



TIME SERIES ANALYSIS & RECURRENT NEURAL NETWORKS

#10

- Training RNN (cont.): BPTT
- Vanishing/ exploding grad. problem
- LSTM, GRU, rPLRNN

Main lecture: Daniel Durstewitz

Exercises: Leonard Bereska, Manuel Brenner,
Daniel Kramer, Georgia Koppe

Heidelberg University

Training RNN by grad. descent

$$x_t = \phi(w_{x_{t-1}} + h + c_t)$$

$$D = \{ \{ s_t^{(p)} \}, \{ \tilde{x}_t^{(p)} \}, p: 1 \dots P \}$$

$$L(w, h, c) = \frac{1}{2} \sum_{p=1}^P \sum_{t=1}^T \sum_{i=1}^N I\{\varepsilon_{it}=1\} (\tilde{x}_{it} - x_{it})^2 \frac{1}{PT}$$

$$\nabla_w L := \left(\frac{\partial L}{\partial w_{ij}} \right), \quad r_{kij} := \frac{\partial x_k}{\partial w_{ij}}, O(M^3)$$

Back-Prop. through time (BPT)

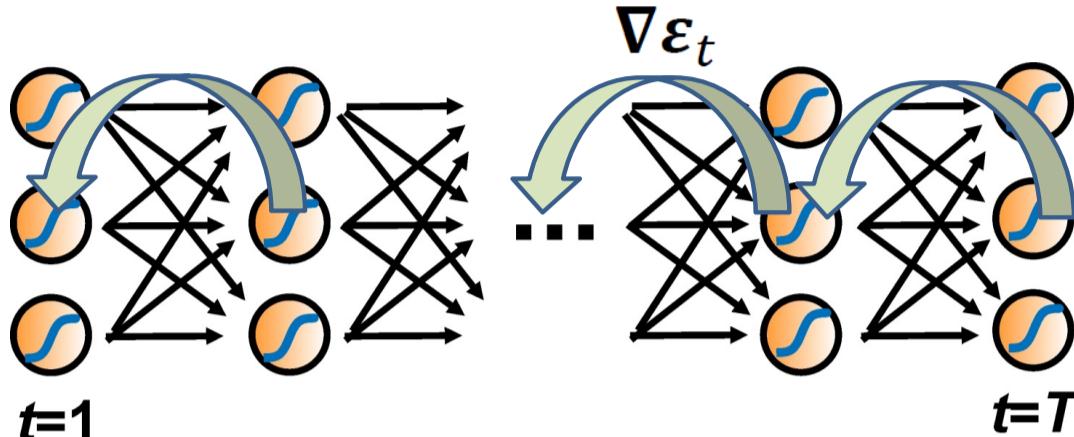
Werbos (88), Willia

Making RNN deep in time: “Vanishing/ exploding gradient problem”

“Inputs”: $S = \{s_t\}$ “Outputs”: $\tilde{\mathbf{z}} = \{\tilde{\mathbf{z}}_t\}, t = 1 \dots T$

RNN $\mathbf{z}_t \in \mathbb{R}^M$
 $\mathbf{z}_t = \phi(\mathbf{W}\mathbf{z}_{t-1} + \mathbf{h} + \mathbf{C}s_t)$

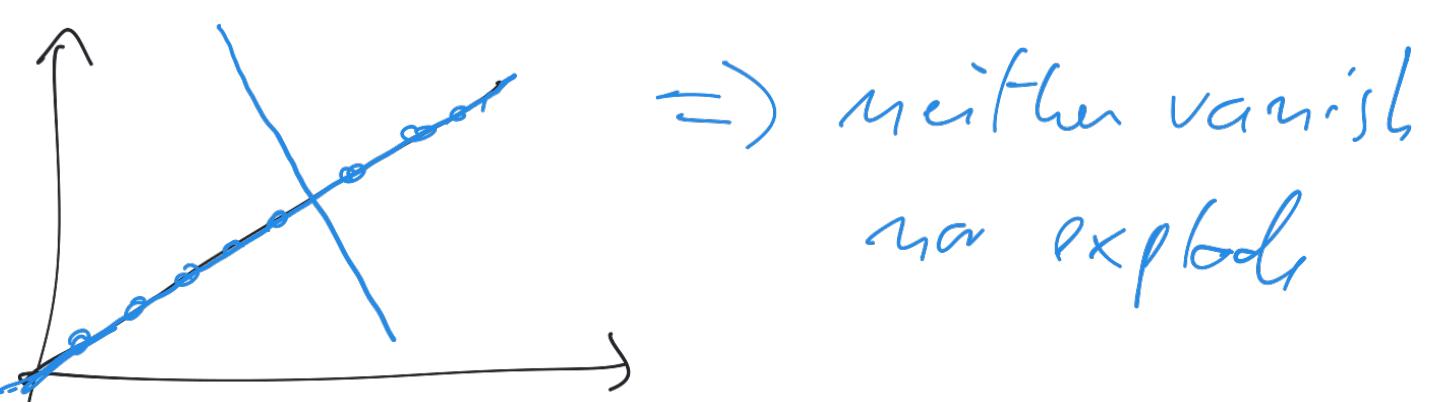
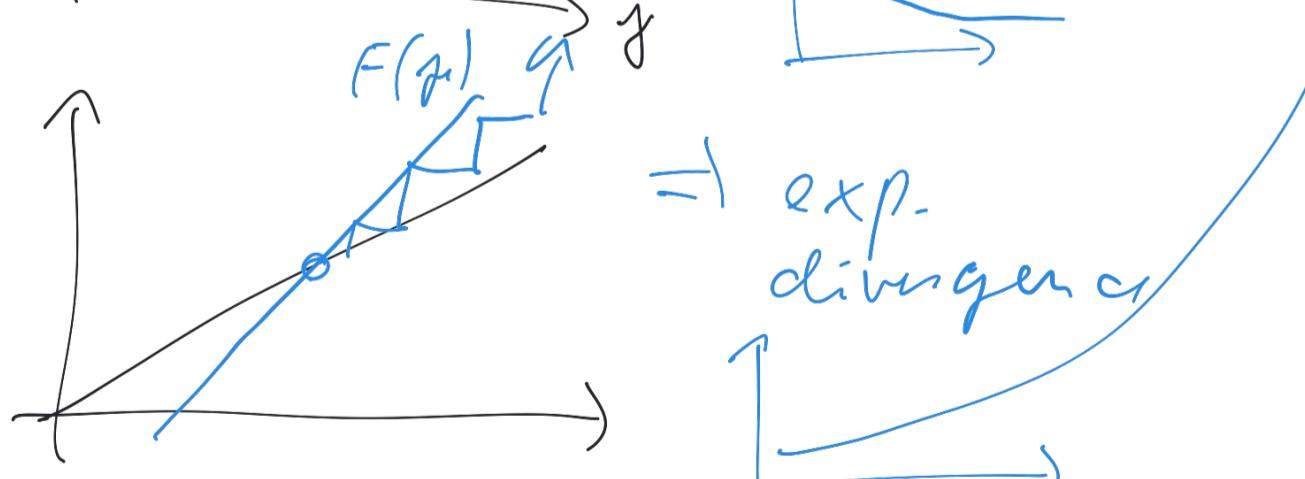
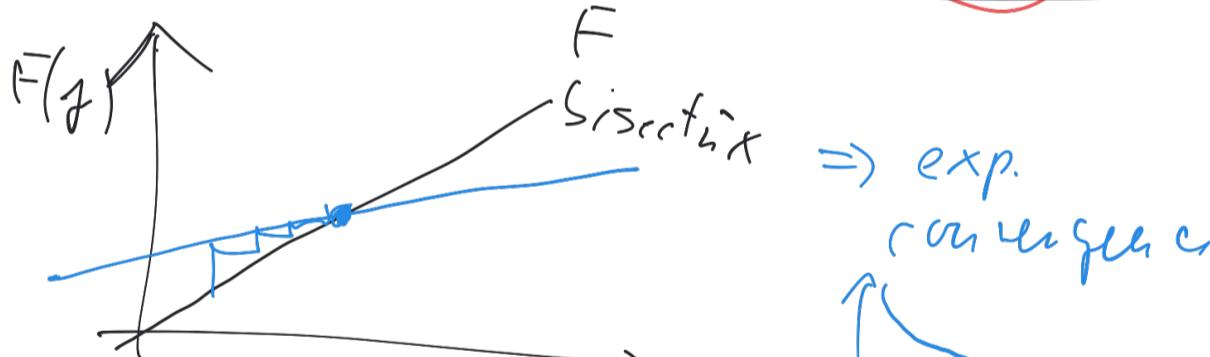
Loss function
 $\varepsilon = \sum_{r=1}^R \sum_{t=1}^T \|\tilde{\mathbf{z}}_t - \mathbf{z}_t\|_2^2$



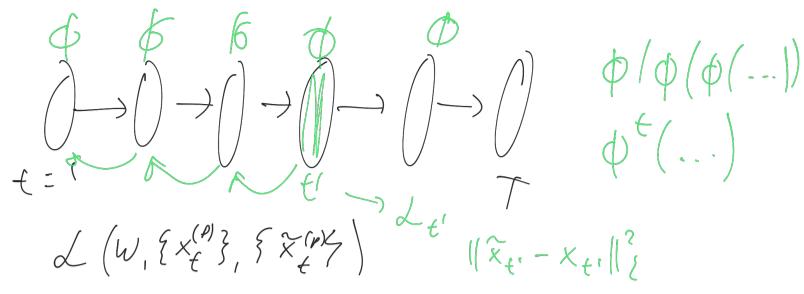
RTKL

$$\mathcal{L} = \frac{1}{2} \sum_i \sum_j \sum_k \mathbb{I}\{\varepsilon_{it}=1\} (\hat{x}_{it} - x_{it})^2$$

$$\gamma_{kij}^{(t)} := \frac{\partial x_{kt}}{\partial w_{ij}} = \phi' \left[\delta_{ki} x_{j,t-1} + \sum_{l=1}^M w_{kl} \frac{\partial x_{lt-1}}{\partial w_{ij}} \right]$$



\Rightarrow neither vanish
nor explode



$$\phi / \phi(\phi(\dots))$$

$$\phi^t(\dots)$$

$$\mathcal{L}(w, \{x_t^{(n)}\}, \{\hat{x}_t^{(n)}\}) \parallel \|\hat{x}_{t+1} - x_{t+1}\|_2^2$$

$$\frac{\partial \mathcal{L}}{\partial w_{ij}} = \sum_{t=1}^T \frac{\partial \mathcal{L}_t}{\partial w_{ij}}$$

$$\frac{\partial \mathcal{L}_t}{\partial w_{ij}} = \sum_{\tau=1}^t \frac{\partial \mathcal{L}_t}{\partial x_\tau} \frac{\partial x_\tau}{\partial x_\tau} \frac{\partial x_\tau}{\partial w_{ij}}$$

$$\begin{matrix} M \times T \\ \left(\begin{array}{cccc} \frac{\partial x_{1c}}{\partial x_{1c}} & \frac{\partial x_{1c}}{\partial x_{2c}} & \dots & \frac{\partial x_{1c}}{\partial x_{Mc}} \\ \vdots & \ddots & \ddots & \vdots \\ \frac{\partial x_{Mc}}{\partial x_{1c}} & \frac{\partial x_{Mc}}{\partial x_{2c}} & \dots & \frac{\partial x_{Mc}}{\partial x_{Mc}} \end{array} \right) \end{matrix}$$

$$\frac{\partial x_\tau}{\partial x_\tau} = \frac{\partial x_\tau}{\partial x_{\tau-1}} \frac{\partial x_{\tau-1}}{\partial x_{\tau-2}} \dots \frac{\partial x_{2+1}}{\partial x_2} = \prod_{u=\tau+1}^T \frac{\partial x_u}{\partial x_{u-1}}$$

$$x_u = \phi(Wx_{u-1} + b) = \text{ReLU}(W \text{diag}[\frac{\partial \phi}{\partial x_{iu}}])$$

$\Rightarrow \max |\text{eig}(W\phi')| > 1$: divergence!
explode!

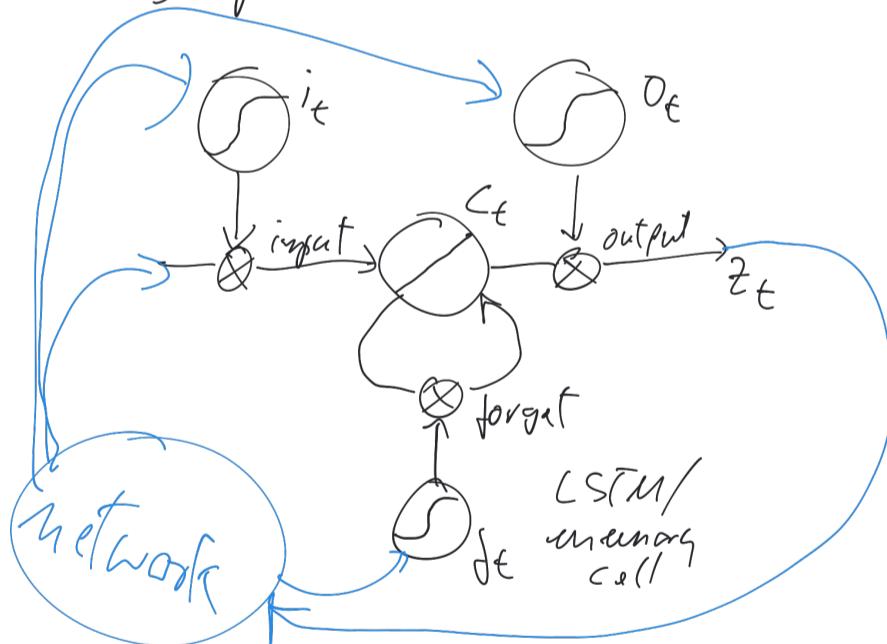
$\max |\text{eig}(W\phi')| < 1$: convergence!
vanish!

solutions

Long-Short Term Memory (LSTM)

Hochreiter & Schmidhuber (97) NC

→ Bengio



$$c_t = f_t \odot c_{t-1} + i_t \odot \tanh [W_c z_{t-1} + h_c]$$

element-wise
multiplic.

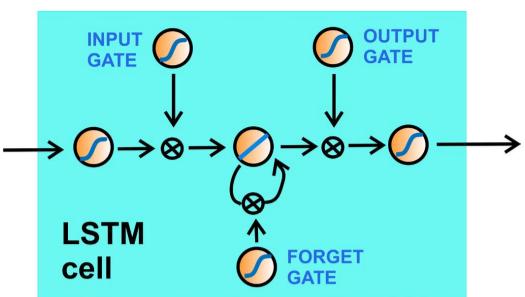
$$z_t = o_t \odot \tanh(c_t)$$

$$i_t = \sigma(W_i z_{t-1} + h_i)$$

$$o_t = \sigma(W_o z_{t-1} + h_o) \in [0, 1],$$

$$f_t = \sigma(W_f z_{t-1} + h_f) \quad \sigma(y) = \frac{1}{1 + e^{-y}}$$

Long Short-Term Memory



Hochreiter & Schmidhuber (1997)

$$\begin{aligned} i_t^l &= \sigma(W_i^l[X_t; h_{t-1}^l; h_t^{l-1}] + b_i^l) \\ f_t^l &= \sigma(W_f^l[X_t; h_{t-1}^l; h_t^{l-1}] + b_f^l) \\ s_t^l &= f_t^l s_{t-1}^l + i_t^l \tanh(W_s^l[X_t; h_{t-1}^l; h_t^{l-1}] + b_s^l) \\ o_t^l &= \sigma(W_o^l[X_t; h_{t-1}^l; h_t^{l-1}] + b_o^l) \\ h_t^l &= o_t^l \tanh(s_t^l) \end{aligned}$$

Gated Recurrent Unit (GRU)

(Cho... Bengio (2014))

$$\text{update gate: } u_t = \sigma(W_u z_{t-1} + h_u) \in [0, 1]$$

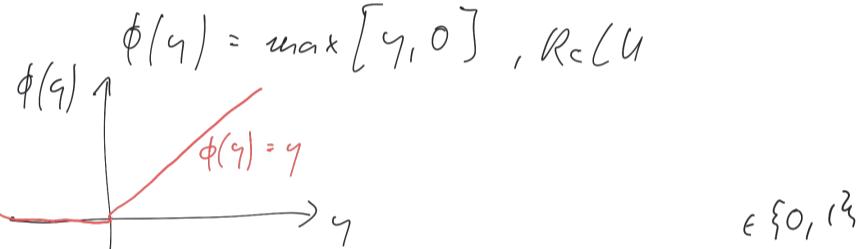
$$\text{reset gate: } r_t = \sigma(W_r z_{t-1} + h_r)$$

$$\begin{aligned} \text{hidden state: } z_t &= u_t z_{t-1} + (1 - u_t) \times \\ &\quad \tanh \left[W_z \underbrace{(r_t z_{t-1})}_{+ C_{z_t}} + h_z \right] \end{aligned}$$

Initialization / Regularization

(Le... Hinton (2015) av X. L)

$$z_t = \phi(W z_{t-1} + h + C_{z_t})$$



$$\frac{\partial z_t}{\partial z_{t-1}} = \prod_{u=t+1}^T \frac{\partial z_u}{\partial z_{u-1}} = \prod_{u=t+1}^T W \underset{\epsilon \{0, 1\}}{\text{diag}} [\phi']$$

$$W = I$$

$$z_t = A z_{t-1} + \underbrace{W \phi(z_{t-1}) + h + C_{z_t}}_{\max[0, z_{t-1}]}$$

$$A = \begin{pmatrix} 1 & & & & \\ 0 & -1 & & & \\ & & \ddots & & \\ & & & 0 & \\ 0 & & & & 0 \end{pmatrix}, \quad W = \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots \\ \text{---} & \text{---} & \text{---} & \text{---} \\ \text{---} & \text{---} & \text{---} & \text{---} \\ \text{---} & \text{---} & \text{---} & \text{---} \end{pmatrix}$$

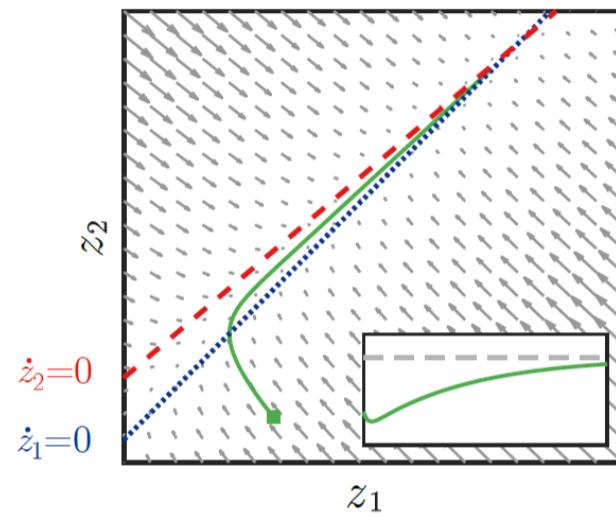
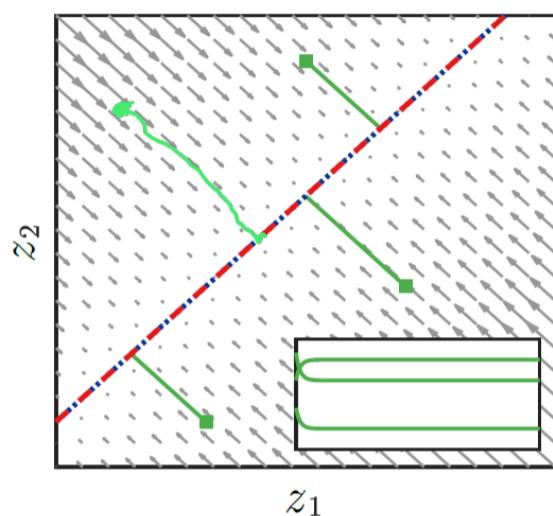
$$h = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

$$z_t \in \mathbb{R}^M, A \in \mathbb{R}^{M \times M} \text{ diag. unfx.}, W \in \mathbb{R}^{M \times M}$$

$$M_{\text{new}} < M$$

$$C = \mathcal{L}_{\text{MSE}} + \lambda \left[\sum_{i=1}^{M_{\text{new}}} (1 - a_{ii})^2 + \sum_{i=1}^{M_{\text{new}}} \sum_{j=1}^M w_{ij}^2 + \sum_{i=1}^{M_{\text{new}}} h_i^2 \right]$$

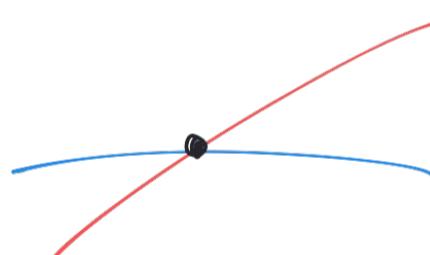
Schmidt et al. (2021) ICLR



nullclines

$$\dot{z}_1 = 0$$

$$\dot{z}_2 = 0$$



$$z_t = (1 - \varepsilon) z_{t-1}$$

ε very
small

Manifold-attractor regularization

$$\mathbf{z}_t = \mathbf{A}\mathbf{z}_{t-1} + \mathbf{W}\phi(\mathbf{z}_{t-1}) + \mathbf{C}\mathbf{s}_t + \mathbf{h}$$

$$\mathbf{A} = \text{diag}([a_{11} \dots a_{MM}]) \in \mathbb{R}^{M \times M}$$

$$\text{L}_{\text{reg}} = \tau_A \sum_{i=1}^{M_{\text{reg}}} (A_{i,i} - 1)^2 + \tau_W \sum_{i=1}^{M_{\text{reg}}} \sum_{\substack{j=1 \\ j \neq i}}^M W_{i,j}^2 + \tau_h \sum_{i=1}^{M_{\text{reg}}} h_i^2$$

penalty term $M_{\text{reg}} \leq M$

$$\begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & \times & 0 & 0 \\ 0 & 0 & 0 & 0 & \times & 0 \\ 0 & 0 & 0 & 0 & 0 & \times \end{pmatrix}$$

\mathbf{A}

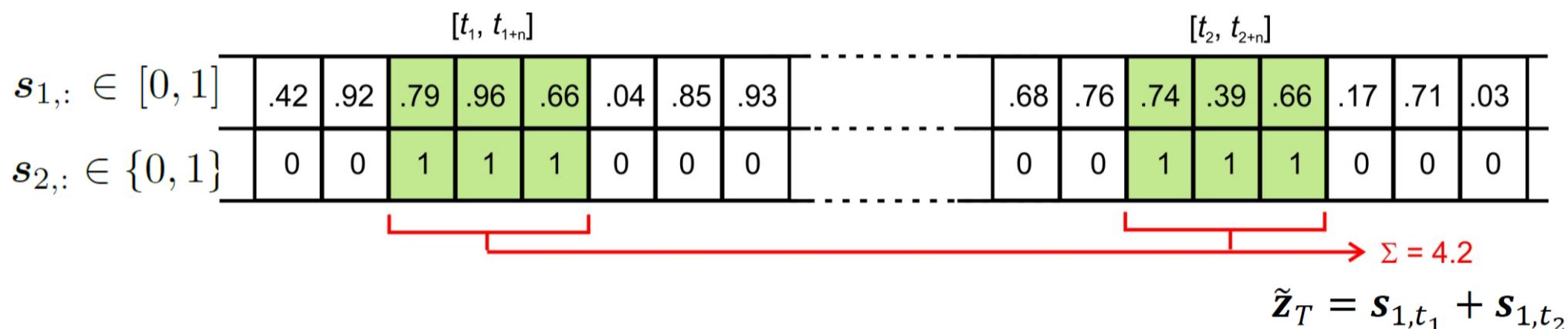
$$\begin{pmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ \times & \times & \times & 0 & \times & \times \\ \times & \times & \times & \times & 0 & \times \\ \times & \times & \times & \times & \times & 0 \end{pmatrix}$$

\mathbf{W}

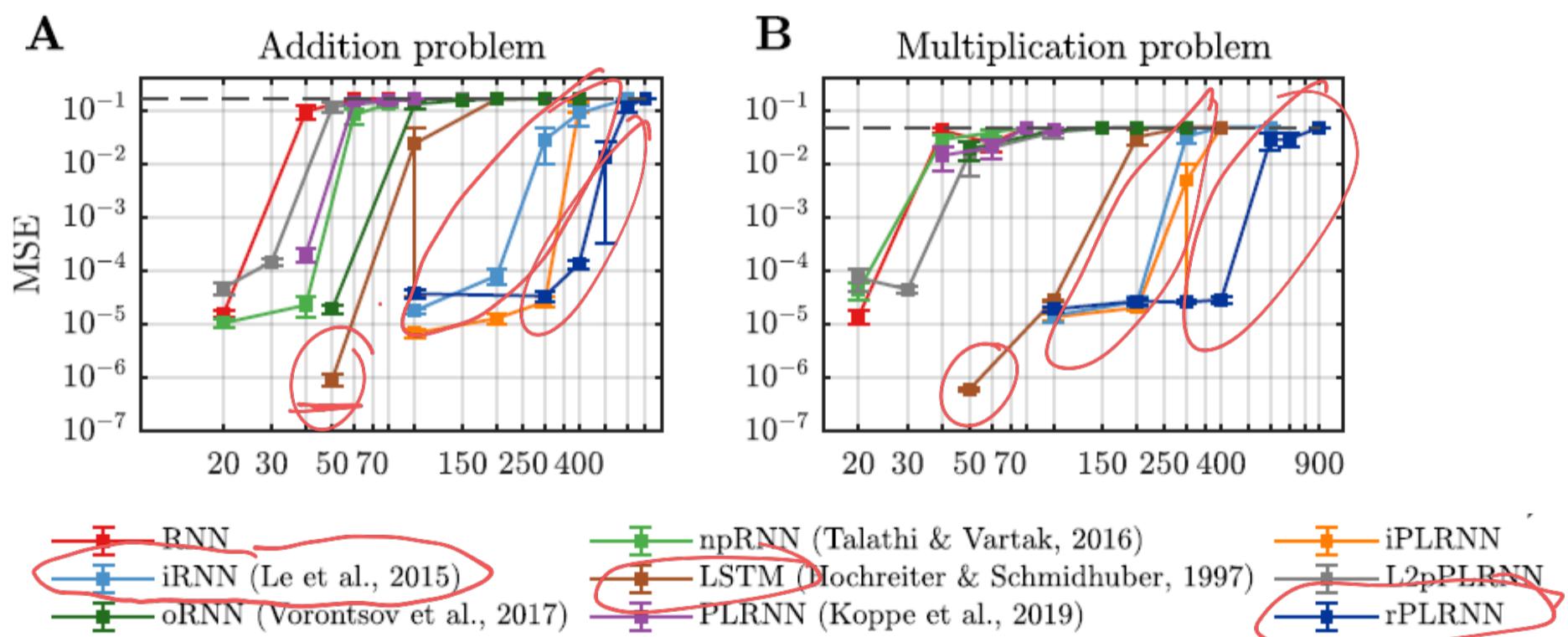
$$\begin{pmatrix} 0 \\ 0 \\ 0 \\ \times \\ \times \\ \times \end{pmatrix}$$

\mathbf{h}

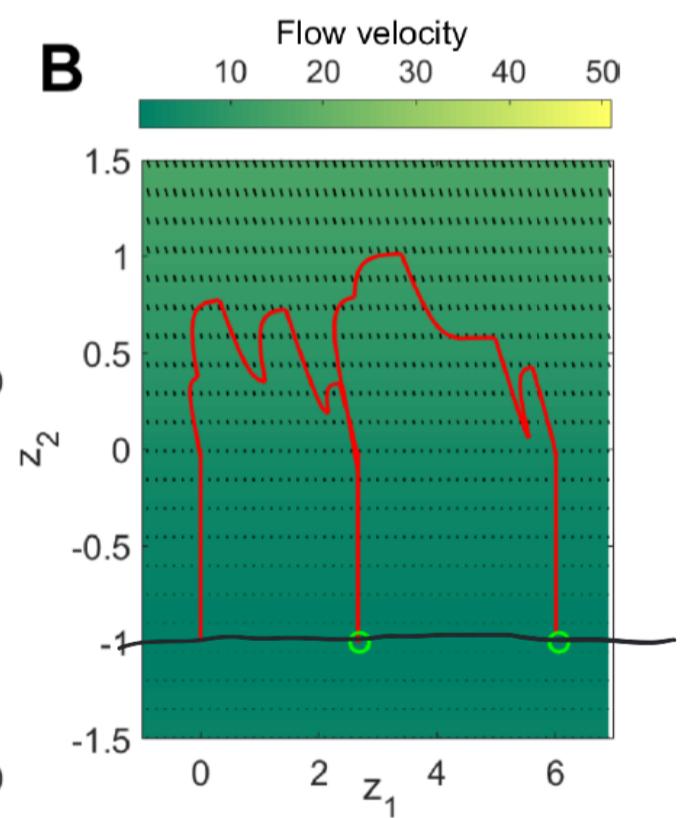
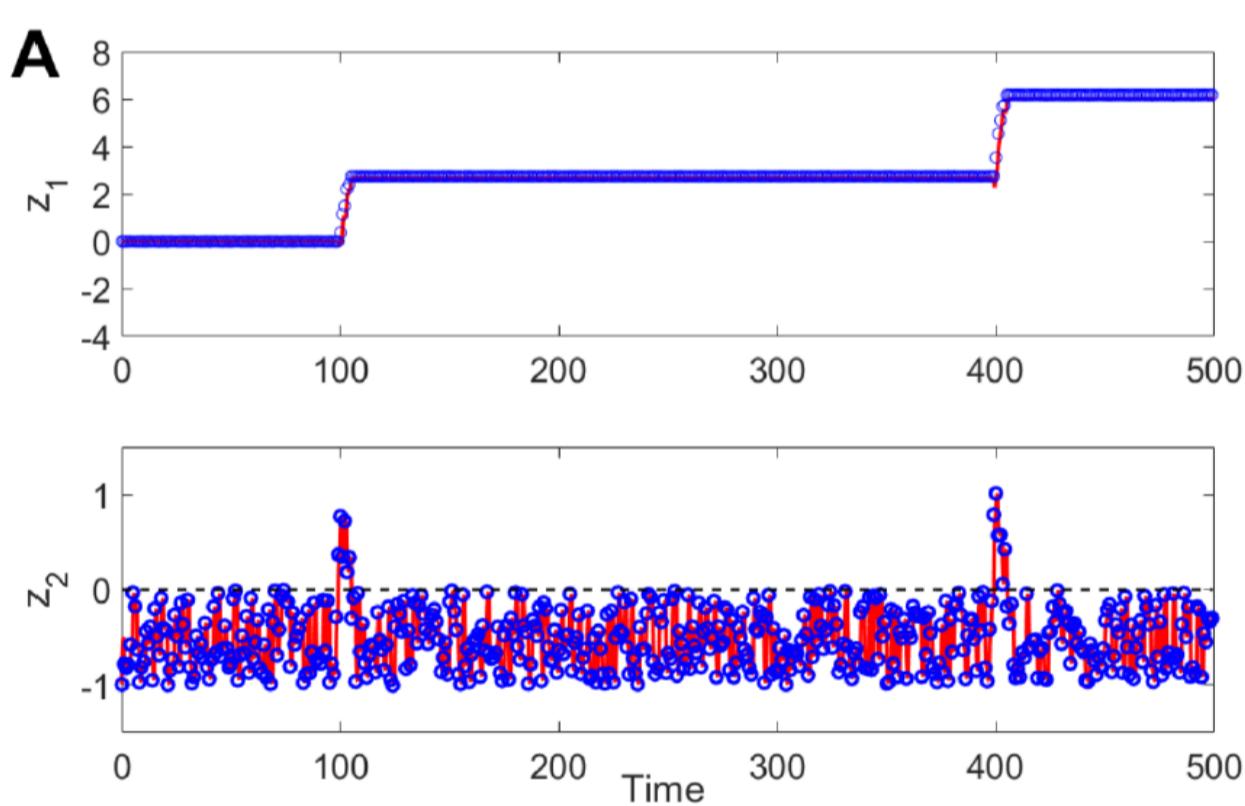
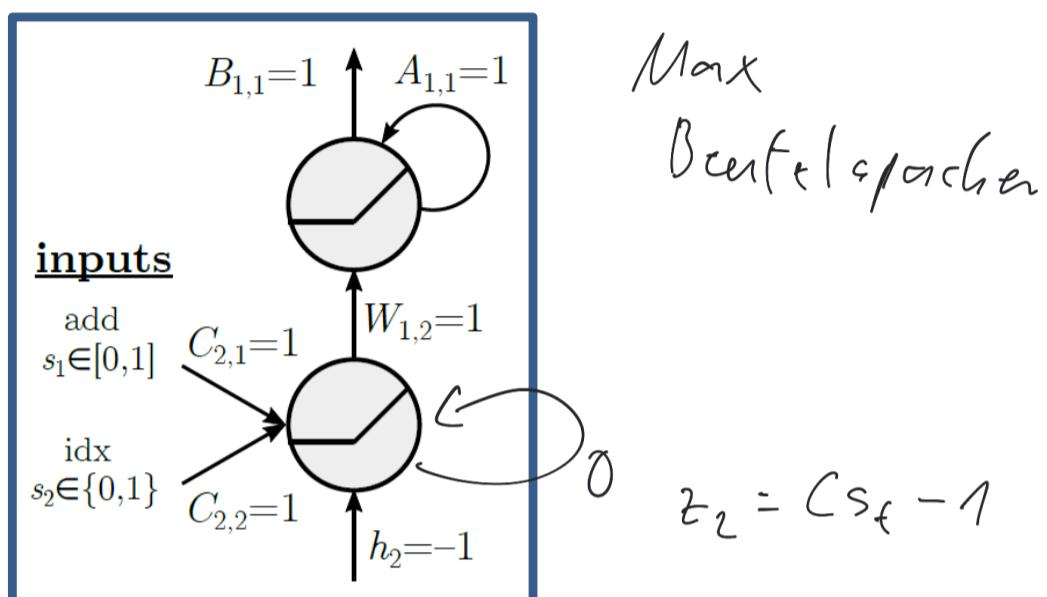
Machine Learning benchmark tasks for RNN



from Monfared & D (2020) ICML; after Hochreiter & Schmidhuber (1997) NC



Schmidt et al. (2021) ICLR

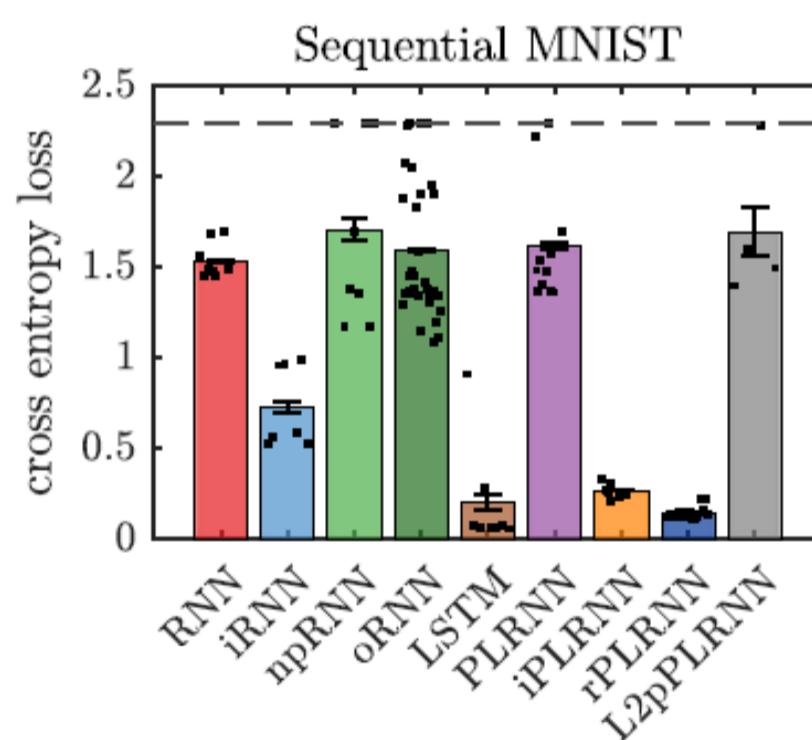
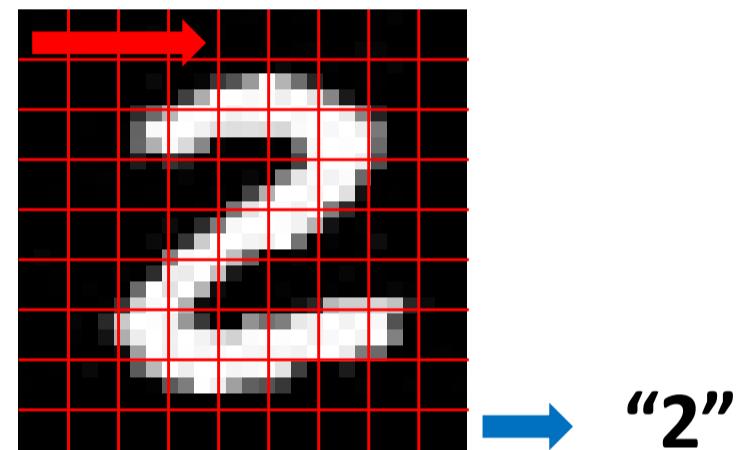


Schmidt et al. (2021) ICLR, Monfared & Durstewitz (2020) ICML

Sequential MNIST benchmark

0 0 0 0 0 0 0 0 0 0 0 0
1 1 1 1 1 1 1 1 1 1 1 1
2 2 2 2 2 2 2 2 2 2 2 2
3 3 3 3 3 3 3 3 3 3 3 3
4 4 4 4 4 4 4 4 4 4 4 4
5 5 5 5 5 5 5 5 5 5 5 5
6 6 6 6 6 6 6 6 6 6 6 6
7 7 7 7 7 7 7 7 7 7 7 7
8 8 8 8 8 8 8 8 8 8 8 8
9 9 9 9 9 9 9 9 9 9 9 9

By Josef Steppan - Own work, CC BY-SA 4.0,
<https://commons.wikimedia.org/w/index.php?curid=64810040>



w ortho-
orthogonal