



# TIME SERIES ANALYSIS & RECURRENT NEURAL NETWORKS

#9

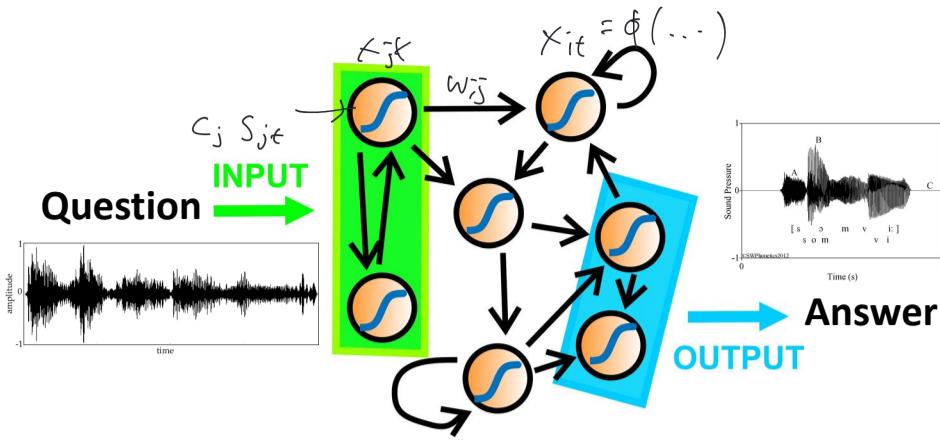
- Universal approx. theorem for RNN
- Training RNN by grad. descent,  
Real Time Recurrent Learning

**Main lecture:** Daniel Durstewitz

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

Heidelberg University

## Recurrent Neural Network



$$x_t = \phi(wx_{t-1} + h + Cs_t)$$

RNN are

- computat. universal: can implement any Th (Siegelmann & Sontag 95, Kozen 94, Cai et al. 2017)
- dynamically universal: can approx. any DS to arbitrary degree

Univ. Approx. Theorem (UAT) for RNN

- 1) UAT for FF networks → apply to the func.  $F(x)$  of DS
- 2) Reformulate the rec. FFN as a RNN

UAT for FFN (Cybenko 89)

Any contin., smooth func.  $F(x)$ ,  $x \in \mathbb{R}^N$ , can be approx. to precis.  $\varepsilon$  by the series

$$f_\theta(x) = \sum_{j=1}^n \alpha_j \sigma(b_j x + h_j)$$

$$\sigma(y) \rightarrow \begin{cases} 1 & \text{for } y \rightarrow +\infty \\ 0 & \text{for } y \rightarrow -\infty \end{cases}$$

sigmoid, differentiable

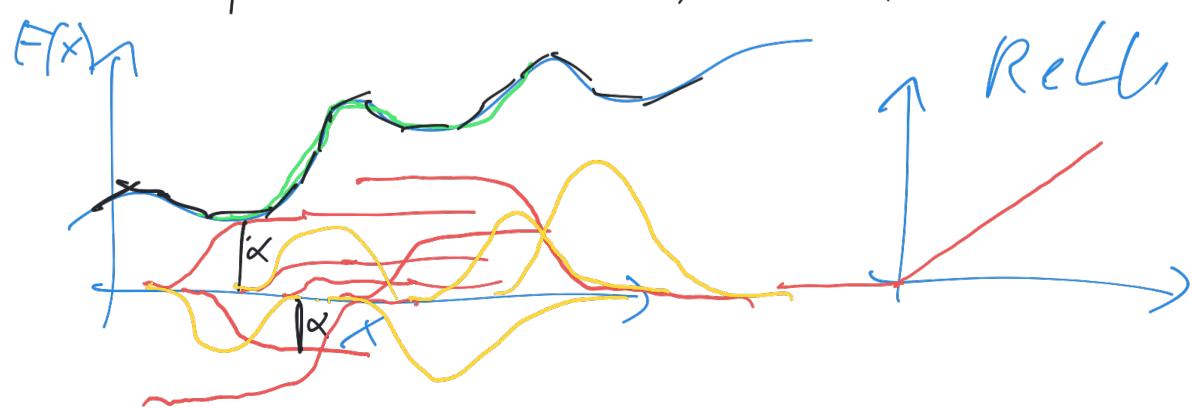
such that  $|F(x) - f_\theta(x)| < \varepsilon \quad \forall \varepsilon > 0$

for approp. choices of  $N, M$   
on a compact set of  $\mathbb{R}^N$

- Kolmogorov, Stone - Weierstrass

- Hornik (91)

- Chatzifotis et al. (2020) (MC, ReLU



## VAT for RNN

- approx. any DS wifv. closely
  - Funahashi & Nakamura (93) NN
  - Ueda & Nakano (98) NN
  - Tomso & Rajivsky (2020) PMLR

$$x_t = \phi(\underbrace{Wx_{t-1} + b}_{=: z_{t-1}})$$

$$\Rightarrow x_t = W^{-1}z_t - W^{-1}b$$

$$\Rightarrow W^{-1}(z_t - b) = \phi(z_{t-1})$$

$$\Rightarrow z_t = W\phi(z_{t-1}) + b$$

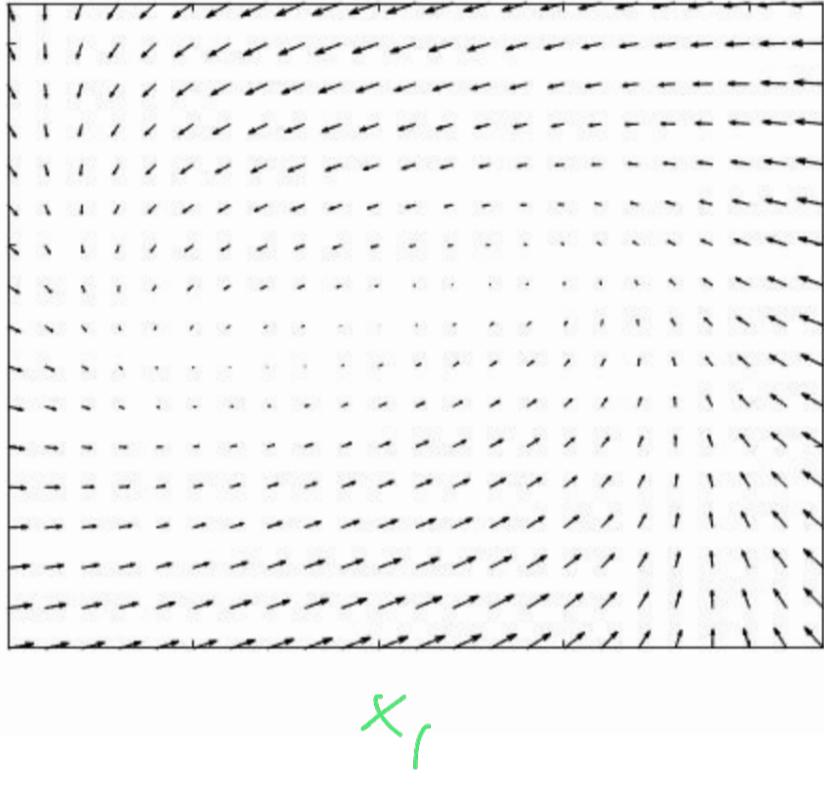
$$z_t = A^{\text{diag}} z_{t-1} + W\phi(z_{t-1}) + b$$

$$z_t = (1 - \Delta t) z_{t-1} + \Delta t W\phi(z_{t-1}) + \Delta t b$$

note  $\Delta t = 1$

$$\Rightarrow \frac{z_t - z_{t-1}}{\Delta t} = -z_{t-1} + W\phi(z_{t-1}) + b$$

lim  $\dot{z} = -z + W\phi(z) + b$   
 $\Delta t \rightarrow 0$



$x_1$

$x_2$

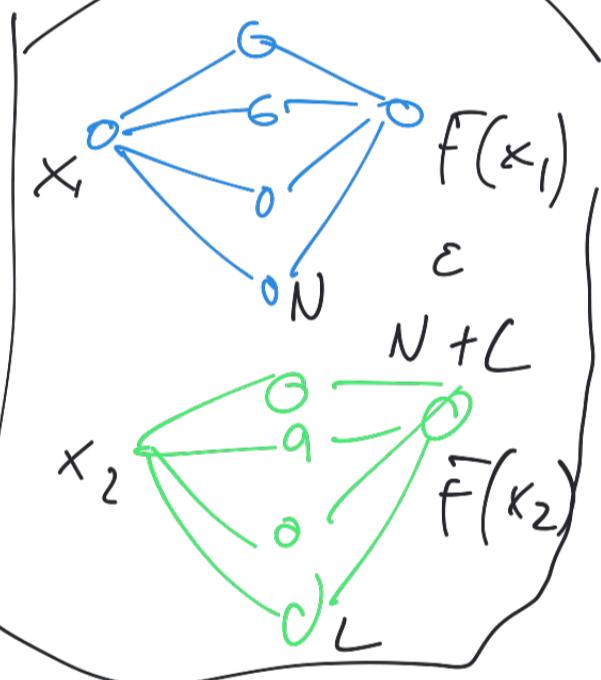
$$\dot{x} \equiv \frac{dx}{dt} = F(x)$$

1 hidden layer

FFN

$F_1(x)$

$F_2(x)$



Let  $F(x)$  be the flow (vector) field  
of a  $N$ -dim. DS  $\dot{x} = F(x)$ ,  $x \in \mathbb{R}^N$

looking for  $\underbrace{N \times M}_{\text{acc-UAT for FFNN}}$   $\underbrace{M \times N}_{\text{acc-UAT for FFNN}}$

$$|F(x) - A\phi(\beta x + b)| < \epsilon$$

everywhere,  $\forall \epsilon > 0$

Let  $\dot{y} = -\frac{1}{\tau}y + A\phi(\beta y + b)$ ,  $\epsilon$  very large

Training RNN

Supervised tasks

$$\left\{ s_t^{(p)} \right\}_{t \in T} \xrightarrow{\quad} \left\{ \tilde{x}_t \right\}_{t \in T}$$

# Real-Time Recurrent Learning (RTRL)

→ online learning algo.

Williams & Zipser (1990)

Pearlmutter (1989)

$$D = \{ \{ s_{t \in T}^{(p)} \}, \{ \tilde{x}_{t \in T}^{(p)} \} \}, p = 1 \dots P$$

$$x_t = \phi(w_{x_{t-1}} + b + \cancel{s_t}), i = 1 \dots M$$

M-dim. col. vec.  
units

Loss func.

$$\text{error signal } \varepsilon_{it} = \begin{cases} \tilde{x}_{it} - x_{it} & \text{if } \tilde{x}_{it} \\ 0 & \text{otherwise} \end{cases}$$

$$\cancel{L(w, b)} := \frac{1}{2} \sum_{p=1}^P \sum_{t=1}^T \sum_{i=1}^M \varepsilon_{it}^2$$

$$= \frac{1}{2} \sum_{t=1}^T \sum_{i=1}^M I\{\varepsilon_{it} \neq 0\} (\tilde{x}_{it} - x_{it})^2$$

$$\frac{\partial L(w)}{\partial w_{ij}} = \frac{1}{2} \sum_{t=1}^T \frac{\partial L_t(w)}{\partial w_{ij}}$$

$$\frac{\partial L_t(w)}{\partial w_{ij}} = \sum_{k=1}^M I\{\varepsilon_{kt} \neq 0\} (\tilde{x}_{kt} - x_{kt}) \frac{\partial x_{kt}}{\partial w_{ij}}$$

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

$(fg)' = f'g + fg'$   
 for  $i = k$

Diagram illustrating the computation of  $\frac{\partial x_{kt}}{\partial w_{ij}}$ :

- The diagram shows a node at time  $t-1$  with index  $i$ . It receives inputs from node  $k$  at time  $t-1$  via weight  $w_{kj}$  and from node  $j$  at time  $t-1$  via weight  $w_{ji}$ .
- The output of this node is  $\phi'(\sum_l w_{kl} x_{l,t-1} + b_i)$ .
- The diagram also shows the calculation of  $\frac{\partial x_{l,t-1}}{\partial w_{ij}}$ , which involves summing contributions from all nodes  $k$  via weight  $w_{kj}$ .

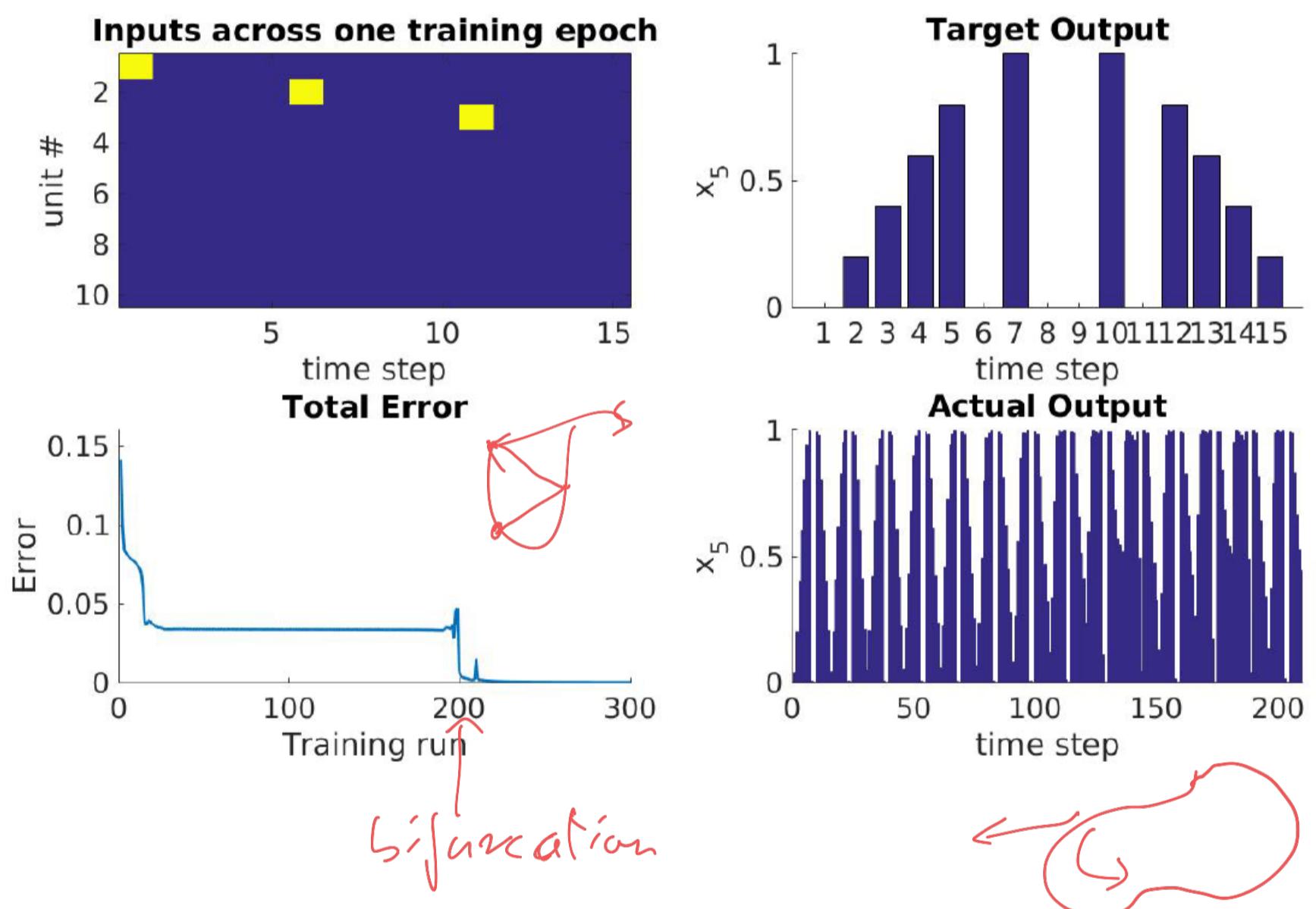
## RTRL algo.

- init. cond.:  $\gamma_{lij}^{(0)} = 0 \quad \forall (l, i, j)$

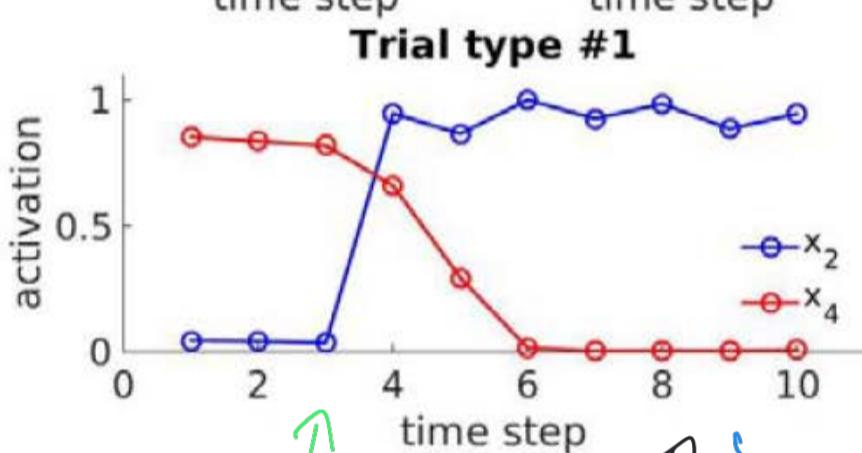
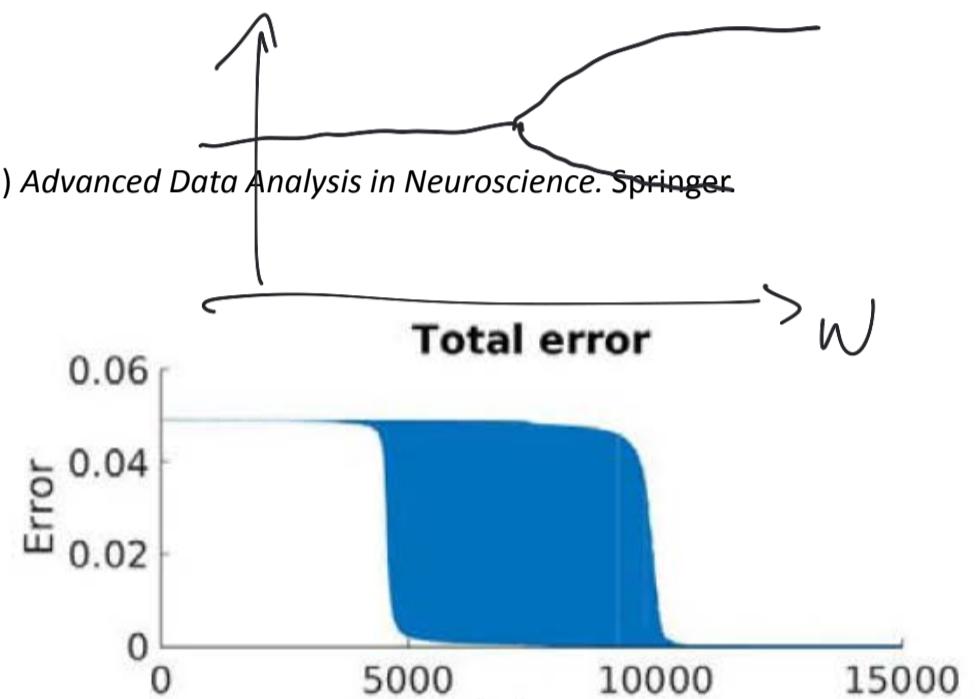
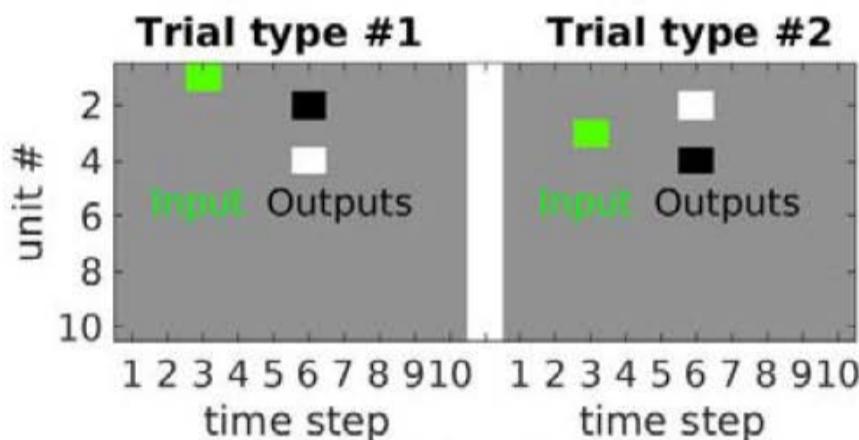
- update  $\gamma_{lij}^{(t)}$  along with act. val.  $x_l^{(t)}$

- update  $w_{ij}^{(t+1)} = w_{ij}^{(t)} - \alpha \frac{\partial L_t}{\partial w_{ij}}$

Source: Durstewitz (2017) *Advanced Data Analysis in Neuroscience*. Springer.

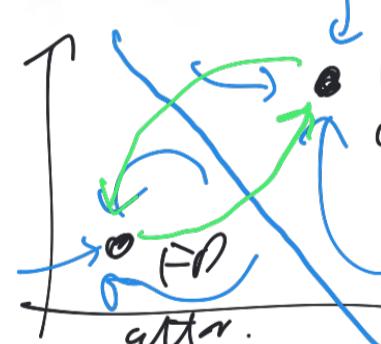


Source: Durstewitz (2017) *Advanced Data Analysis in Neuroscience*. Springer.



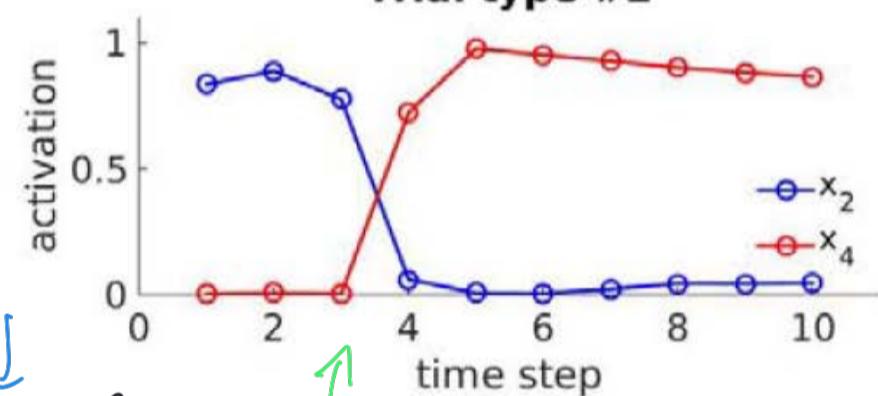
↑  
u#1

$x_2$



F P  
att. -

$x_1$



↑  
u#3