

Introduction

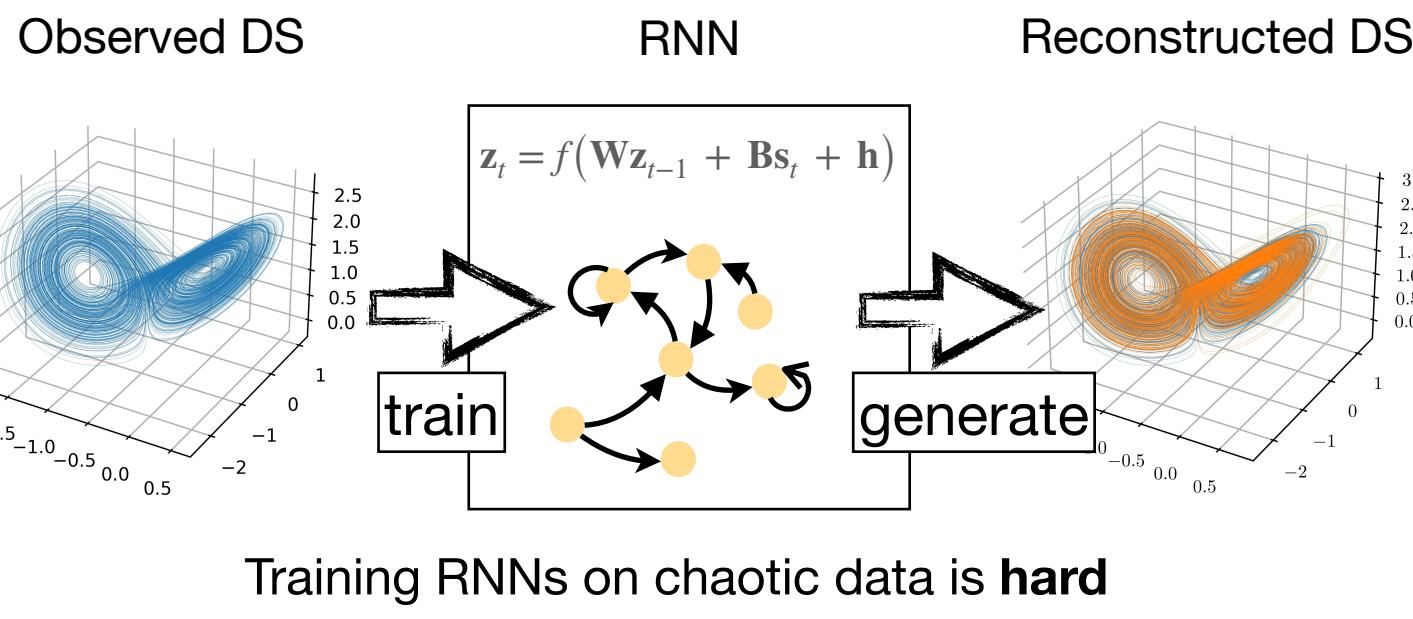
RNNs are a promising and increasingly popular tool for the analysis of **neural recordings** [2,3]

To trust in and interpret results obtained by training RNNs on neural data, we need to be certain they recover the underlying dynamics
-> **dynamical systems reconstruction**

Training RNNs is notoriously hard, especially if the data are governed by long and complicated temporal dependencies
-> **exploding and vanishing gradient problem** [1,4]

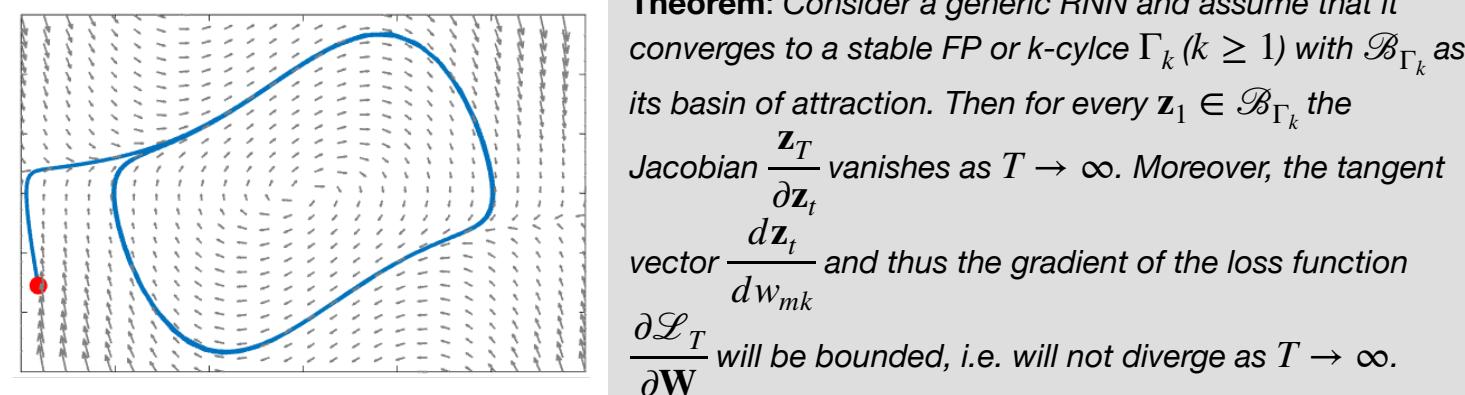
Problem setting

Reconstruction of **chaotic** dynamical systems (DS)



chaotic dynamics → **exploding gradients during training**

Equilibrium or cyclic behaviour
lead to vanishing gradients



How to train RNNs on chaotic neural data

Jonas Mikhaeil*, Zahra Monfared*, Daniel Durstewitz

Department of Theoretical Neuroscience, Central Institute of Mental Health, Medical Faculty Mannheim, Heidelberg University



Zentralinstitut
für
Seelische Gesundheit



UNIVERSITÄT
HEIDELBERG
ZUKUNFT
SEIT 1386

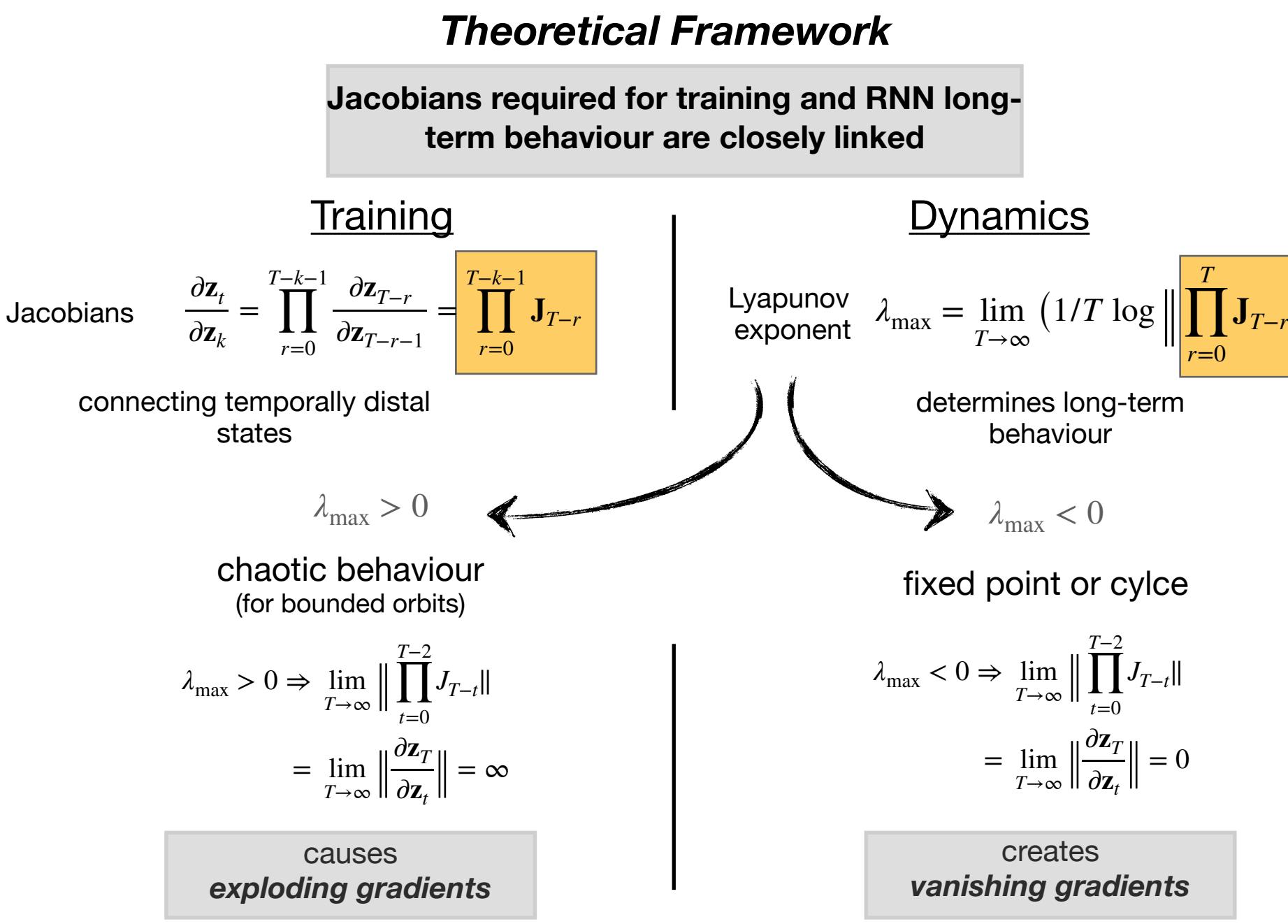


STRUCTURES
CLUSTER OF
EXCELLENCE

bccn
heidelberg-mannheim

Chaotic dynamics lead to exploding gradients

Theorem: Suppose that a generic RNN has a chaotic attractor Γ^* with \mathcal{B}_{Γ^*} as its basin of attraction. Then, for every orbit with $\mathbf{z}_1 \in \mathcal{B}_{\Gamma^*}$ the Jacobians $\frac{\partial \mathbf{z}_T}{\partial \mathbf{z}_t}$ connecting temporally distal states \mathbf{z}_T and \mathbf{z}_t ($T \gg t$), the tangent vectors $\frac{d\mathbf{z}_T}{dw_{mk}}$ and thus the gradients of the loss function $\frac{\partial \mathcal{L}_T}{\partial \mathbf{W}}$ will explode for $T \rightarrow \infty$.



Training with exploding gradients: Sparsely forced BPTT

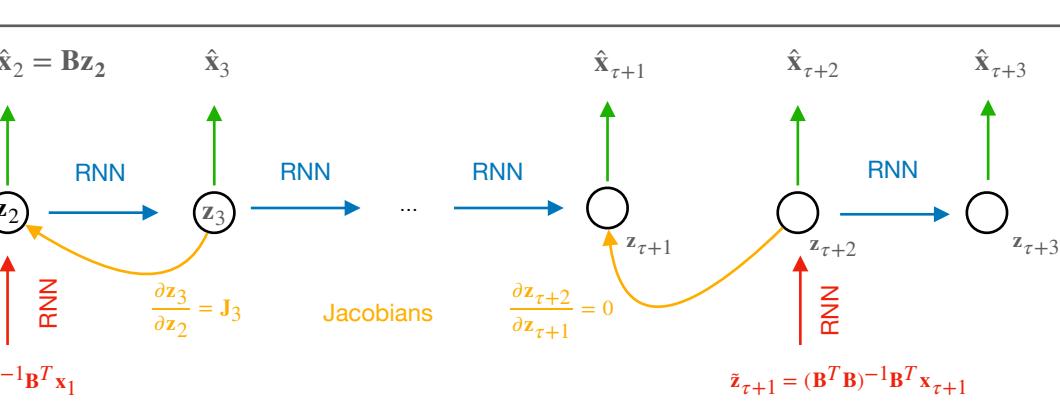
Force the RNN every τ time-steps onto observations

$$\mathbf{z}_{t+1} = \begin{cases} RNN(\hat{\mathbf{z}}_t) & \text{if } t \in \mathcal{T} = \{n\tau + 1\}_{n \in \mathbb{N}_0} \\ RNN(\mathbf{z}_t) & \text{else.} \end{cases}$$

This training procedure
1) breaks trajectory divergence
2) alleviates the exploding gradient problem

The forcing signal is obtained by "inverting" the output layer

$$\hat{\mathbf{z}}_t = (\mathbf{B}^T \mathbf{B})^{-1} \mathbf{B}^T \mathbf{x}_t$$

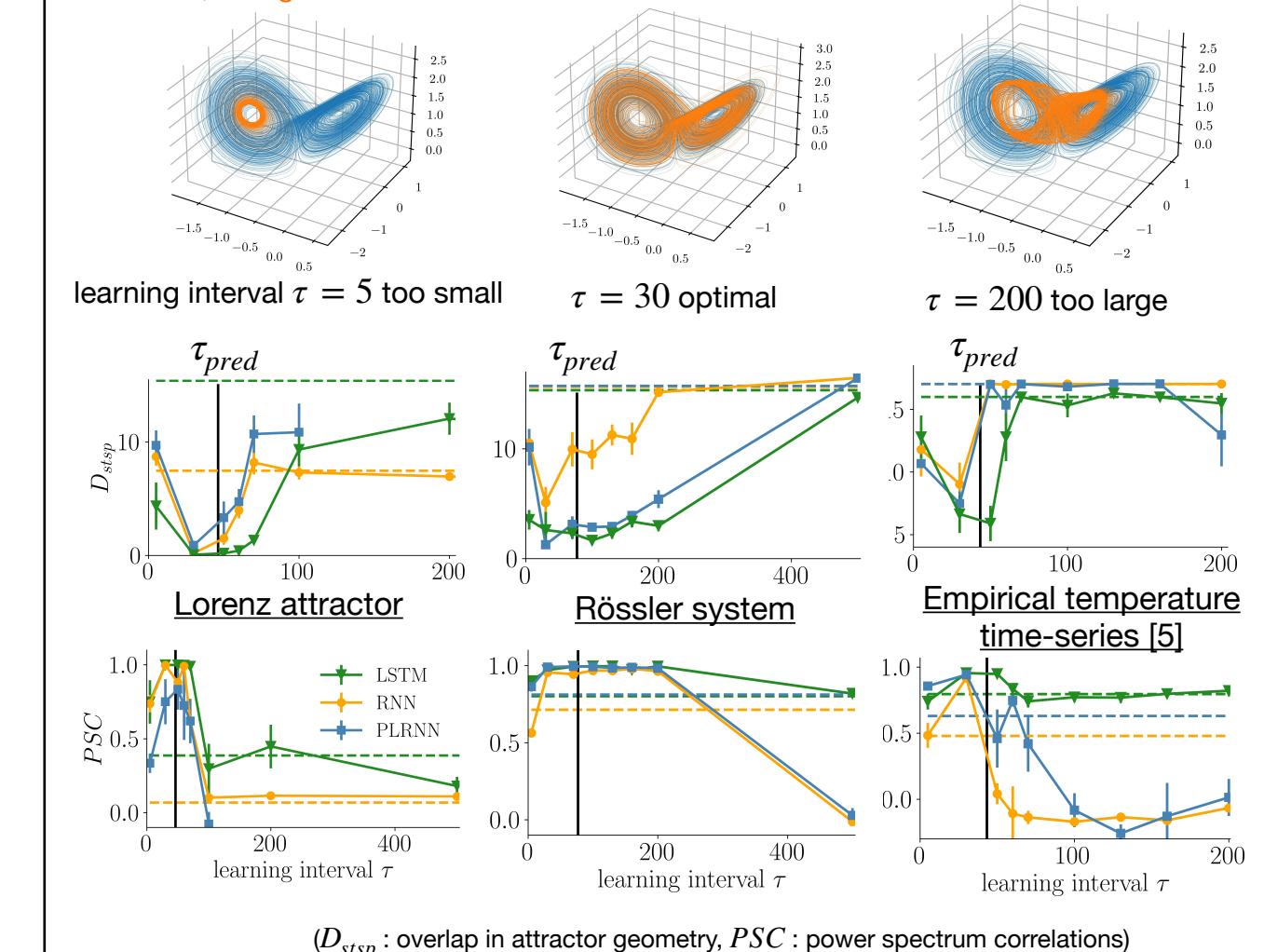


The learning interval τ can be chosen in accordance with the Lyapunov spectrum

$$\tau_{pred} = \log(2)/\lambda_{\max}$$

Results:

Reconstruction examples: blue: data, orange: LSTM reconstruction



Sparsely forced BPTT allows the reconstruction of chaotic dynamical systems (continuous lines)

and significantly outperforms the classical approach of gradient clipping (dashed lines)

References

- [1] Y. Bengio, P. Simard, and P. Frasconi, 'Learning long-term dependencies with gradient descent is difficult', IEEE Transactions on Neural Networks, vol. 5, no. 2, pp. 157–166, Mar 1994, 10.1109/72.279181.
- [2] C. Pandarinath et al., 'Inferring single-trial neural population dynamics using sequential auto-encoders', Nat Methods, vol. 15, no. 10, pp. 805–815, Oct 2018, 10.1038/s41592-018-0109-9.
- [3] M. G. Perich et al., 'Inferring brain-wide interactions using data-constrained recurrent neural network models', Neuroscience, preprint, Dec 2020, 10.1101/2020.12.18.423348.
- [4] D. Schmidt et al., 'Identifying nonlinear dynamical systems with multiple time scales and long-range dependencies', ICLR, 2021, <https://openreview.net/forum?id=XYzwxPlQu6>
- [5] recorded at the Weather Station at the Max Planck Institute for Biogeochemistry in Jena, Germany spanning the time period between 2009 and 2016. Data can be accessed at <https://www.kaggle.com/pankrziyski/weather-archive-jena>