate of *b*. The  $\odot$  operator is not associative. To simplify notation, we take the convention that it is evaluated from left to right, that is, we write  $a \odot b \odot c$  for  $(a \odot b) \odot c$ .

**Definition 5.** Let  $a \in (\mathbb{R}^e)^{\otimes k}$ . For a given permutation  $\pi$  of  $\{1, \ldots, k\}$ , we denote by  $\pi(a)$  the permuted tensor in  $(\mathbb{R}^e)^{\otimes k}$  such that

$$\pi(a)_{(i_1,\dots,i_k)} = a_{(i_{\Pi(1)},\dots,i_{\Pi(k)})}.$$

**Example 5.** If A is a matrix, then  $A^T = \pi(A)$ , with  $\pi$  defined by  $\pi(1) = 2, \pi(2) = 1$ .

## C.2 Computation rules

We need to obtain two computation rules for the tensor dot product: bounding the norm (Lemma 6) and differentiating (Lemma 7).

**Lemma 6.** Let  $a \in (\mathbb{R}^e)^{\otimes k}$ ,  $b \in (\mathbb{R}^e)^{\otimes \ell}$ . Then, for all p, q,

$$\|a \odot_{p,q} b\|_{(\mathbb{R}^e)^{\otimes k+\ell-2d}} \le \|a\|_{(\mathbb{R}^e)^{\otimes k}} \|b\|_{(\mathbb{R}^e)^{\otimes \ell}}.$$

#### Proof. By the Cauchy-Schwartz inequality,

$$\begin{split} |a \odot_{p,q} b||_{(\mathbb{R}^{e})^{\otimes k+\ell-2}}^{2} &= \sum_{1 \leq i_{1}, \dots, i_{k-1}, j_{1}, \dots, j_{\ell-1} \leq e} (a \odot_{p,q} b)_{(i_{1}, \dots, i_{k-1}, j_{1}, \dots, j_{\ell-1})}^{2} \\ &= \sum_{1 \leq i_{1}, \dots, i_{k-1}, j_{1}, \dots, j_{\ell-1} \leq e} \left(\sum_{1 \leq j \leq e} a_{(i_{1}, \dots, i_{p-1}, j, i_{p}, \dots, i_{k-1})} b_{(j_{1}, \dots, j_{q-1}, j, j_{q}, \dots, j_{\ell-1})}\right)^{2} \\ &\leq \sum_{i_{1}, \dots, i_{k-1}, j_{1}, \dots, j_{\ell-1}} \left(\sum_{j} a_{(i_{1}, \dots, i_{p-1}, j, i_{p}, \dots, i_{k-1})}\right) \left(\sum_{j} b_{(j_{1}, \dots, j_{q-1}, j, j_{q}, \dots, j_{\ell-1})}\right) \\ &\leq \sum_{i_{1}, \dots, i_{k-1}, j} a_{(i_{1}, \dots, i_{p-1}, j, i_{p}, \dots, i_{k-1})} \sum_{j_{1}, \dots, j_{\ell-1}, j} b_{(j_{1}, \dots, j_{q-1}, j, j_{q}, \dots, j_{\ell-1})} \\ &\leq \|a\|_{(\mathbb{R}^{e})^{\otimes k}}^{2} \|b\|_{(\mathbb{R}^{e})^{\otimes \ell}}^{2}. \end{split}$$

**Lemma 7.** Let  $A : \mathbb{R}^e \to (\mathbb{R}^e)^{\otimes k}$ ,  $B : \mathbb{R}^e \to (\mathbb{R}^e)^{\otimes \ell}$  be smooth vector fields,  $p \in \{1, \ldots, k\}$ ,  $q \in \{1, \ldots, \ell\}$ . Let  $A \odot_{p,q} B : \mathbb{R}^e \to (\mathbb{R}^e)^{\otimes k+\ell-2}$  be defined by  $A \odot_{p,q} B(h) = A(h) \odot_{p,q} B(h)$ . Then there exists a permutation  $\pi$  such that

$$J(A \odot_{p,q} B) = \pi(J(A) \odot_{p,q} B) + A \odot_{p,q} J(B).$$

Proof. The left-hand side takes the form

$$(J(A \odot_{p,q} B))_{i_1,\dots,i_{k-1},j_1,\dots,j_{\ell-1},m} = \sum_j \Big[ \frac{\partial A}{\partial h_m} \Big[_{(i_1,\dots,i_{p-1},j,i_p,\dots,i_{k-1})} B_{(j_1,\dots,j_{q-1},j,j_q,\dots,j_{\ell-1})} + A_{(i_1,\dots,i_{p-1},j,i_p,\dots,i_{k-1})} \frac{\partial B}{\partial h_m} \Big]_{(j_1,\dots,j_{q-1},j,j_q,\dots,j_{\ell-1})} \Big].$$

The first term of the right-hand side writes

$$(J(A) \odot_{p,q} B)_{i_1,\dots,i_{k-1},m,j_1,\dots,j_{\ell-1}} = \sum_j \Big[ \frac{\partial A}{\partial h_m} \Big|_{(i_1,\dots,i_{p-1},j,i_p,\dots,i_{k-1})} B_{(j_1,\dots,j_{q-1},j,j_q,\dots,j_{\ell-1})} \Big],$$

and the second one

$$(A \odot_{p,q} J(B))_{i_1,\dots,i_{k-1},j_1,\dots,j_{\ell-1},m} = \sum_j \Big[ A_{(i_1,\dots,i_{p-1},j,i_p,\dots,i_{k-1})} \frac{\partial B}{\partial h_m}_{(j_1,\dots,j_{q-1},j,j_q,\dots,j_{\ell-1})} \Big].$$

Let us introduce the permutation  $\pi$  which keeps the first (k-1) axes unmoved, and rotates the remaining  $\ell$  ones such that the last axis ends up in kth position. Then

$$\pi(J(A) \odot_{p,q} B)_{i_1,\dots,i_{k-1},j_1,\dots,j_{\ell-1},m} = \sum_j \Big[ \frac{\partial A}{\partial h_m} \Big[_{(i_1,\dots,i_{p-1},j,i_p,\dots,i_{k-1})} B_{(j_1,\dots,j_{q-1},j,j_q,\dots,j_{\ell-1})} \Big].$$

Hence  $J(A \odot_{p,q} B) = \pi(J(A) \odot_{p,q} B) + A \odot_{p,q} J(B)$ , which concludes the proof.

The following two lemmas show how to compose the Jacobian and the tensor dot operations with permutations. Their proofs follow elementary operations and are therefore omitted.

**Lemma 8.** Let  $A : \mathbb{R}^e \to (\mathbb{R}^e)^{\otimes k}$  and  $\pi$  a permutation of  $\{1, \ldots, k\}$ . Then there exists a permutation  $\tilde{\pi}$  of  $\{1, \ldots, k+1\}$  such that

$$J(\pi(A)) = \tilde{\pi}(J(A)).$$

**Lemma 9.** Let  $a \in (\mathbb{R}^e)^{\otimes k}$ ,  $b \in (\mathbb{R}^e)^{\otimes \ell}$ ,  $p \in \{1, \ldots, k\}$ ,  $q \in \{1, \ldots, \ell\}$ ,  $\pi$  a permutation of  $\{1, \ldots, k\}$ . Then there exists  $\tilde{p} \in \{1, \ldots, k\}$ ,  $\tilde{q} \in \{1, \ldots, \ell\}$ , and a permutation  $\tilde{\pi}$  of  $\{1, \ldots, k + \ell - 2\}$  such that

$$\pi(a) \odot_{p,q} b = \tilde{\pi}(a \odot_{\tilde{p},\tilde{q}} b).$$

The following result is a generalization of Lemma 7 to the case of a dot product of several tensors.

**Lemma 10.** For  $\ell \in \{1, \ldots, k\}$ ,  $n_{\ell} \in \mathbb{N}$ , let  $A_{\ell} : \mathbb{R}^e \to (\mathbb{R}^e)^{\otimes n_{\ell}}$  be smooth tensor fields. For any  $(p_{\ell})_{1 \leq \ell \leq k-1}$  and  $(q_{\ell})_{1 \leq \ell \leq k-1}$  such that  $p_{\ell} \in \{1, \ldots, n_{\ell}\}$ ,  $q_{\ell} \in \{1, \ldots, n_{\ell+1}\}$ , there exist k permutations  $(\pi_{\ell})_{1 \leq \ell \leq k}$  such that

$$J(A_1 \odot_{p_1,q_1} A_2 \odot_{p_2,q_2} \cdots \odot_{p_{k-1},q_{k-1}} A_k) = \sum_{\ell=1}^k \pi_\ell \left[ A_1 \odot A_2 \odot \cdots \odot J(A_\ell) \odot \cdots \odot A_k \right],$$

where the dot products of the right-hand side are along some axes that are not specify for simplicity.

*Proof.* The proof is done by induction on k. The formula for k = 1 is straightforward. Assume that the formula is true at order k. As before, we do not specify indexes for tensor dot products as we are only interested in their existence. By Lemma 9, we have

$$J(A_{1} \odot \cdots \odot A_{k+1})$$

$$= J((A_{1} \odot \cdots \odot A_{k}) \odot A_{k+1})$$

$$= \pi(J(A_{1} \odot \cdots \odot A_{k}) \odot A_{k+1}) + A_{1} \odot \cdots \odot A_{k} \odot J(A_{k+1})$$

$$= \pi\left[\sum_{\ell=1}^{k} \pi_{\ell} [A_{1} \odot A_{2} \odot \cdots \odot J(A_{\ell}) \odot \cdots \odot A_{k}] \odot A_{k+1}\right] + A_{1} \odot \cdots \odot A_{k} \odot J(A_{k+1})$$

$$= \pi\left[\sum_{\ell=1}^{k} \tilde{\pi}_{\ell} [A_{1} \odot A_{2} \odot \cdots \odot J(A_{\ell}) \odot \cdots \odot A_{k} \odot A_{k+1}]\right] + A_{1} \odot \cdots \odot A_{k} \odot J(A_{k+1})$$

$$= \sum_{\ell=1}^{k} \hat{\pi}_{\ell} [A_{1} \odot A_{2} \odot \cdots \odot J(A_{\ell}) \odot \cdots \odot A_{k} \odot A_{k+1}] + A_{1} \odot \cdots \odot A_{k} \odot J(A_{k+1})$$
(where  $\hat{\pi} = \pi \circ \tilde{\pi}$ )
$$= \sum_{\ell=1}^{k+1} \hat{\pi}_{\ell} [A_{1} \odot A_{2} \odot \cdots \odot J(A_{\ell}) \odot \cdots \odot A_{k} \odot A_{k+1}].$$

# **D** Experimental details

All the code to reproduce the experiments is available on GitHub at https://github.com/ afermanian/rnn-kernel. Our experiments are based on the PyTorch (Paszke et al., 2019) framework. When not specified, the default parameters of PyTorch are used.

**Convergence of the Taylor expansion.** For Figure 1,  $10^3$  random RNN with 2 hidden units are generated, with the default weight initialization. The activation is either the logistic or the hyperbolic tangent. In Figure 1b, only the results with the logistic activation are plotted. The process X is taken as a 2-dimensional spiral. The reference solution to the ODE (3) is computed with a numerical integration method from SciPy (Virtanen et al., 2020, scipy.integrate.solve\_ivp with the 'LSODA' method). The signature in the step-N Taylor expansion is computed with the package Signatory (Kidger and Lyons, 2021).

The step-N Taylor expansion requires computing higher-order derivatives of tensor fields (up to order N). This is a highly non-trivial task since standard deep learning frameworks are optimized for first-order differentiation only. We refer to, for example, Kelly et al. (2020), for a discussion on higher-order differentiation in the context of a deep learning framework. To compute it efficiently, we manually implement forward-mode higher-order automatic differentiation for the operations needed in our context (described in Appendix C). A more efficient and general approach is left for future work. Our code is optimized for GPU.

**Penalization on a toy example.** For Figure 2, the RNN is taken with 32 hidden units and hyperbolic tangent activation. The data are 50 examples of spirals, sampled at 100 points and labeled  $\pm 1$ 

according to their rotation direction. We do not use batching and the loss is taken as the cross entropy. It is trained for 200 epochs with Adam (Kingma and Ba, 2015) with an initial learning rate of 0.1. The learning rate is divided by 2 every 40 epochs. For the penalized RNN, the RKHS norm is truncated at N = 3 and the regularization parameter is selected at  $\lambda = 0.1$ . Earlier experiments show that this order of magnitude is sensible. We do not perform hyperparameter optimization since our goal is not to achieve high performance. The initial hidden state  $h_0$  is learned (for simplicity of presentation, our theoretical results were written with  $h_0 = 0$  but they extend to this case). The accuracy is computed on a test set of size 1000. We generate adversarial examples using 50 steps of projected gradient descent (following Bietti et al., 2019). The whole methodology (data generation + training) is repeated 20 times. The average training time on a Tesla V100 GPU for the RNN is 8.5 seconds and for the penalized RNN 12 seconds.

Figure 3 is obtained by selecting randomly one run among the 20 of Figure 2.

**Libraries.** We use PyTorch (Paszke et al., 2019) as our overall framework, Signatory (Kidger and Lyons, 2021) to compute the signatures, and SciPy (Virtanen et al., 2020) for ODE integration. We use Sacred (Klaus Greff et al., 2017) for experiment management. The links and licences for the assets are given in the following table:

Name	Homepage link	License
PyTorch	GitHub repository	BSD-style License
Sacred	GitHub repository	MIT License
SciPy	GitHub repository	BSD 3-Clause "New" or "Revised" License
Signatory	GitHub repository	Apache License 2.0