



STRUCTURES  
CLUSTER OF  
EXCELLENCE



# Bifurcations and loss jumps in RNN training

Lukas Eisenmann<sup>1,2,\*</sup>, Zahra Monfared<sup>1,\*</sup>, Niclas Göring<sup>1,2</sup>, and Daniel Durstewitz<sup>1,2,3</sup>

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

2 Faculty of Physics and Astronomy, Heidelberg University, Heidelberg, Germany

3 Interdisciplinary Center for Scientific Computing, Heidelberg University

\* These authors contributed equally



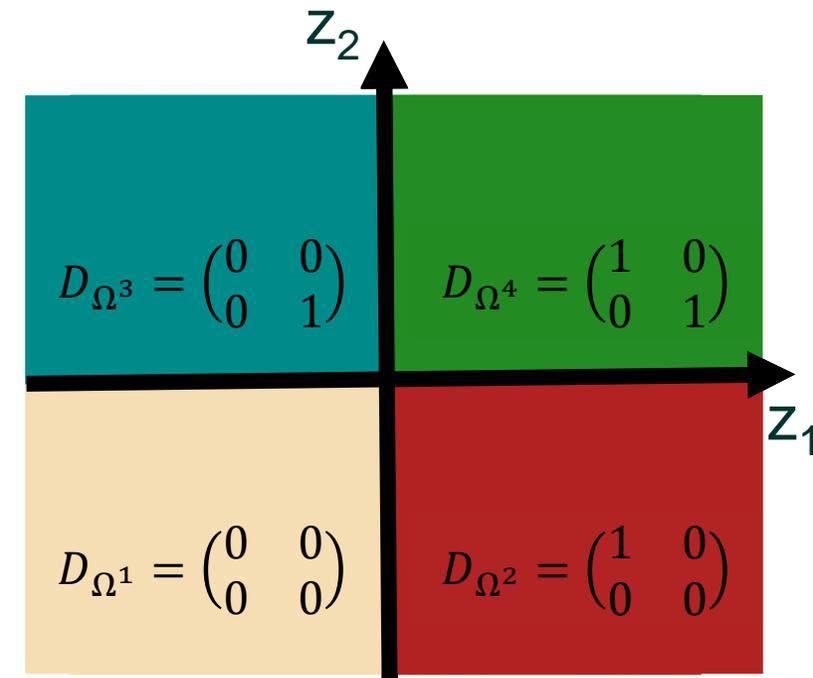
# PLRNN

$$z_t = F(z_{t-1}) = A z_{t-1} + W \text{ReLU}(z_{t-1}) + C s_t + h$$

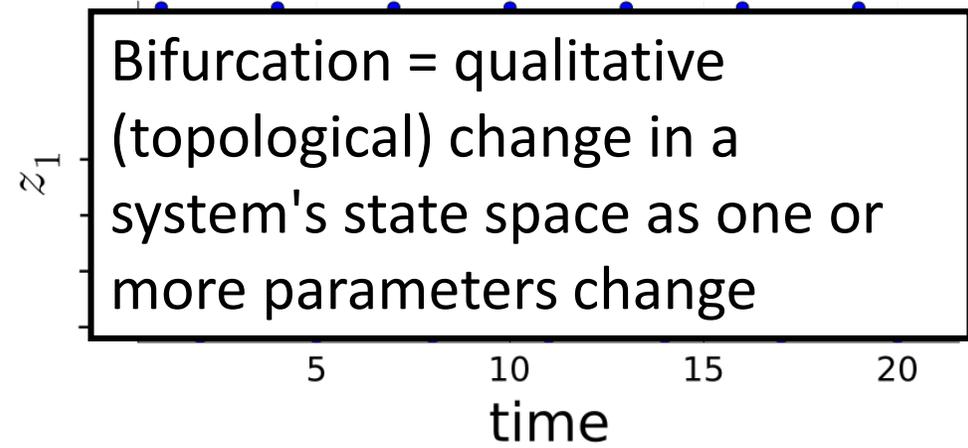
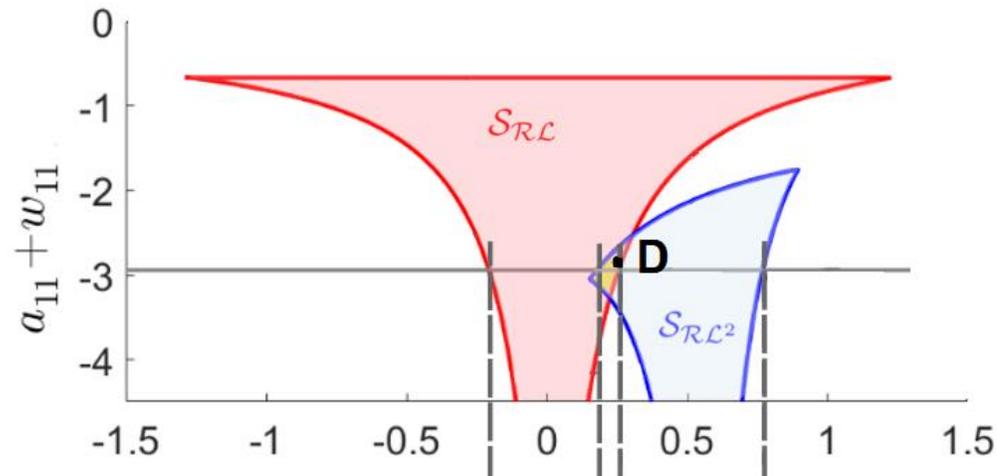
$$A \in \mathbb{R}^{M \times M} \quad h \in \mathbb{R}^M \quad W \in \mathbb{R}^{M \times M} \quad C \in \mathbb{R}^{M \times K} \quad s_t \in \mathbb{R}^K \quad z_t \in \mathbb{R}^M$$

$$D_{\Omega(t)} := \text{diag}(d_{\Omega,t}) \quad \text{with} \quad d_{m,t} = \begin{cases} 1, & z_{m,t} > 0 \\ 0, & \text{else} \end{cases}$$

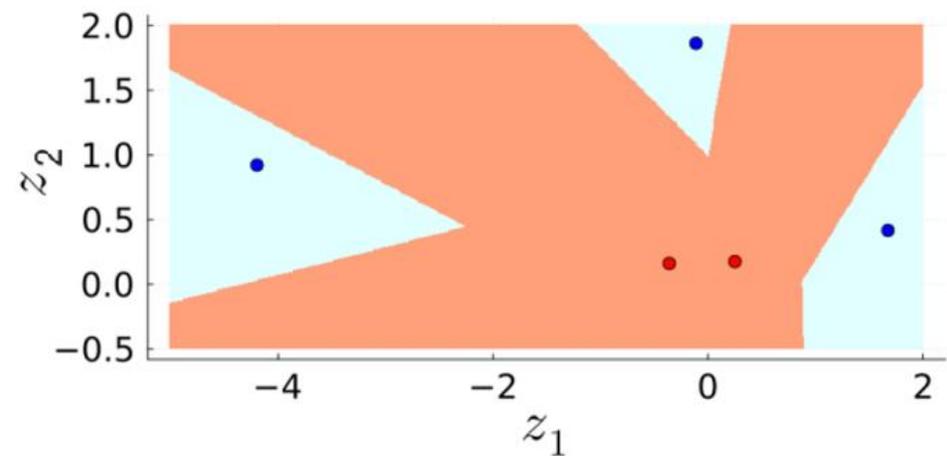
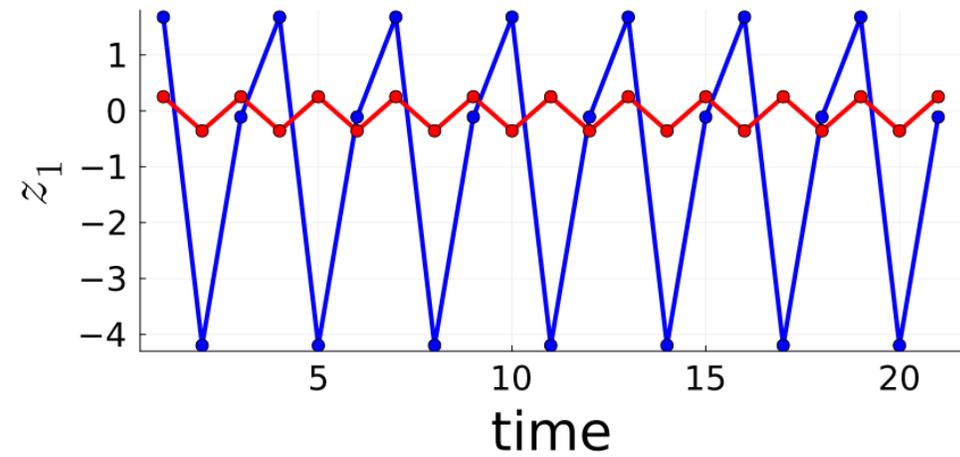
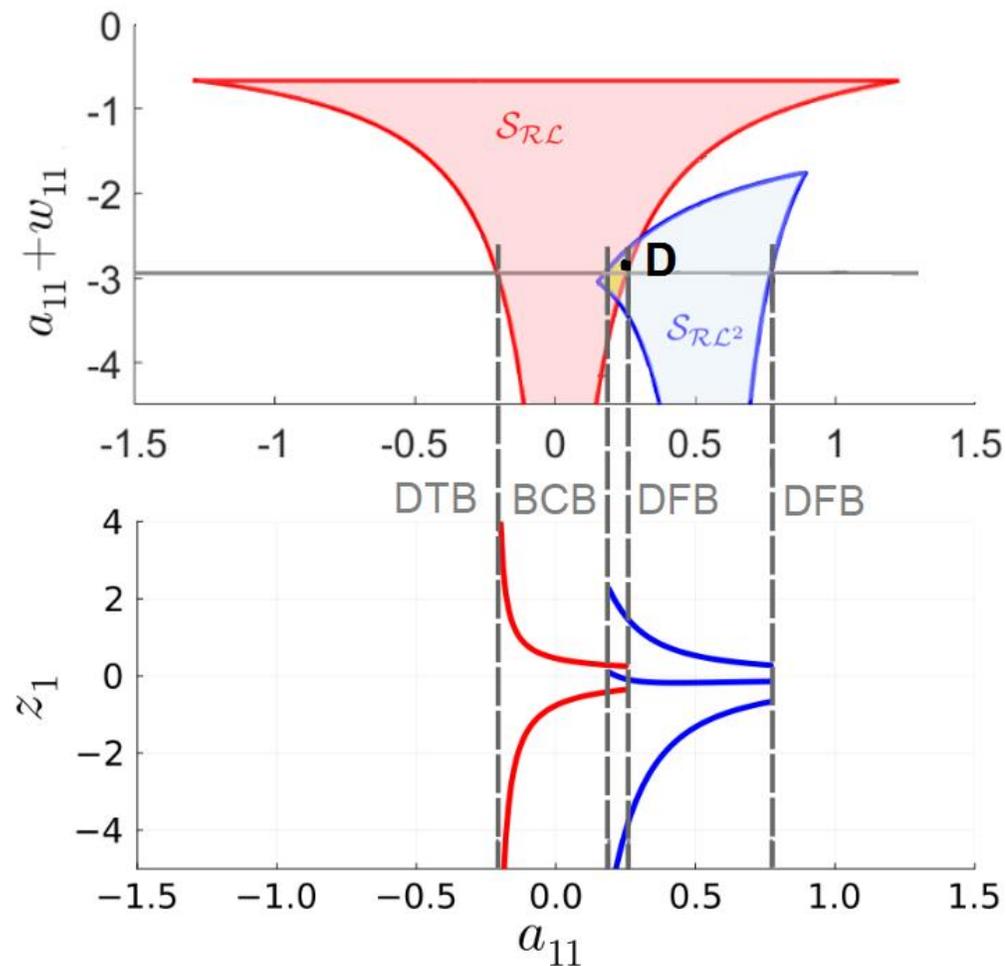
$$\begin{aligned} z_t = F(z_{t-1}) &= (A + W D_{\Omega(t-1)}) z_{t-1} + C s_t + h \\ &=: W_{\Omega(t-1)} z_{t-1} + C s_t + h \end{aligned}$$



# Bifurcation manifolds in PLRNN parameter space

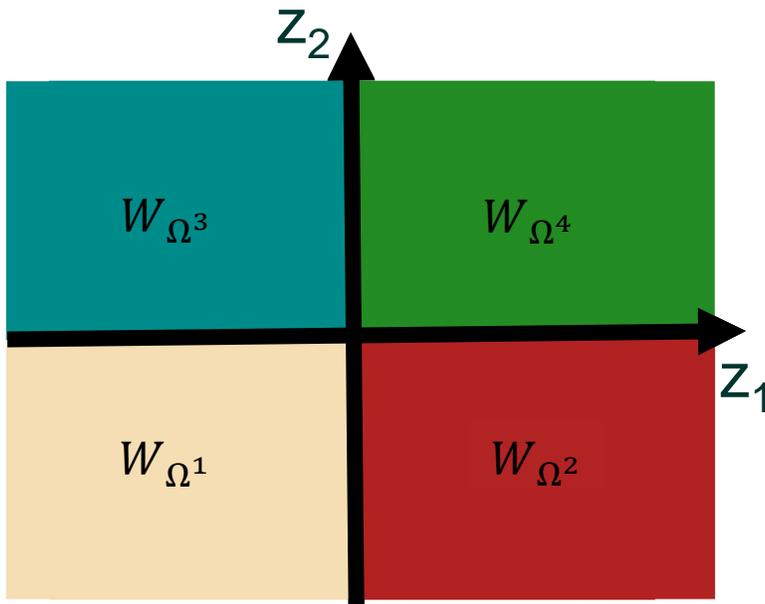


# Bifurcation manifolds in PLRNN parameter space

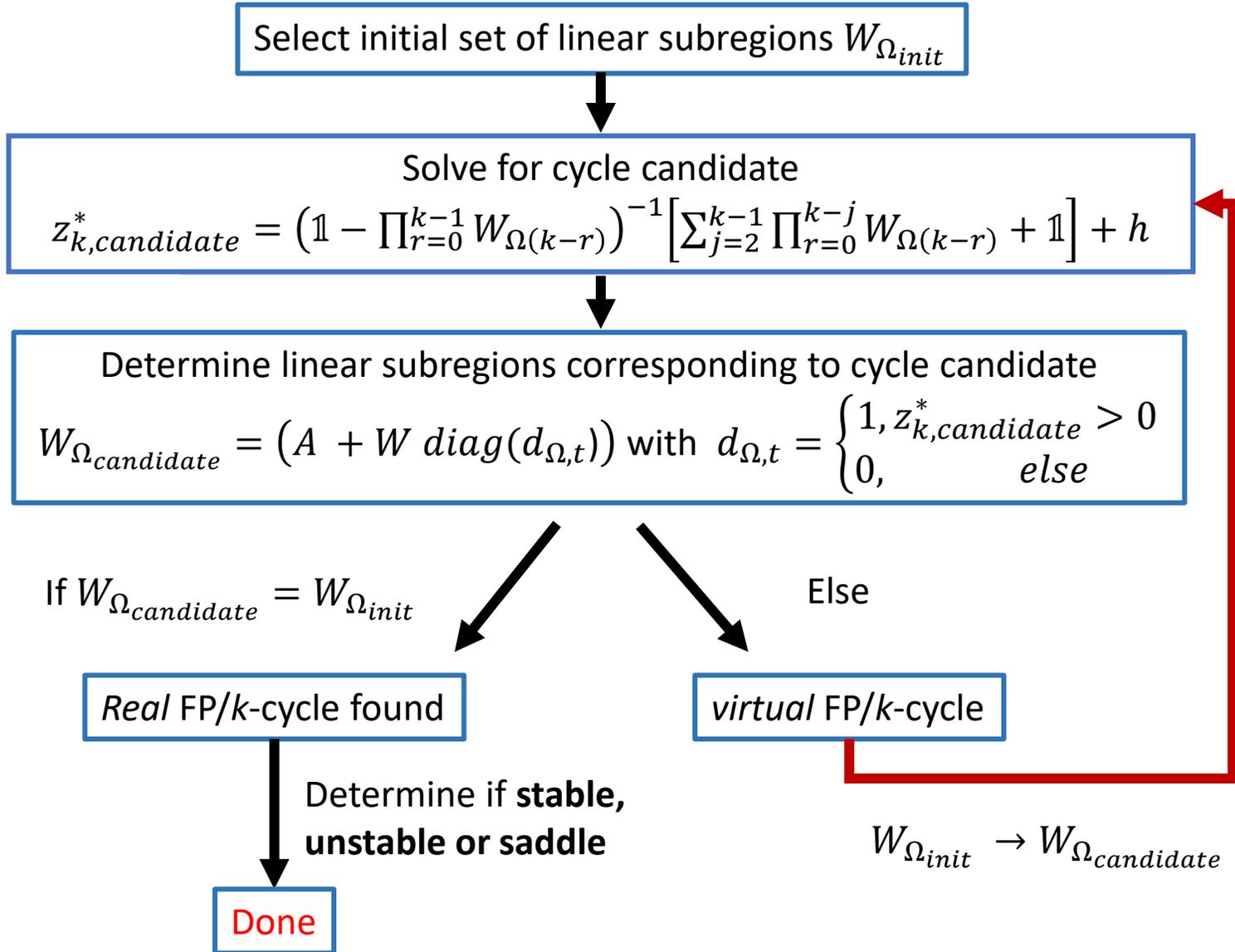


# Searcher for Cycles and Fixed points: SCYFI

- Mathematically tractable: Allows for semi-analytic calculation of fixed points and cycles
- BUT: Combinatorial problem!

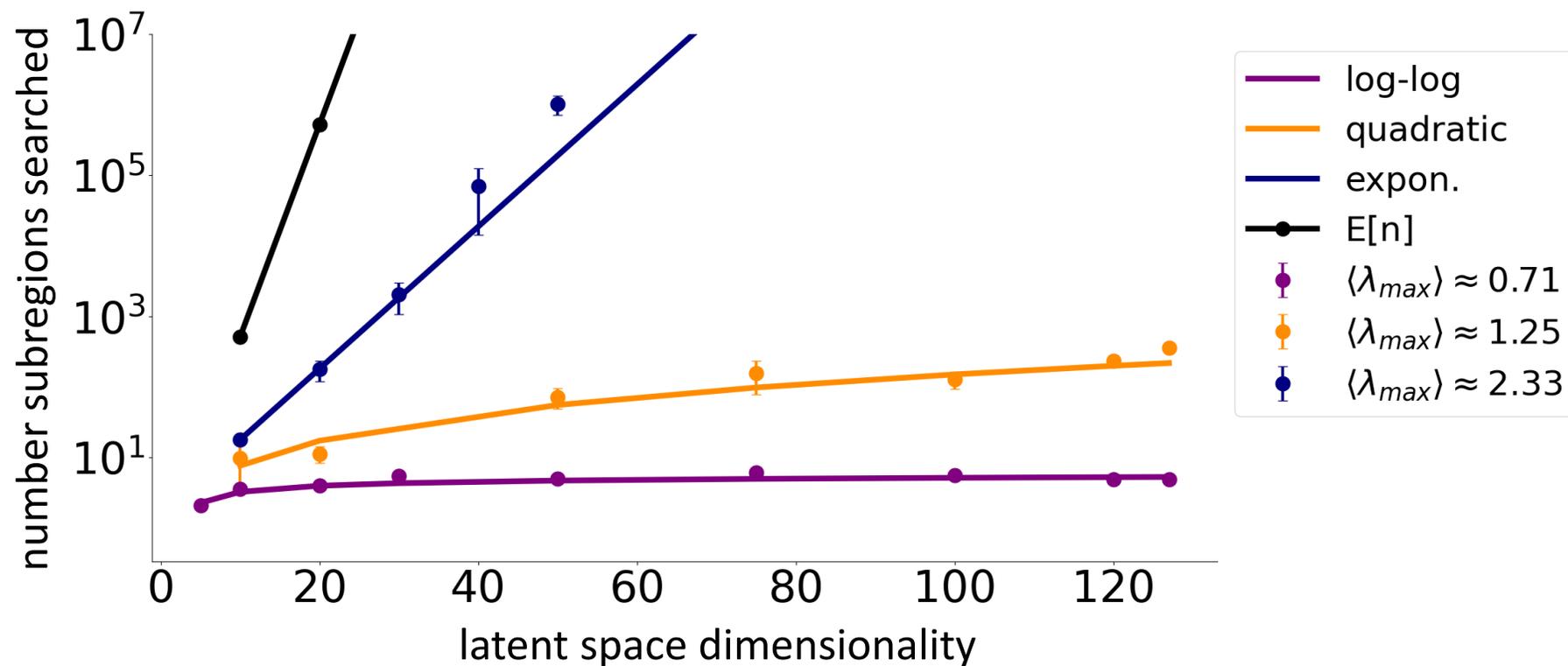


→ Number of linear regions:  $2^{Mk}$



# SCYFI: Scaling

**Theorem 3.** Under the condition  $\|A\| + \|W\| < 1$ , SCYFI will converge in at most linear time

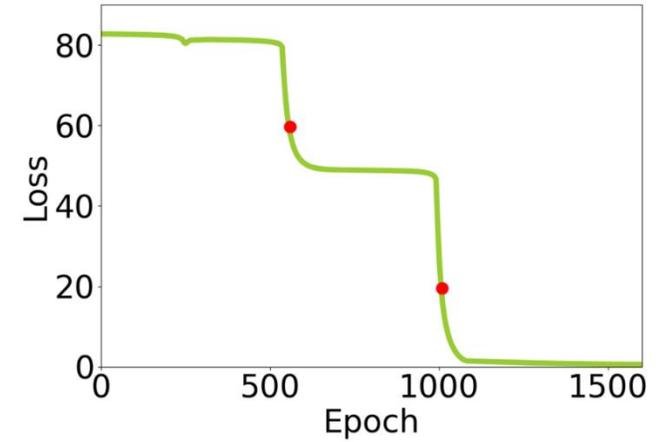
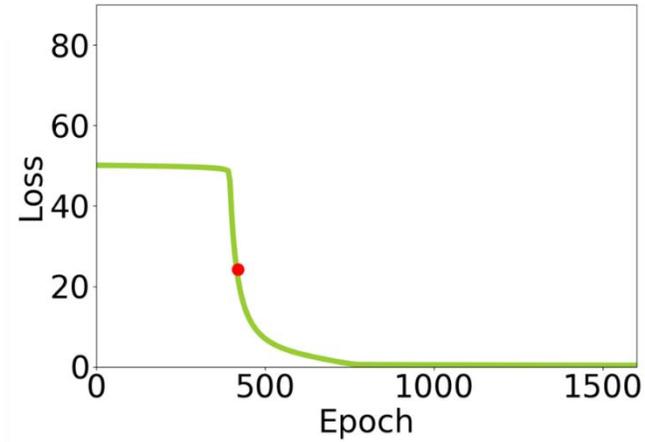
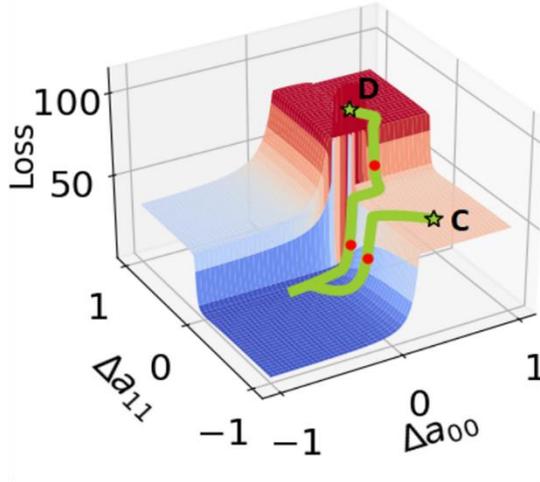
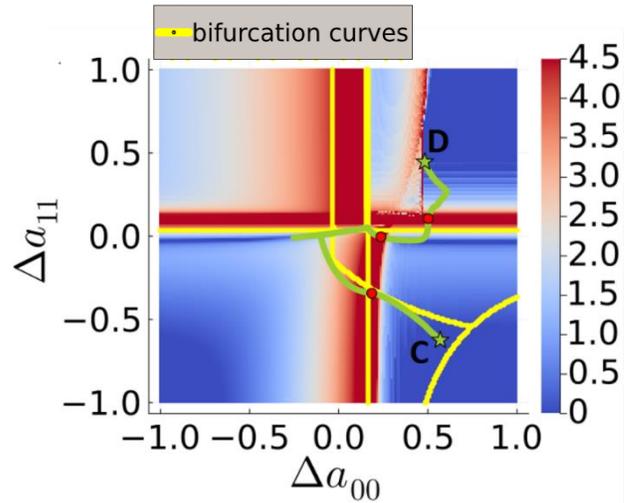


# Theorems: Bifurcations and exploding/ vanishing gradients

**Theorem 1.** If a stable fixed point or a  $k$ -cycle undergoes a degenerate transcritical bifurcation, the norm of the PLRNN loss gradient tends to infinity  $\lim_{t \rightarrow \infty} \left\| \frac{dL}{d\theta} \right\| = \infty$

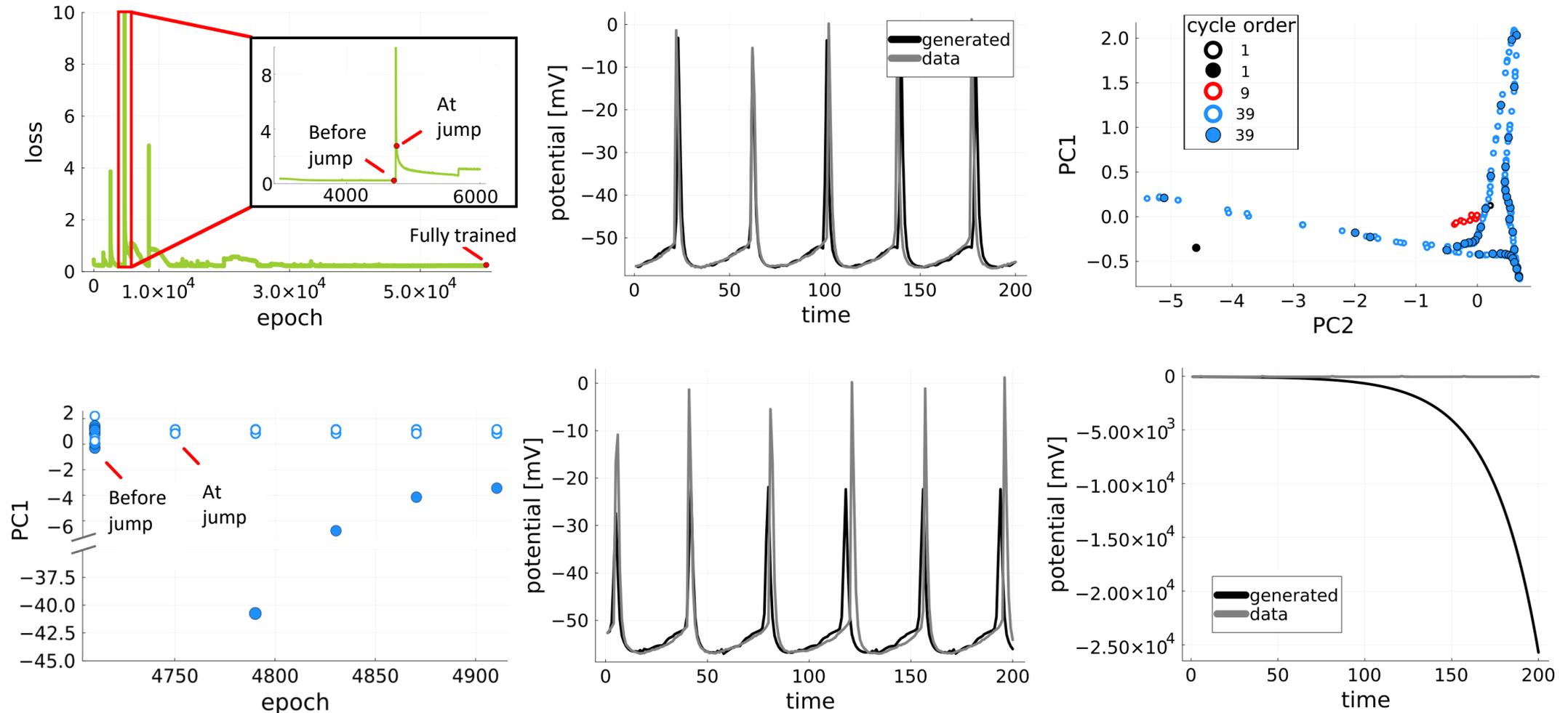
**Theorem 2.** If a stable fixed point or a  $k$ -cycle undergoes a border collision bifurcation, the norm of the PLRNN loss gradients vanishes  $\lim_{t \rightarrow \infty} \left\| \frac{dL}{d\theta} \right\| = 0$

# Toy example

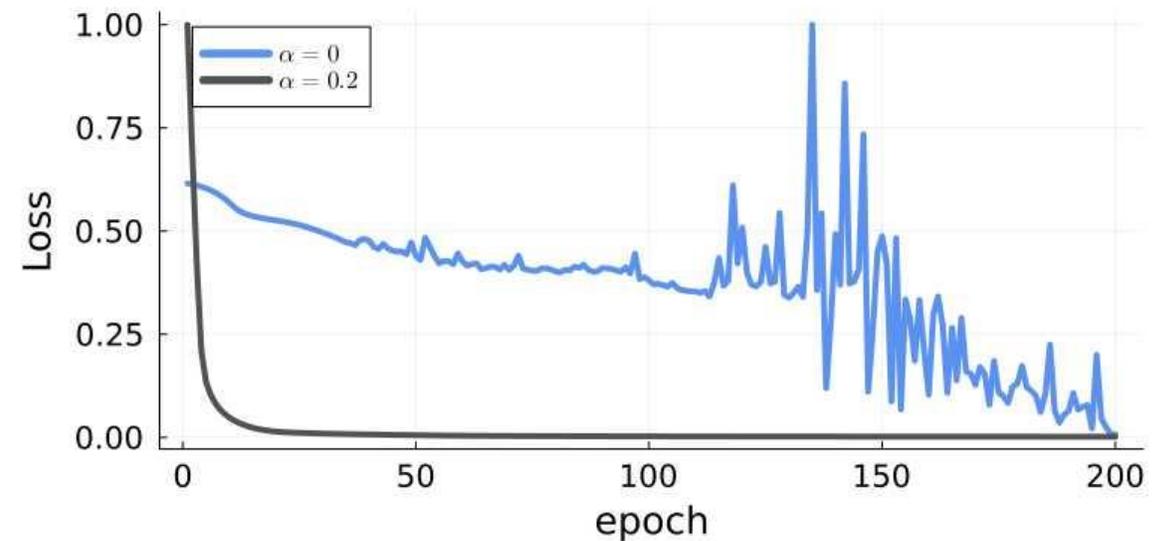
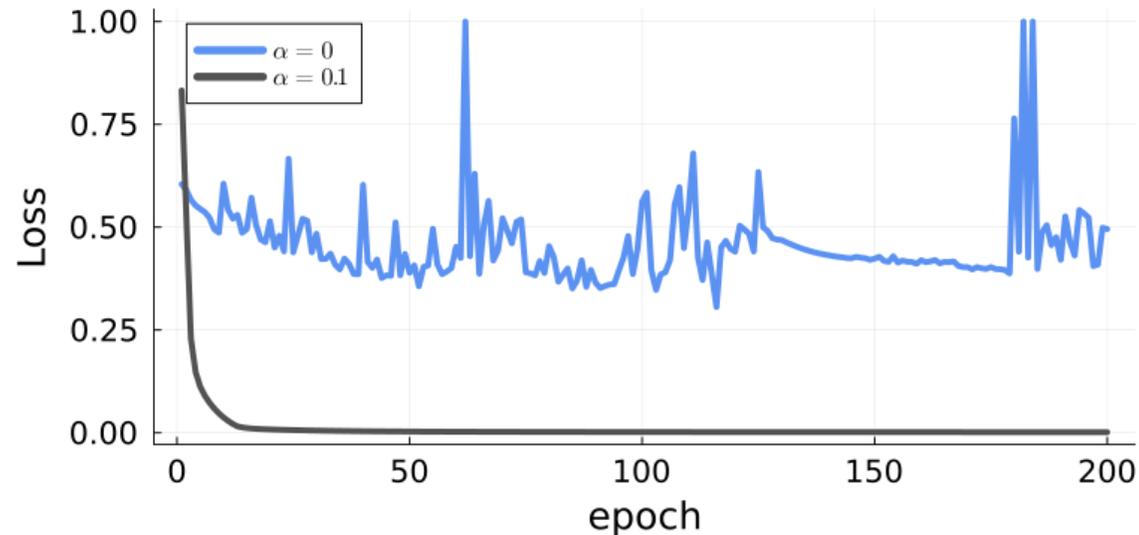


→ Bifurcations can cause jumps in the loss

# Empirical example: Training PLRNN on membrane voltage traces of real cell



# Generalized Teacher Forcing (GTF)<sup>1</sup> prevents bifurcations in training



## Theorem 3.

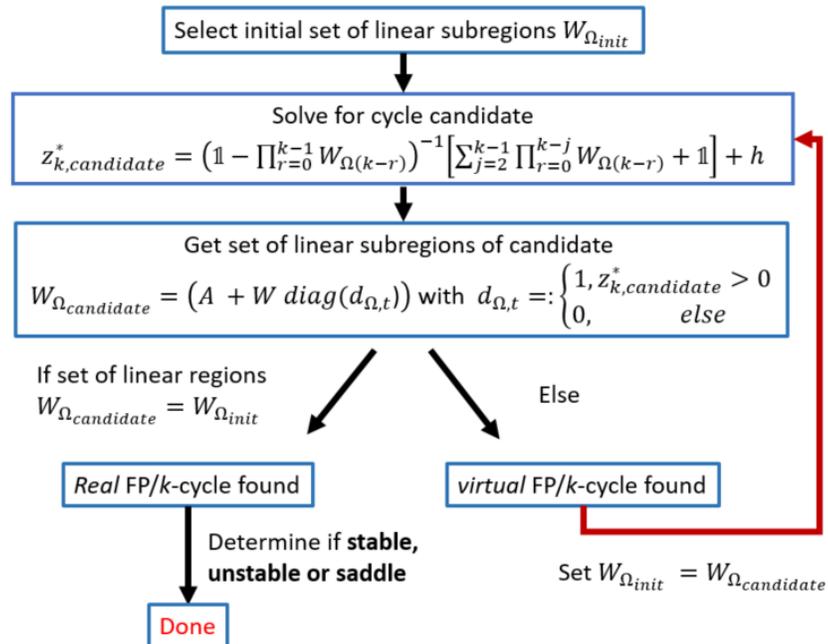
If  $\|A\| + \|W\| < 1$  then for any  $0 < \alpha < 1$  GTF controls the system, preventing degenerate transcritical bifurcations.

If  $\|A\| + \|W\| = r > 1$  then for any  $1 - \frac{1}{r} < \alpha < 1$  GTF prevents degenerate transcritical bifurcations.

# Conclusion

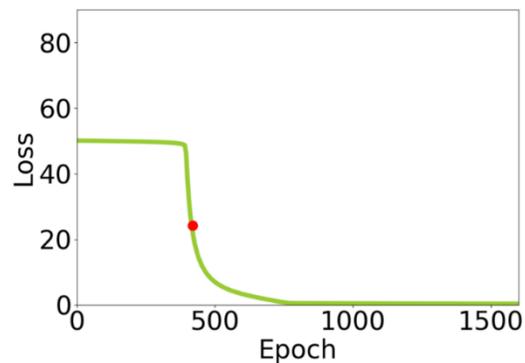
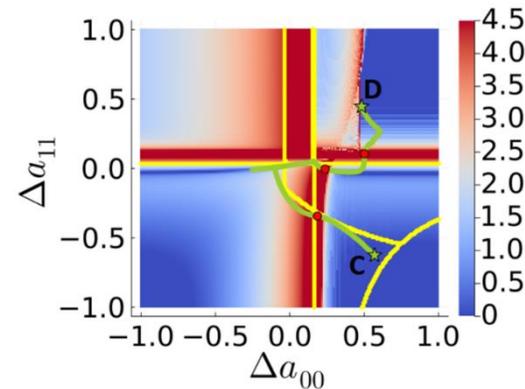
## SCYFI

- Computes fixed points *and* k-cycles exactly
- Efficient: surprisingly good, often linear, scaling



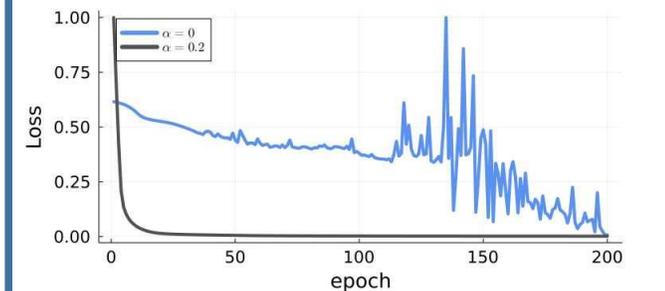
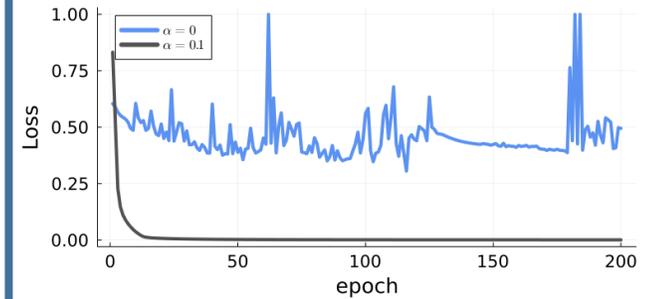
## Theorems

Formal connection between bifurcations and exploding or vanishing gradients



## Implications for RNN training

Generalized teacher forcing provably avoids bifurcations in training





**Thanks for  
your  
attention!**



This work was supported by the German Research Foundation (DFG) through individual grants Du 354/10-1 & Du 354/15-1 to DD, within research cluster FOR-5159 (“Resolving prefrontal flexibility”; Du 354/14-1), and through the Excellence Strategy EXC 2181/1 – 390900948 (STRUCTURES)