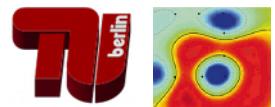
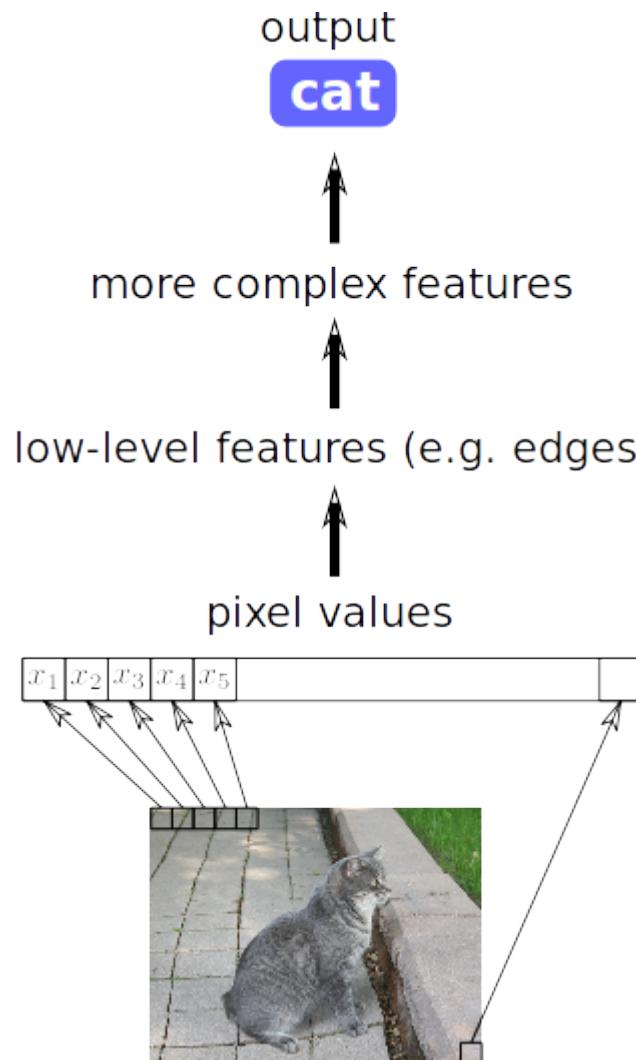

SoSe 2018: Deep Neural Networks

Lecture 8: Stateful DNNs

Machine Learning Group
Technische Universität Berlin



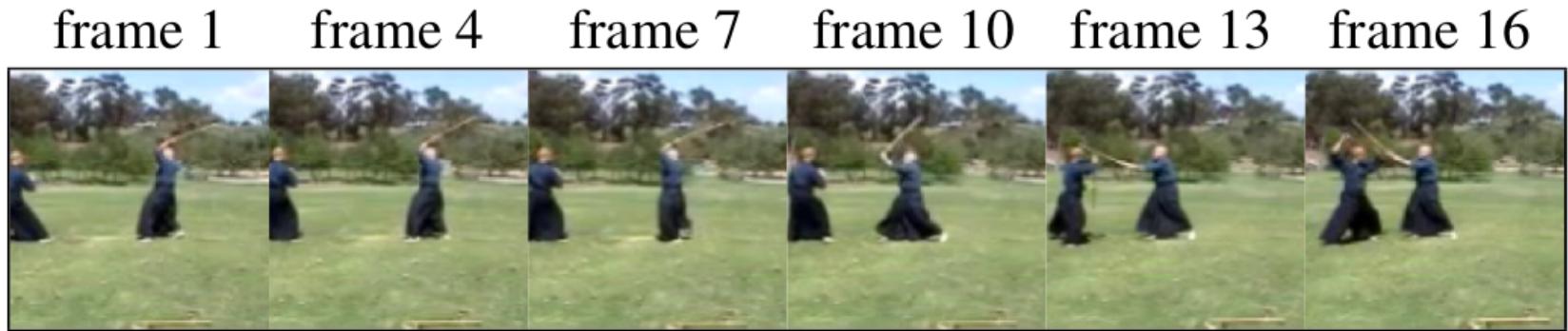
Recap: Feedforward Neural Networks



Example:

Prediction from a single static image.

Neural Networks for Time



Many signals (e.g. videos, motion data, speech, text, ...) have a time component.

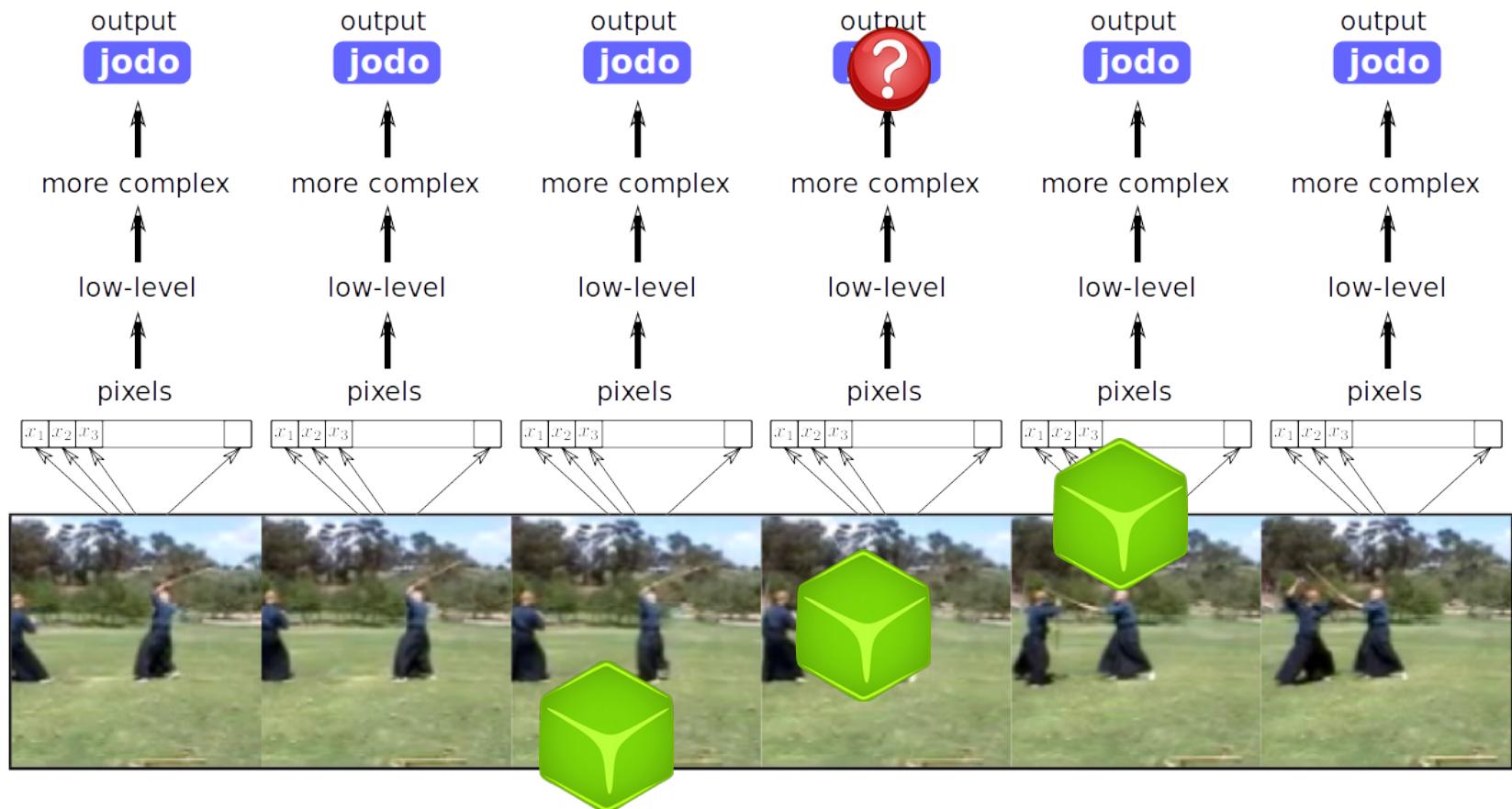
Neural Networks for Time

Basic approach, predict each time step individually.



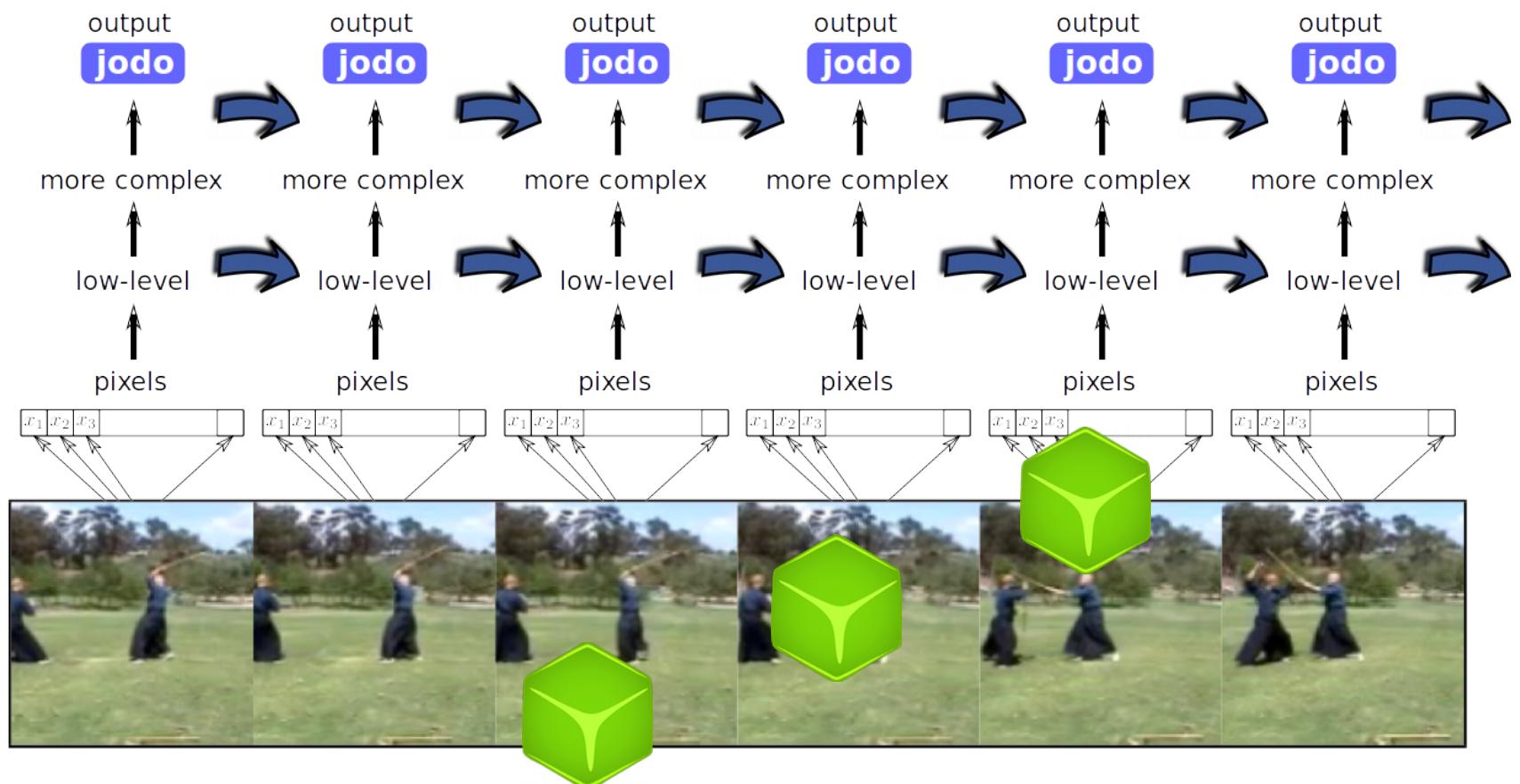
Neural Networks for Time

Basic approach not robust to e.g. occlusion.

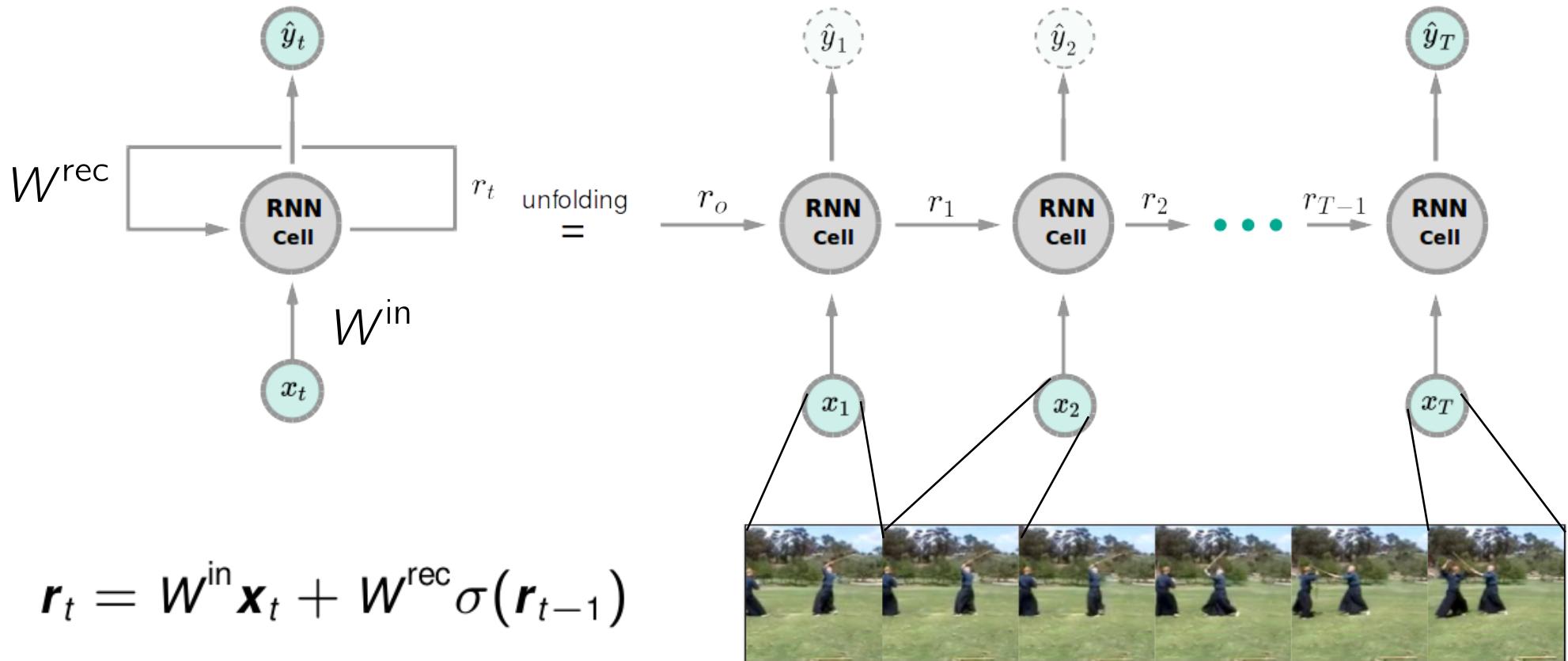


Neural Networks for Time

Idea: incorporate recurrent connections that maintain underlying dynamics.



RNN Basics



RNN Basics (Forward Pass)

Assume that $\theta = \{W^{\text{in}}, W^{\text{rec}}\}$, $r_0 = \mathbf{0}$, and the **forward pass** is constructed as follows:

$$r_1 = W^{\text{in}}x_1 + W^{\text{rec}}\sigma(r_0)$$

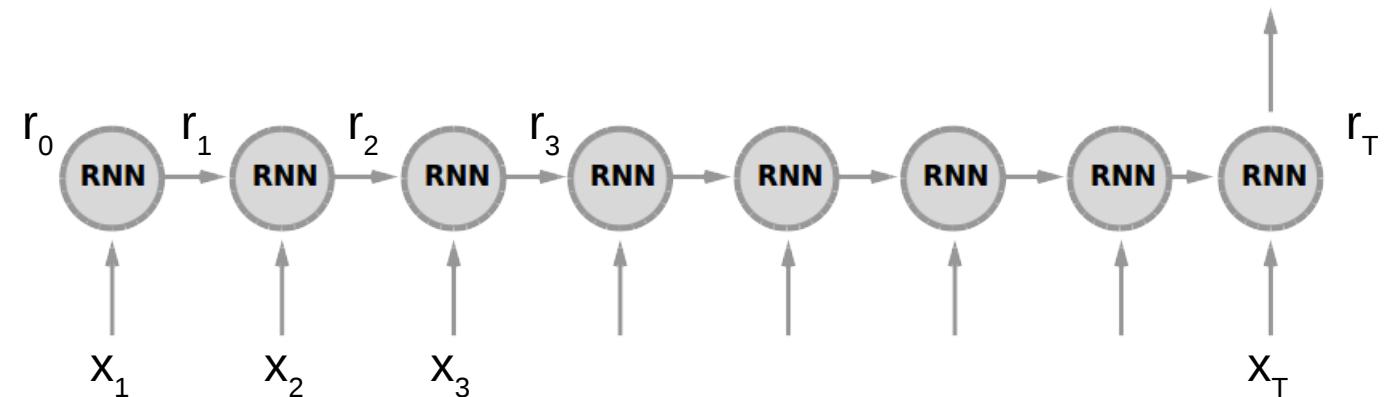
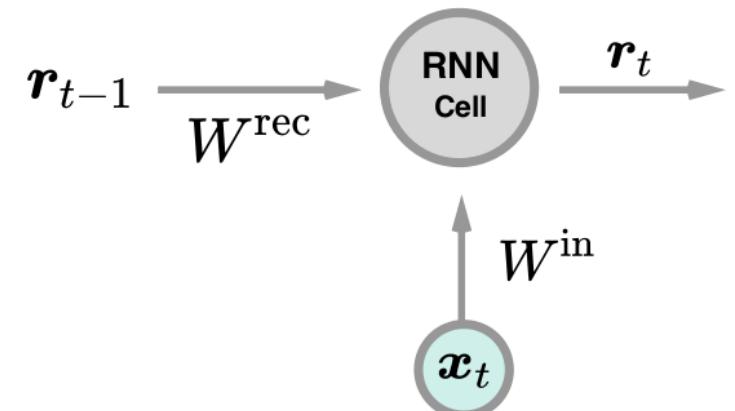
⋮

$$r_t = W^{\text{in}}x_t + W^{\text{rec}}\sigma(r_{t-1})$$

⋮

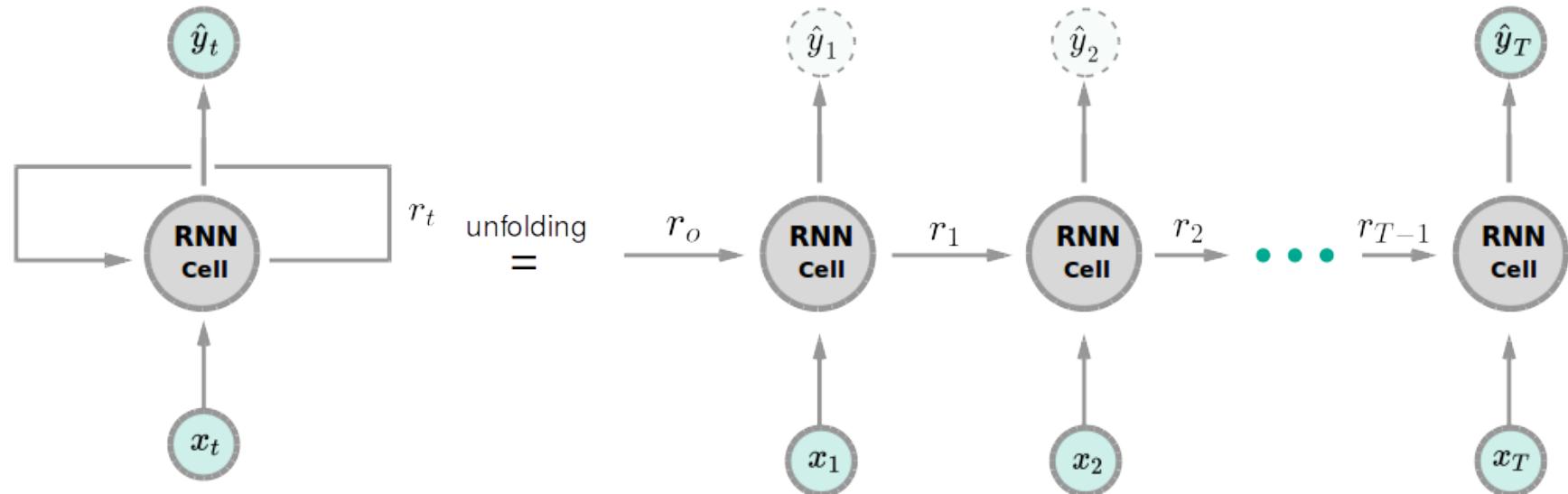
$$r_T = W^{\text{in}}x_T + W^{\text{rec}}\sigma(r_{T-1})$$

$$\hat{y} = r_T$$

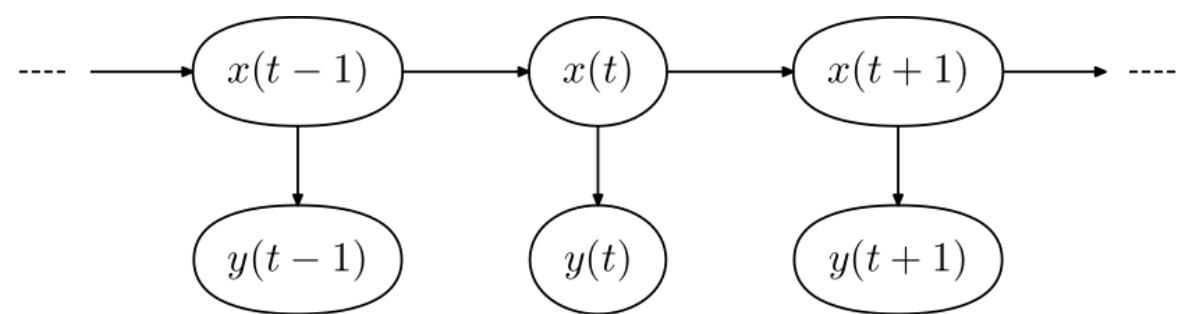
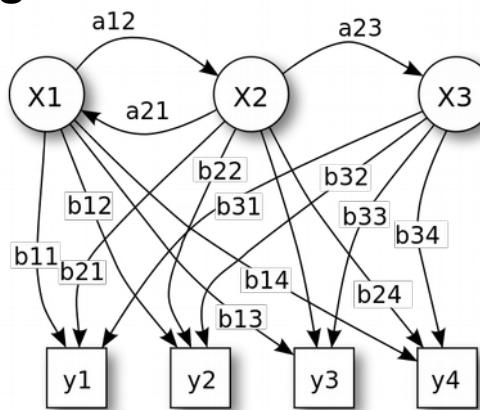


RNNs vs. Hidden Markov Models (HMMs)

RNNs



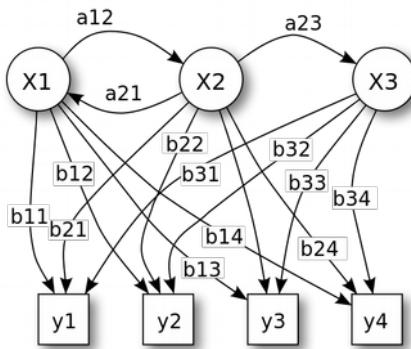
HMMs



Recurrent Neural Networks vs. HMMs

HMMs

Image source
wikipedia.org
user: Tdunning



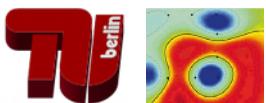
Discrete state space

$$S = \{s_1, \dots, s_h\}$$

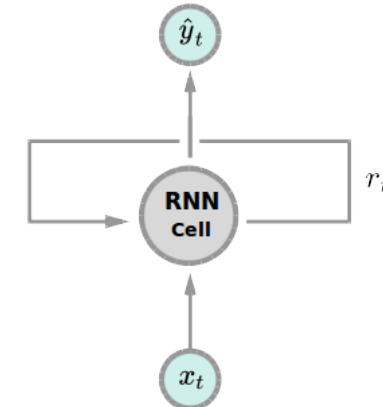
Probabilistic transitions with
transition model encoded in
a matrix

$$A : h \times h$$

No inputs



RNNs



Continuous state space

$$S = \mathbb{R}^h$$

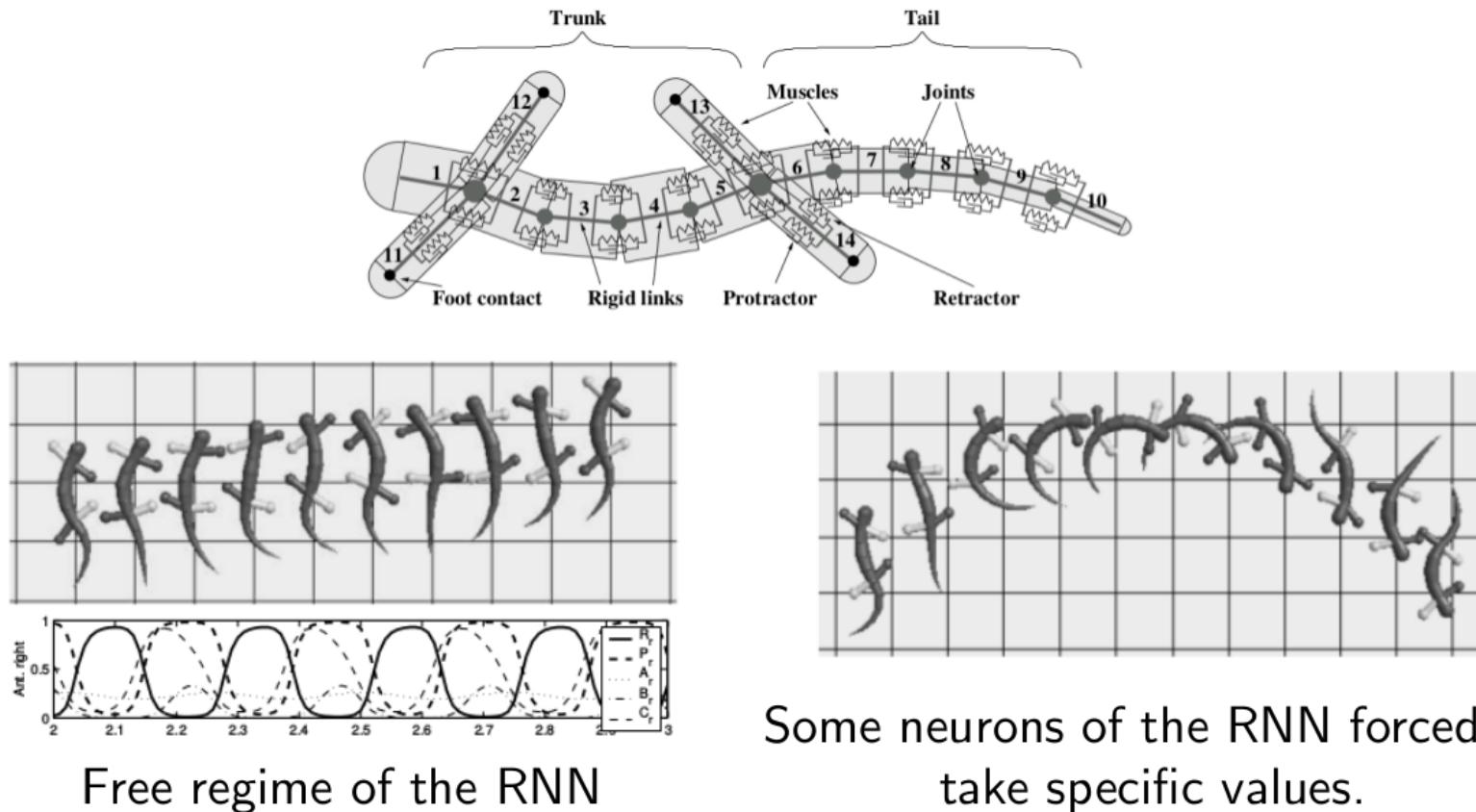
Deterministic transitions
with transition model
encoded by a parameterized
function

$$f : \mathbb{R}^h \rightarrow \mathbb{R}^h$$

Receives inputs

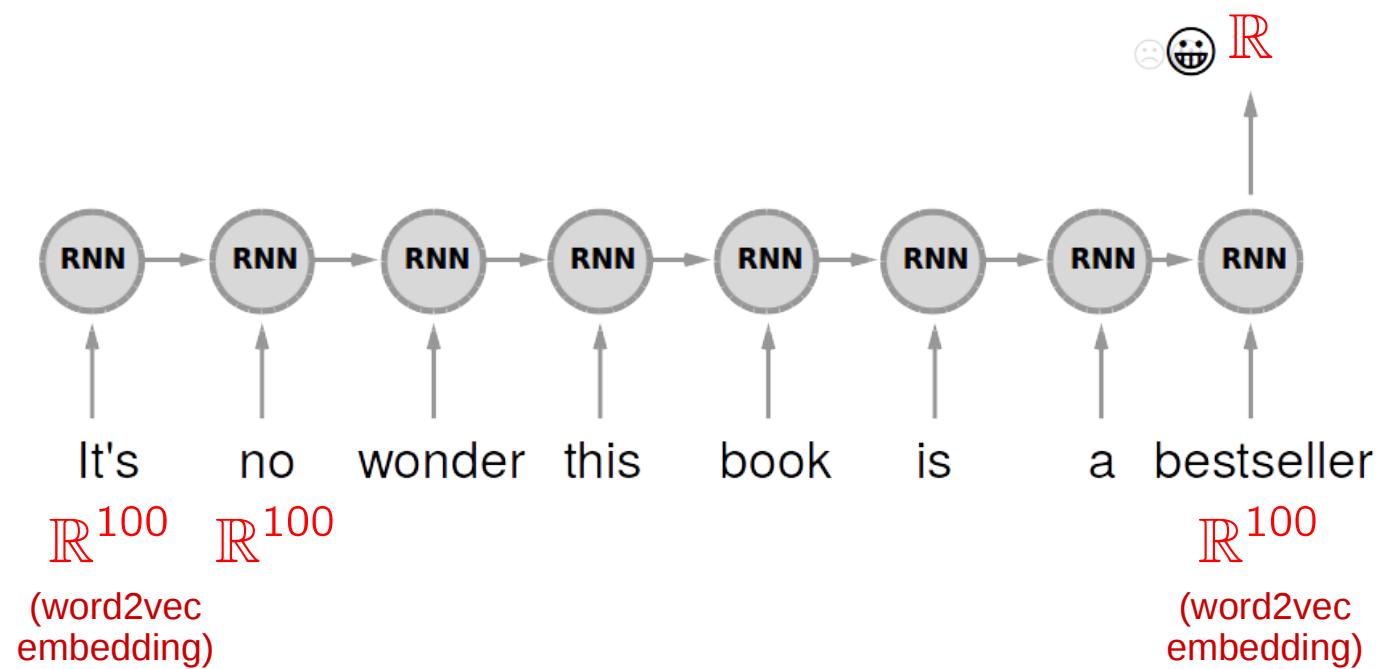
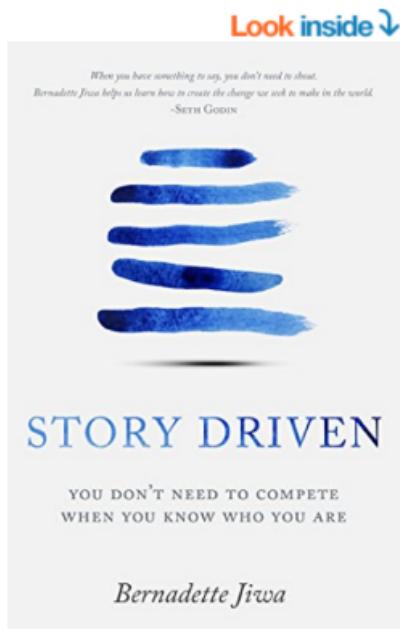
RNN Application: Salamander Motion

Experiments by Ijspeert et al. (2001):

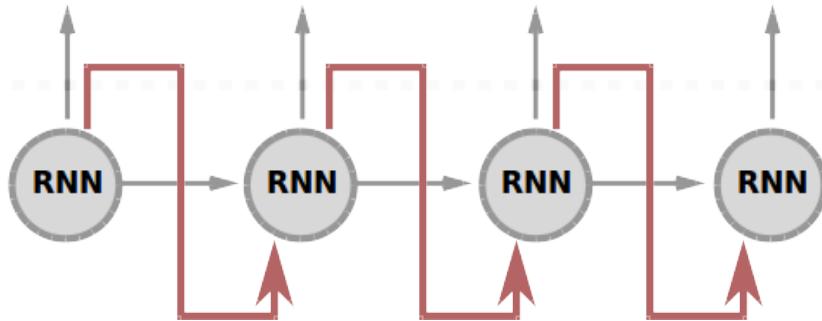


Images from Ijspeert et al. 2001, A Connectionist Central Pattern Generator for the Aquatic and Terrestrial Gaits of a Simulated Salamander

RNN Application: Sentiment Analysis



RNN Application: Random Text Generation



RNN output can be given as input to condition generation on what has been produced at previous time steps.

Lemma 0.3. Let S be a scheme. Let X be a scheme and X is an affine open covering. Let $\mathcal{U} \subset \mathcal{X}$ be a canonical and locally of finite type. Let X be a scheme. Let X be a scheme which is equal to the formal complex.

The following to the construction of the lemma follows.

Let X be a scheme. Let X be a scheme covering. Let

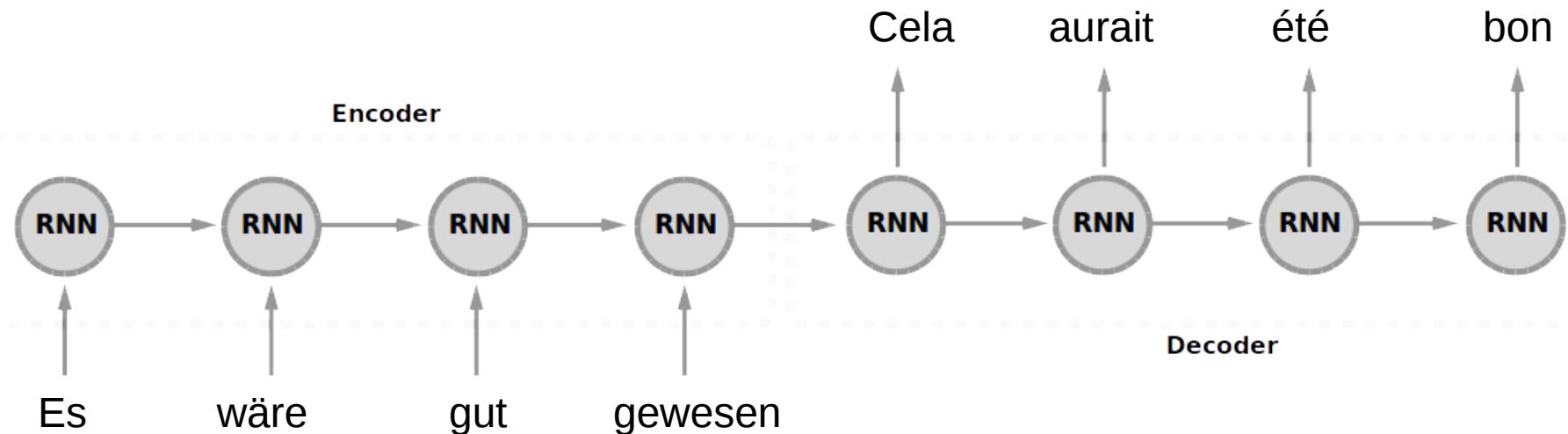
$$b : X \rightarrow Y' \rightarrow Y \rightarrow Y \rightarrow Y' \times_X Y \rightarrow X.$$

be a morphism of algebraic spaces over S and Y .

Proof. Let X be a nonzero scheme of X . Let X be an algebraic space. Let \mathcal{F} be a quasi-coherent sheaf of \mathcal{O}_X -modules. The following are equivalent

- (1) \mathcal{F} is an algebraic space over S .
- (2) If X is an affine open covering.

RNN Application: Machine Translation



\mathbb{R}^{100}

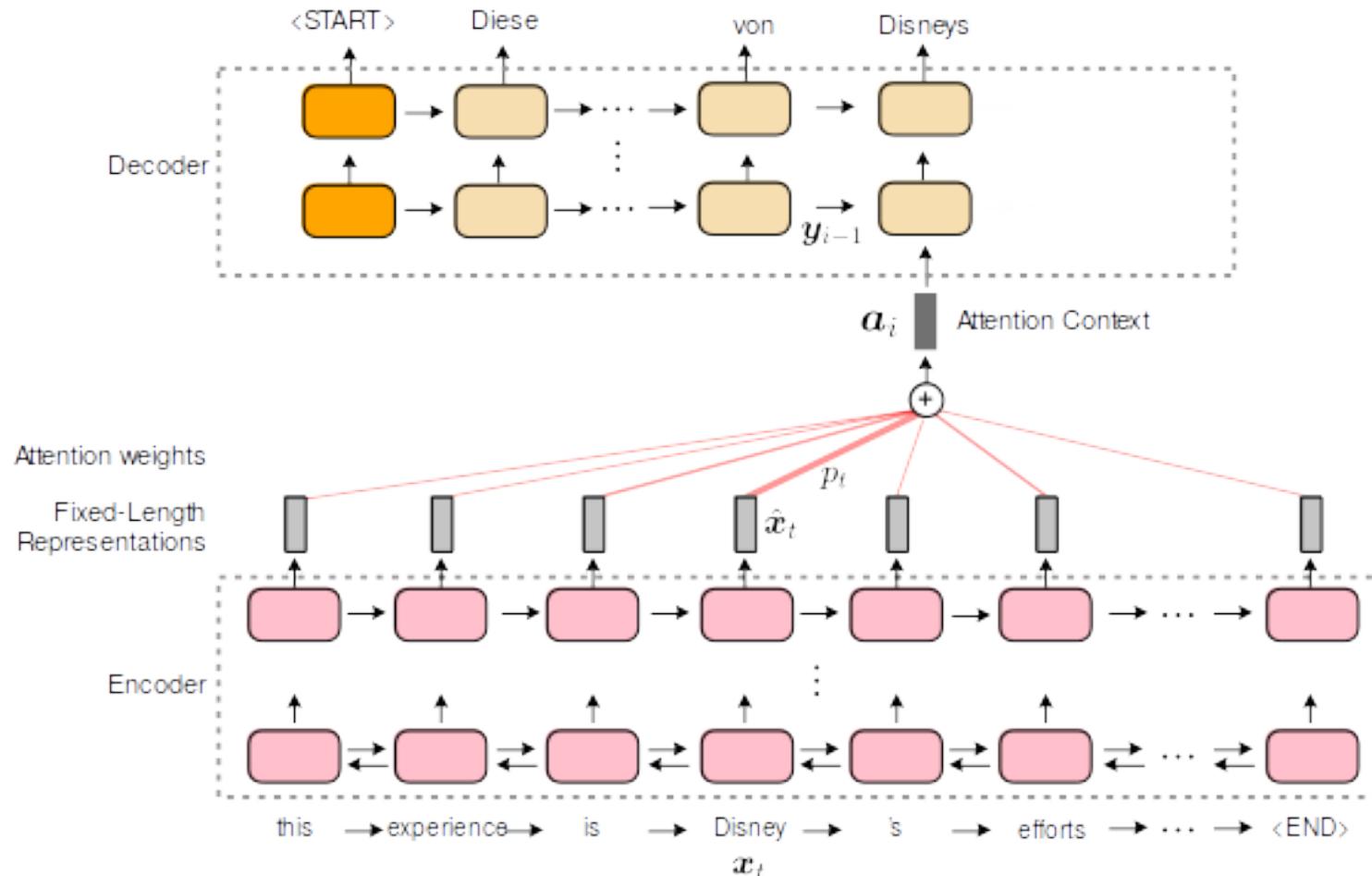
(word2vec
embedding)

Notes:

- Need to finish reading the sentence to start the translation.
- Complete sentence meaning is being encoded in the RNN state.

Machine Translation

- Google's Neural Machine Translation System [WSC + 16].



Source This experience is part of Disney's efforts to enlarge the audience.

Target Diese Erfahrung ist Teil von Disneys Bemühungen, das Publikum zu vergrößern.

Machine Translation

- Google's Neural Machine Translation System [Wu'16].

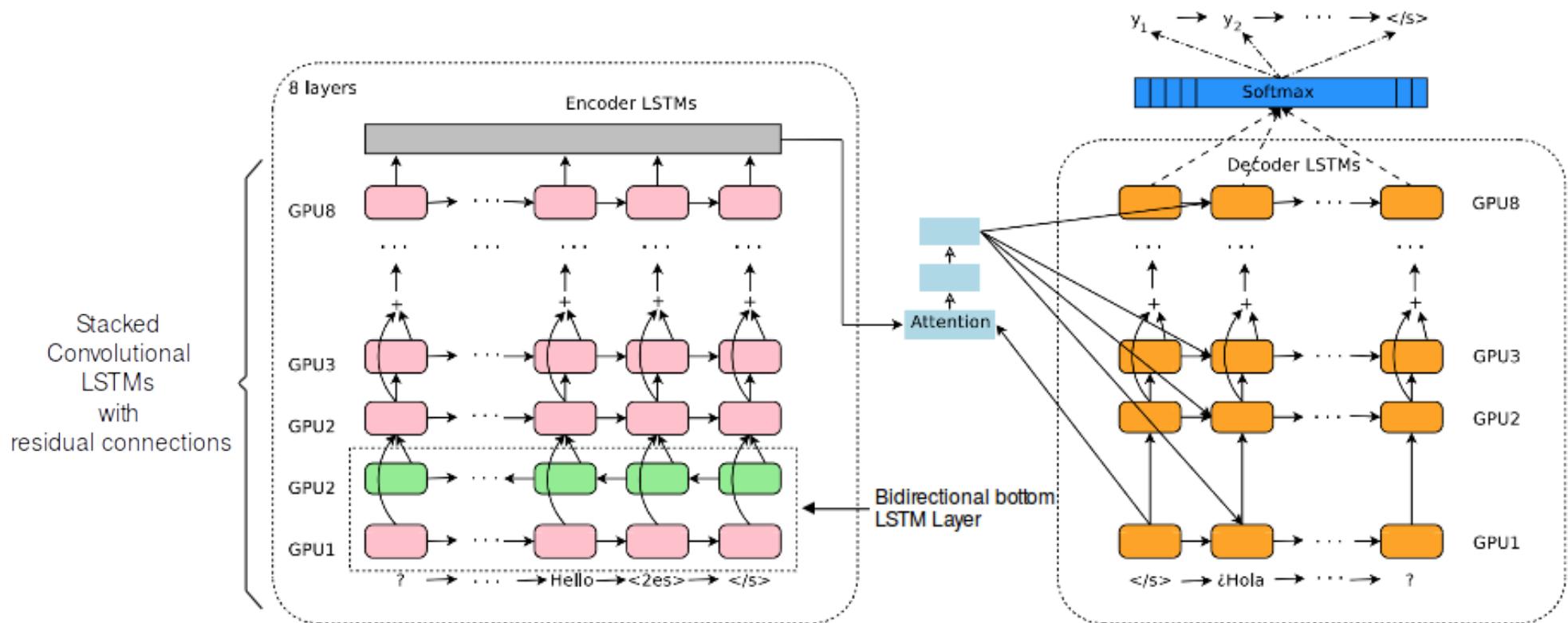


Image from Wu'16: Google's Neural Machine Translation System: Bridging the Gap between Human and Machine Translation

RNNs for Image Segmentation

- An RNN for image segmentation [McIntosh'17].

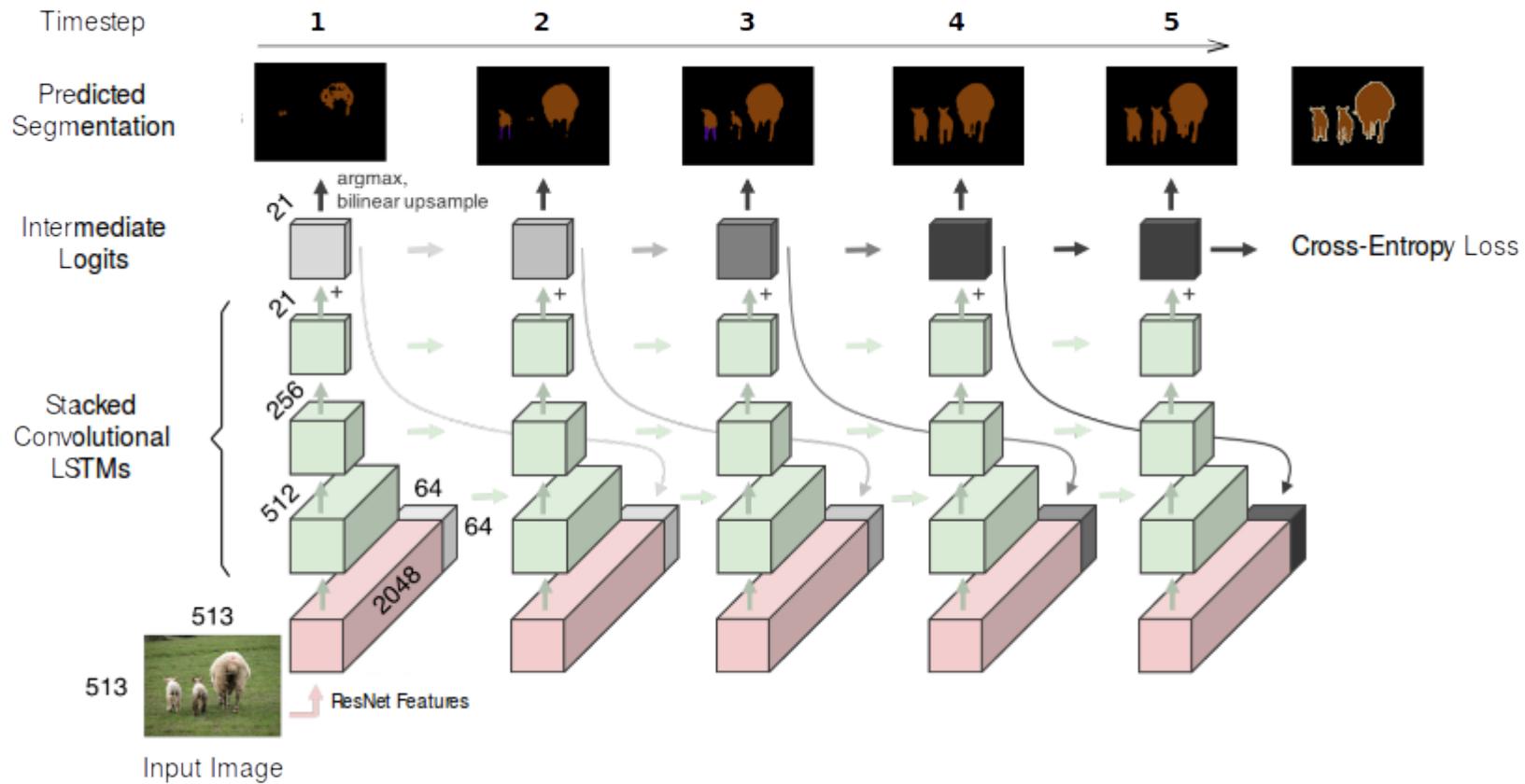


Image from McIntosh'17 Recurrent Segmentation for Variable Computational Budgets

RNN Basics Recap (Forward Pass)

Assume that $\theta = \{W^{\text{in}}, W^{\text{rec}}\}$, $r_0 = \mathbf{0}$, and the **forward pass** is constructed as follows:

$$r_1 = W^{\text{in}}x_1 + W^{\text{rec}}\sigma(r_0)$$

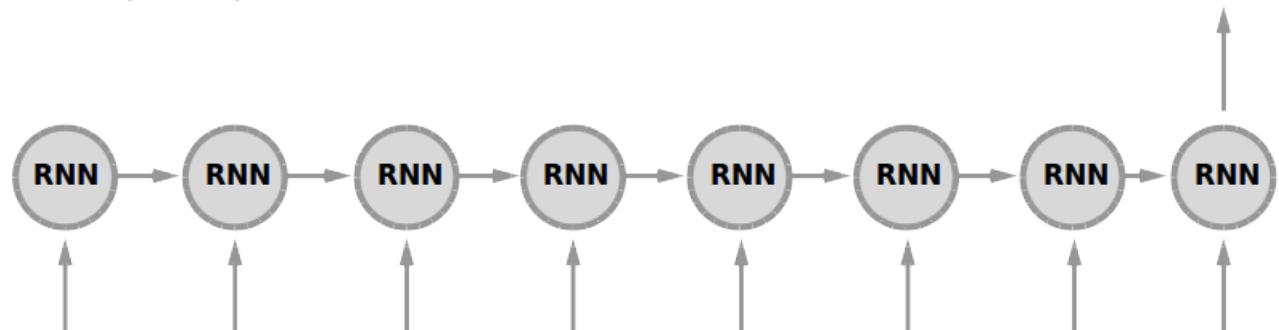
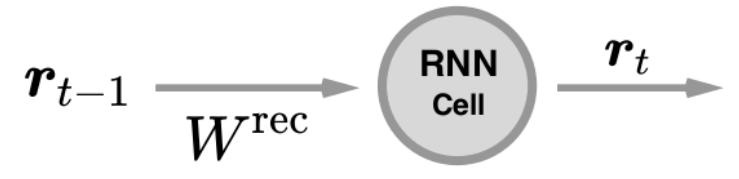
⋮

$$r_t = W^{\text{in}}x_t + W^{\text{rec}}\sigma(r_{t-1})$$

⋮

$$r_T = W^{\text{in}}x_T + W^{\text{rec}}\sigma(r_{T-1})$$

$$\hat{y} = r_T$$



Extracting a RNN from a Dynamical System

Example: Linear dynamical system

$$\frac{\partial}{\partial t} \mathbf{r}(t) = A \cdot \mathbf{r}(t)$$

Step 1: Euler discretization

$$\frac{\mathbf{r}(t_{n+1}) - \mathbf{r}(t_n)}{t_{n+1} - t_n} = A \cdot \mathbf{r}(t_n)$$

Step 2: Isolate the last time step

$$\mathbf{r}(t_{n+1}) = \mathbf{r}(t_n) + (t_{n+1} - t_n) \cdot A \cdot \mathbf{r}(t_n)$$

Step 3: Write in standard format

$$\mathbf{r}(t_{n+1}) = W_{\text{rec}} \cdot \mathbf{r}(t_n) \quad W_{\text{rec}} = I + (t_{n+1} - t_n) \cdot A$$

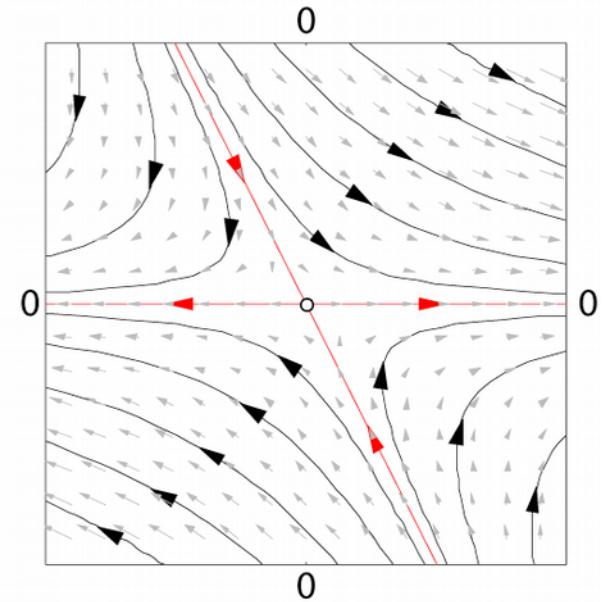


Image source: Wikipedia, user:
XaosBits

Training the RNN: Computing the Gradient

$$\mathbf{r}_1 = W^{\text{in}} \mathbf{x}_1 + W^{\text{rec}} \sigma(\mathbf{r}_0)$$

:

$$\mathbf{r}_t = W^{\text{in}} \mathbf{x}_t + W^{\text{rec}} \sigma(\mathbf{r}_{t-1})$$

:

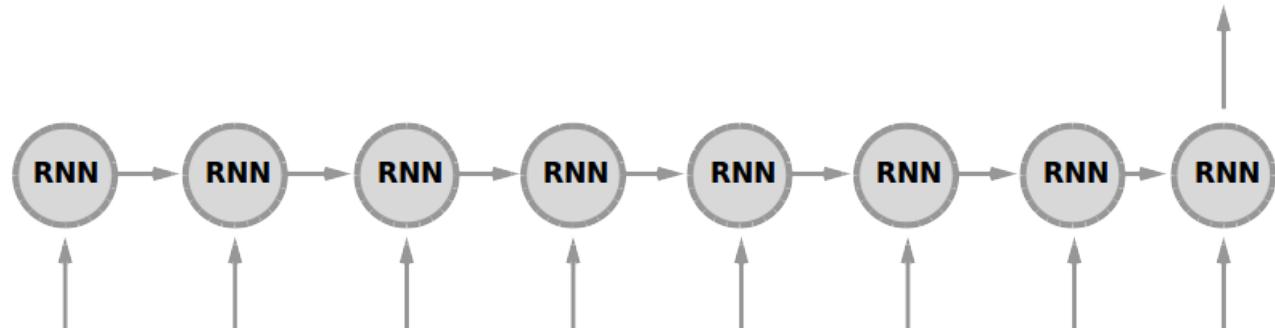
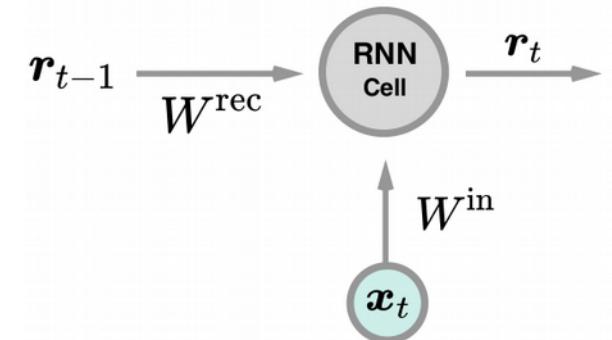
$$\mathbf{r}_T = W^{\text{in}} \mathbf{x}_T + W^{\text{rec}} \sigma(\mathbf{r}_{T-1})$$

$$\hat{\mathbf{y}} = \mathbf{r}_T$$

$$\frac{\partial E}{\partial \theta} = \sum_{k=1}^T \frac{\partial E}{\partial \mathbf{r}_T} \frac{\partial \mathbf{r}_T}{\partial \mathbf{r}_k} \frac{\partial^+ \mathbf{r}_k}{\partial \theta}$$

$$\frac{\partial \mathbf{r}_T}{\partial \mathbf{r}_k} = \prod_{T \geq i > k} \frac{\partial \mathbf{r}_i}{\partial \mathbf{r}_{i-1}}$$

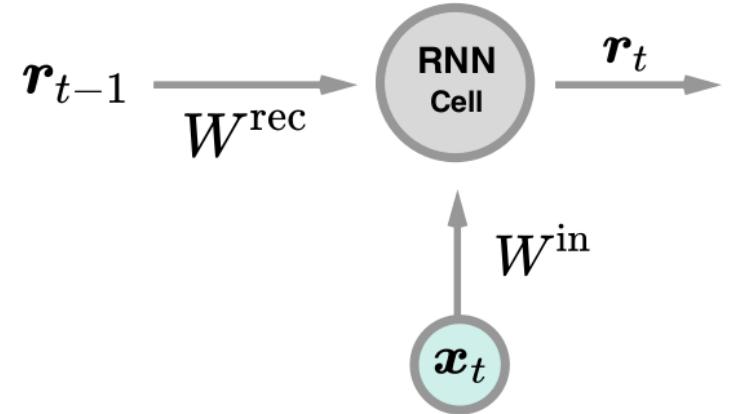
$$= \prod_{T \geq i > k} \underbrace{(W^{\text{rec}})^T}_{!} \text{diag}(\sigma'(\mathbf{r}_{i-1}))$$



Optimizing RNNs: Eigenvalues View

$$\frac{\partial E}{\partial \theta} = \sum_{k=1}^T \frac{\partial E}{\partial \mathbf{r}_T} \frac{\partial \mathbf{r}_T}{\partial \mathbf{r}_k} \frac{\partial^+ \mathbf{r}_k}{\partial \theta}$$

$$\begin{aligned}\frac{\partial \mathbf{r}_T}{\partial \mathbf{r}_k} &= \prod_{T \geq i > k} \frac{\partial \mathbf{r}_i}{\partial \mathbf{r}_{i-1}} \\ &= \prod_{T \geq i > k} \underbrace{(W^{\text{rec}})^T}_{!} \text{diag}(\sigma'(\mathbf{r}_{i-1}))\end{aligned}$$



$\sigma' = 1$ assume the activations are linear

$$\prod_{T \geq i > k} (W^{\text{rec}})^k = U \cdot \Lambda^{(T-k)} \cdot V$$

Optimizing RNNs: Eigenvalues View



$\sigma' = 1$ assume the activations are linear

$$\prod_{T \geq i > k} (W^{\text{rec}})^k = U \cdot \underbrace{\Lambda^{(T-k)} \cdot V}_{\lambda_1, \dots, \lambda_d}$$

Observations:

- Disbalance of gradient magnitude between different time steps

$$\frac{\lambda_1}{\lambda_1^{(T-k)}}$$

- Bad conditioning at lower-time steps

$$\left(\frac{\lambda_1}{\lambda_d}\right)^{(T-k)}$$

Optimizing RNNs: Eigenvalues View

Solution:

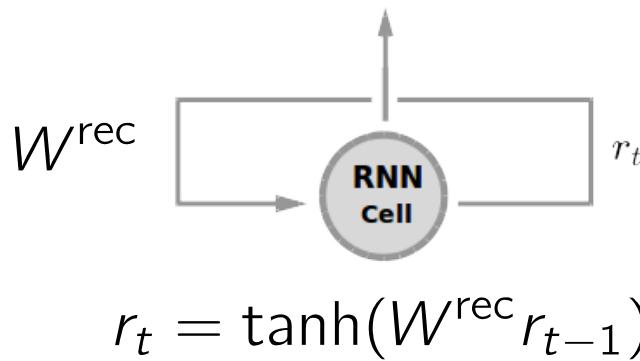
$\frac{\lambda_1}{\lambda_1^{(T-k)}}$ → Normalize weights to make first eigenvalue a bit smaller than 1.

$\left(\frac{\lambda_1}{\lambda_d}\right)^{(T-k)}$ → Add a LOT of momentum.

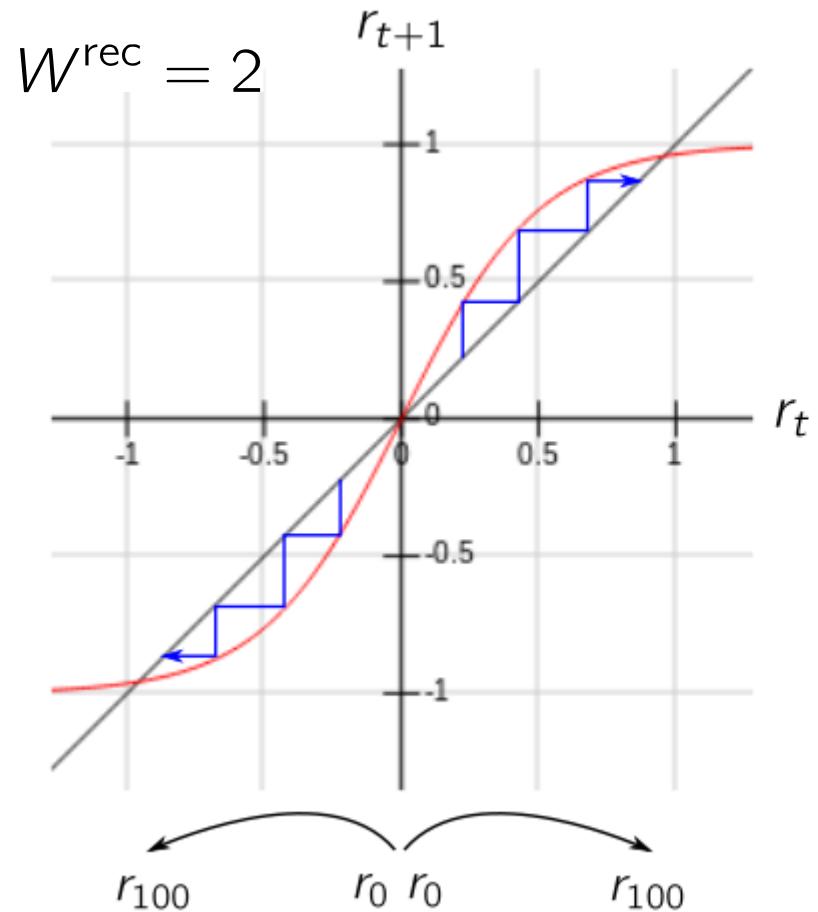
Note: eigenvalues view is only exact for linear models, what about nonlinear models?

Optimizing RNNs: Dynamical System View

Most RNNs are nonlinear, let's take a different view at the problem:

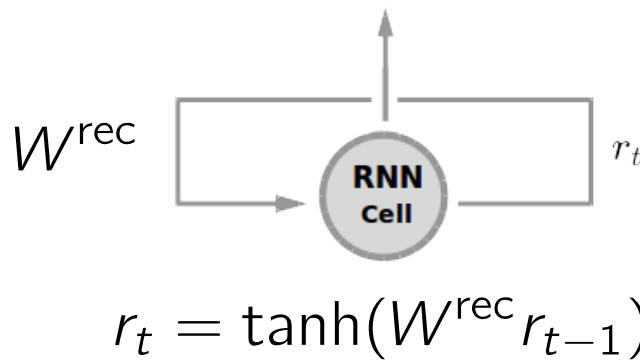


$$\frac{\partial r_{100}}{\partial r_0} = \lim_{r'_0 \rightarrow r_0} \frac{r'_{100} - r_{100}}{r'_0 - r_0}$$

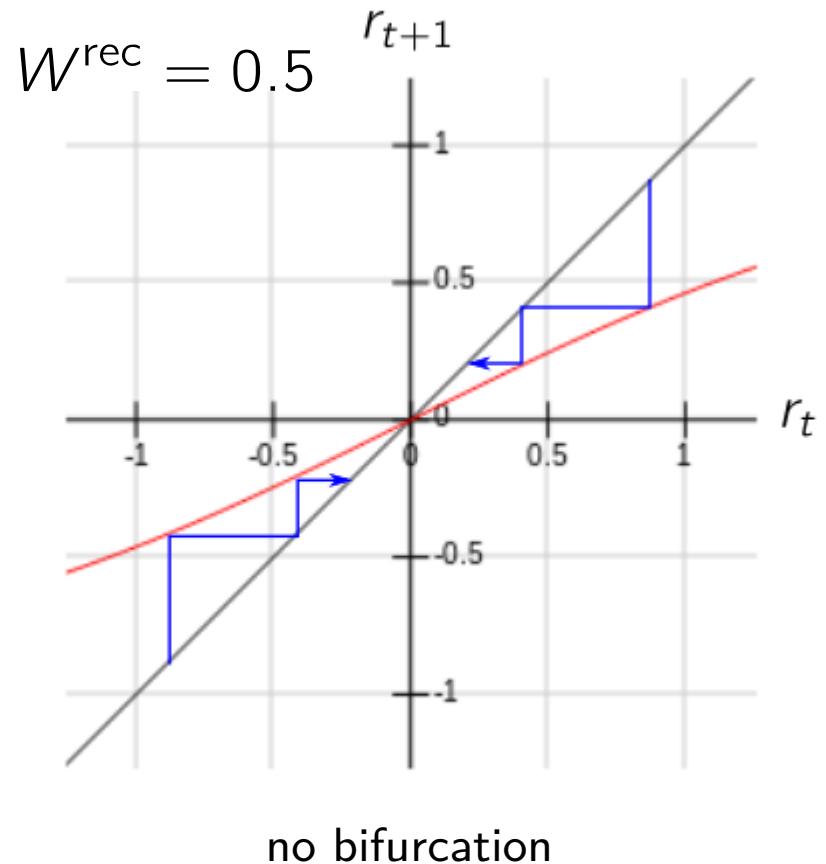


Optimizing RNNs: Dynamical System View

Small weights will reduce the number of bifurcations

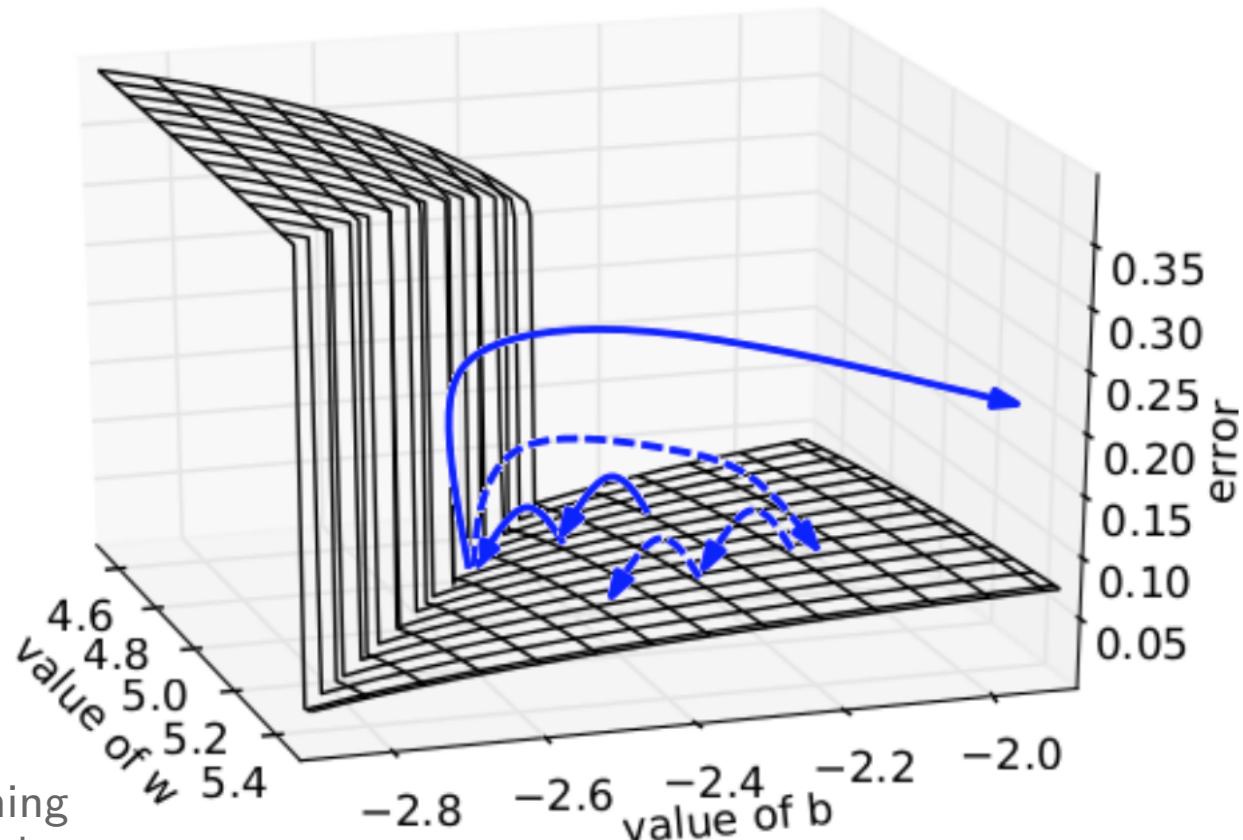


$$\frac{\partial r_{100}}{\partial r_0} = \lim_{r'_0 \rightarrow r_0} \frac{r'_{100} - r_{100}}{r'_0 - r_0}$$



Optimizing RNNs: Dynamical System View

Bifurcations create cliffs in the error function. As a result, gradient descent will not work.



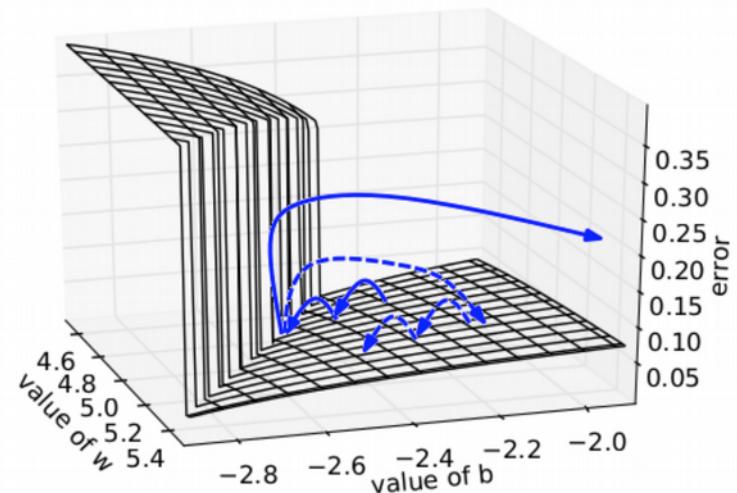
Source: Pascanu '13

On the difficulty of training
Recurrent Neural Networks

Optimizing RNNs: Dynamical System View

Trick: gradient clipping
(Mikolov'12, Pascanu'13):

```
 $\epsilon \leftarrow threshold ;$ 
 $\hat{g} \leftarrow \nabla_{\theta} J;$ 
if  $\|\hat{g}\| \geq threshold$  then
|  $\hat{g} \leftarrow \frac{\epsilon}{\|\hat{g}\|} \hat{g};$ 
end
```



Ignore the magnitude of the gradient, just look at the direction.

Related idea: RProp (Riedmiller'92)

The Difficulty of Training RNNs

RNNs have an unstable gradient due to the recurrent nature of the model. Some tricks exists, but they do not solve the problem completely.

In the following, we consider two alternative approaches:

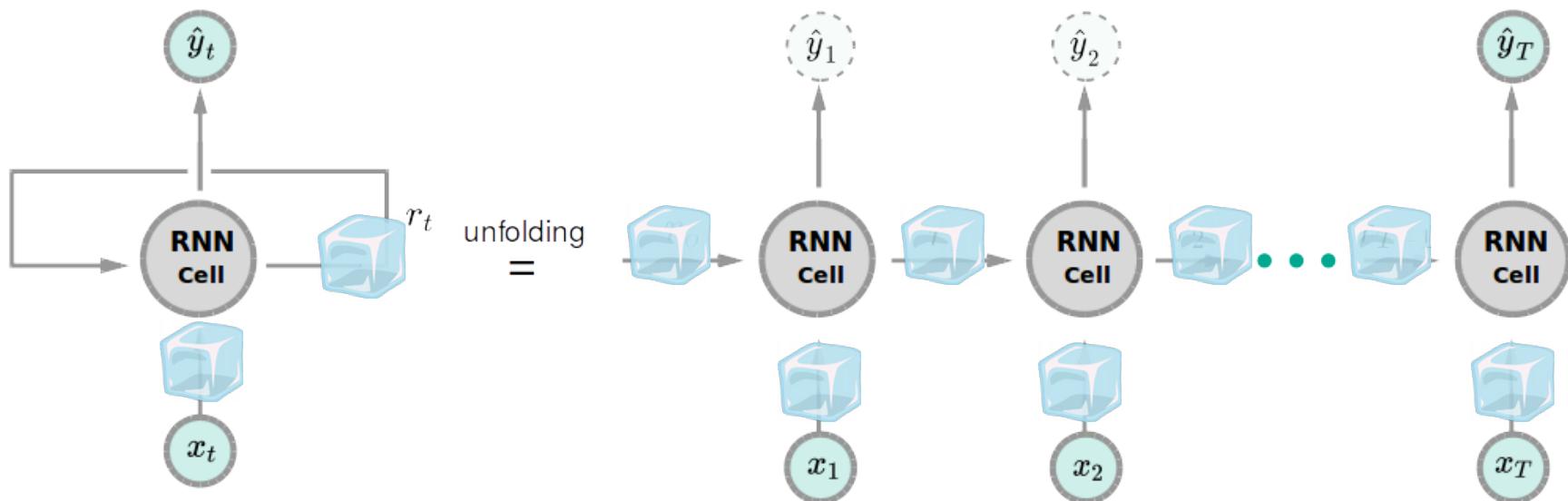
1. **Echo state networks** (only train the top layer)
2. **LSTM** (better RNN architecture)

Echo State Networks (ESNs) [Jaeger'01]

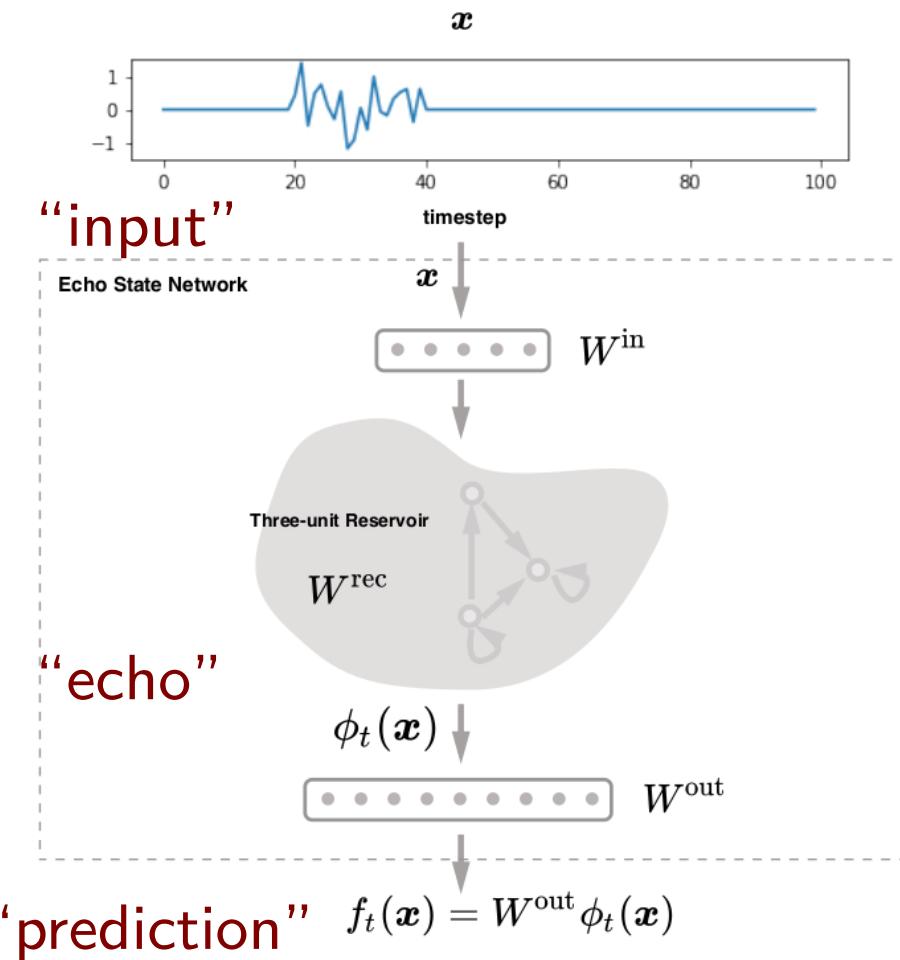
Special type of RNN where recurrent weights are initialized at random and kept fixed throughout training.

Only the output weights W^{out} are adapted.

Adapting the output weights is a simple convex problem, and can be easily optimized.



Echo State Networks (ESNs)



The name “echo state network” reflects the intuition that the network performs a time-dependent expansion (“echo”) of the input time series.

Ideally, the echo representation has produced features that are more suitable for predicting the input time series (e.g. its class).

Echo State Networks (ESNs)

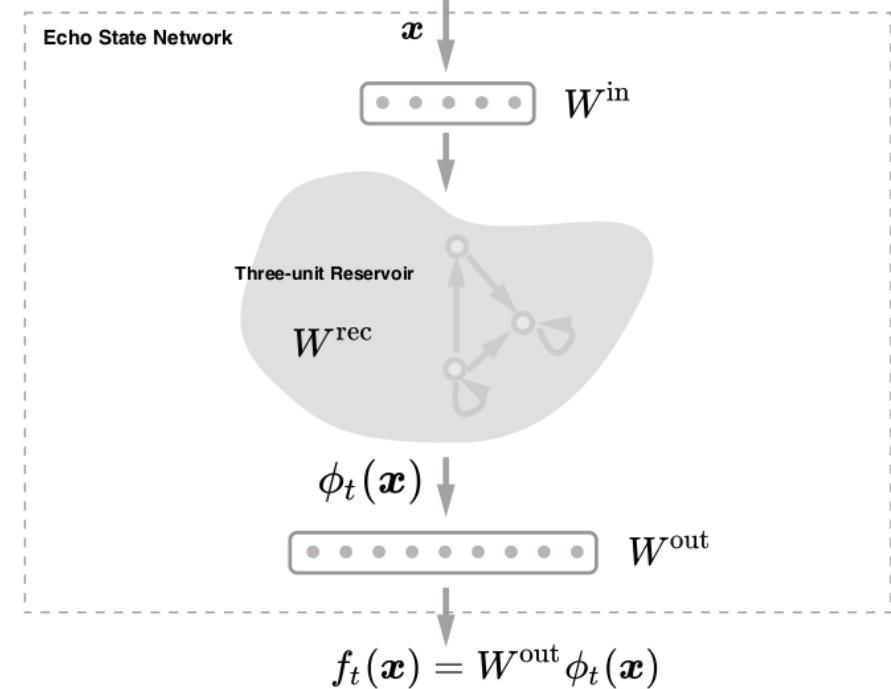
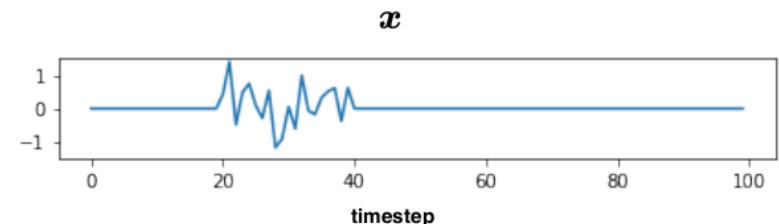
Equations:

$$\tilde{\mathbf{h}}_t = \tanh(W^{\text{in}} \mathbf{x}_t + W^{\text{rec}} \mathbf{h}_{t-1}),$$

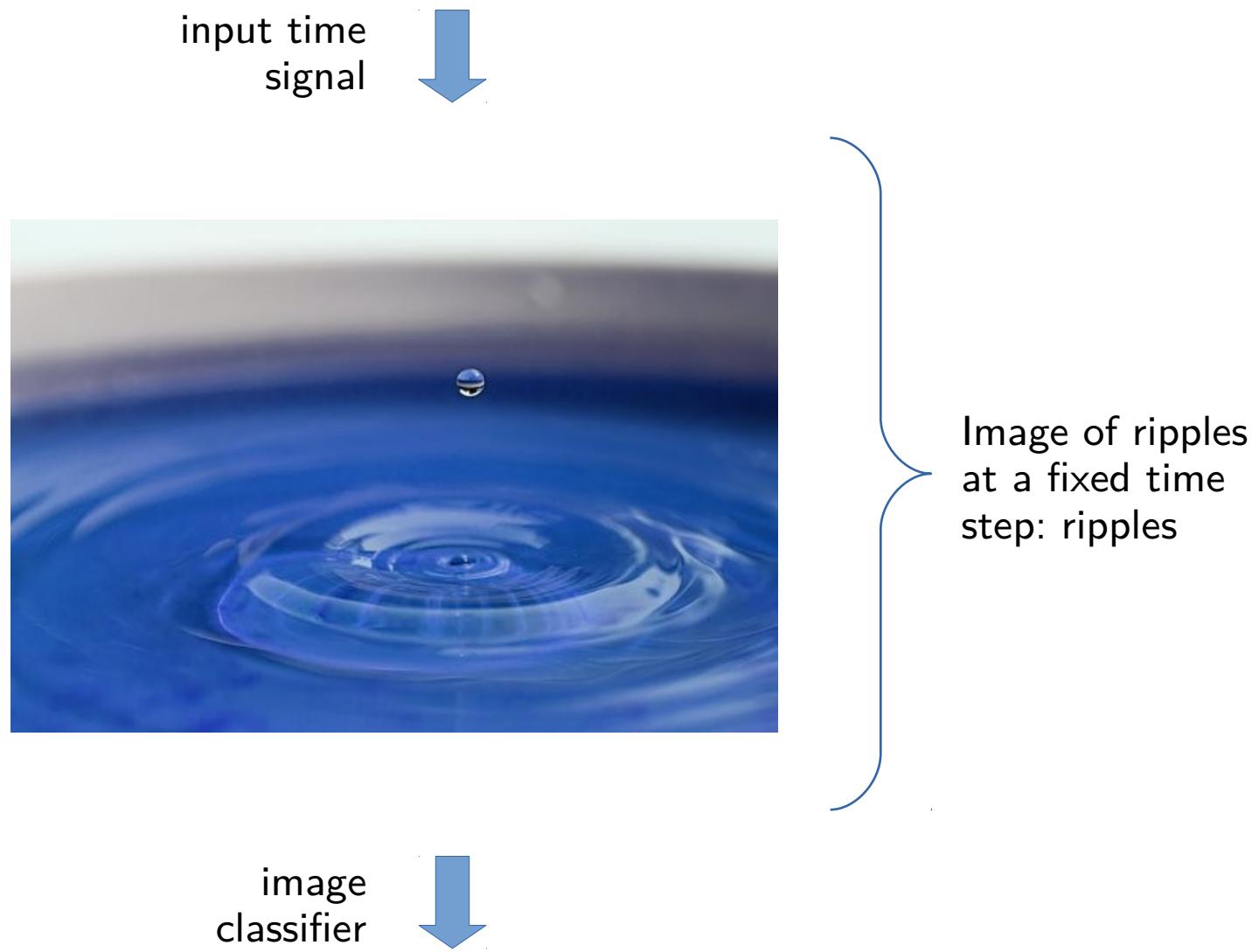
$$\mathbf{h}_t = (1 - \alpha) \mathbf{h}_{t-1} + \alpha \tilde{\mathbf{h}}_t,$$

$$\hat{\mathbf{y}}_t = W^{\text{out}} \text{concat}(\mathbf{x}_t, \mathbf{h}_t),$$

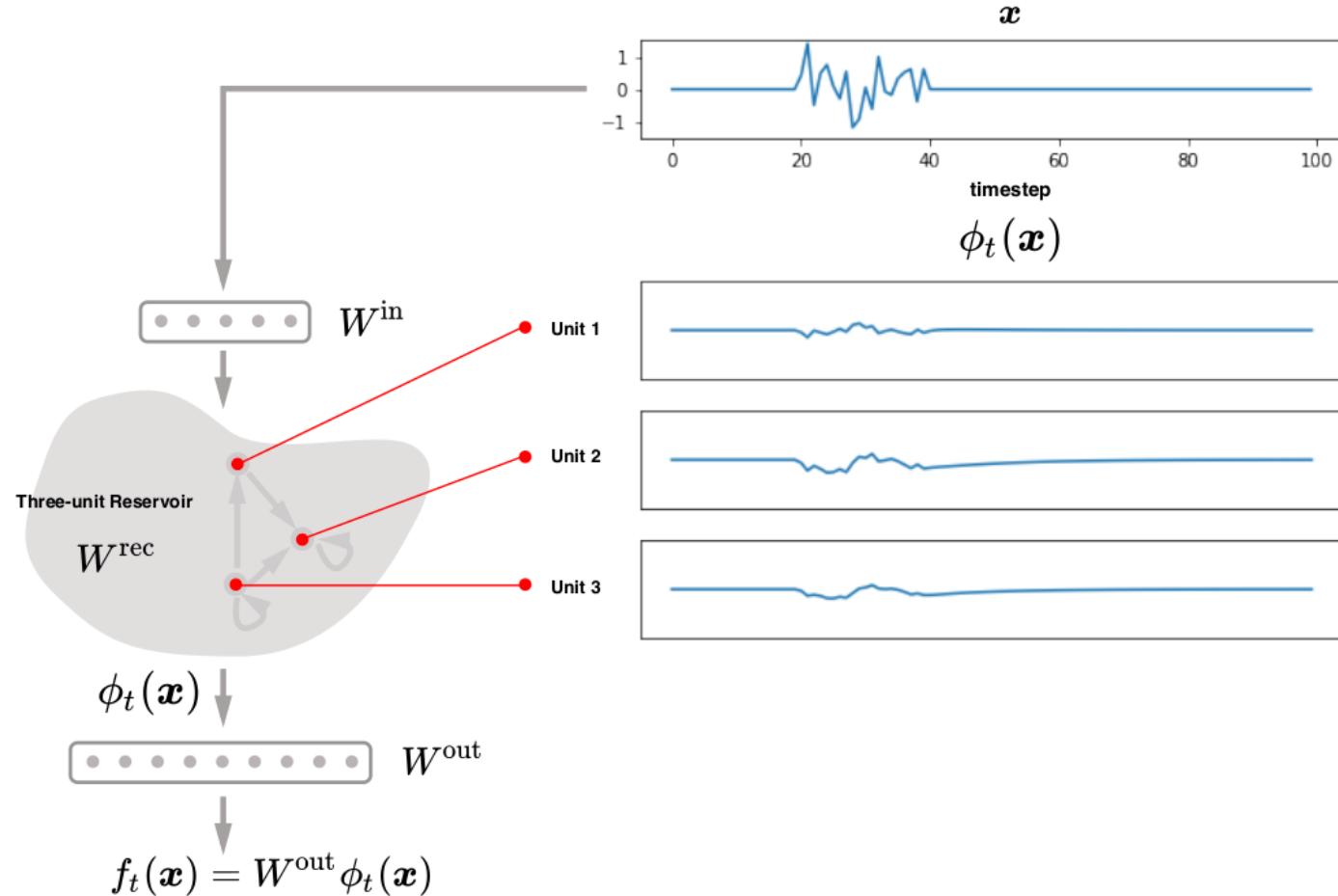
where $W^{\text{rec}} \in \mathbb{R}^{M \times M}$ with M is the number of reservoir units, and $\alpha \in (0, 1]$ is the leaking rate.



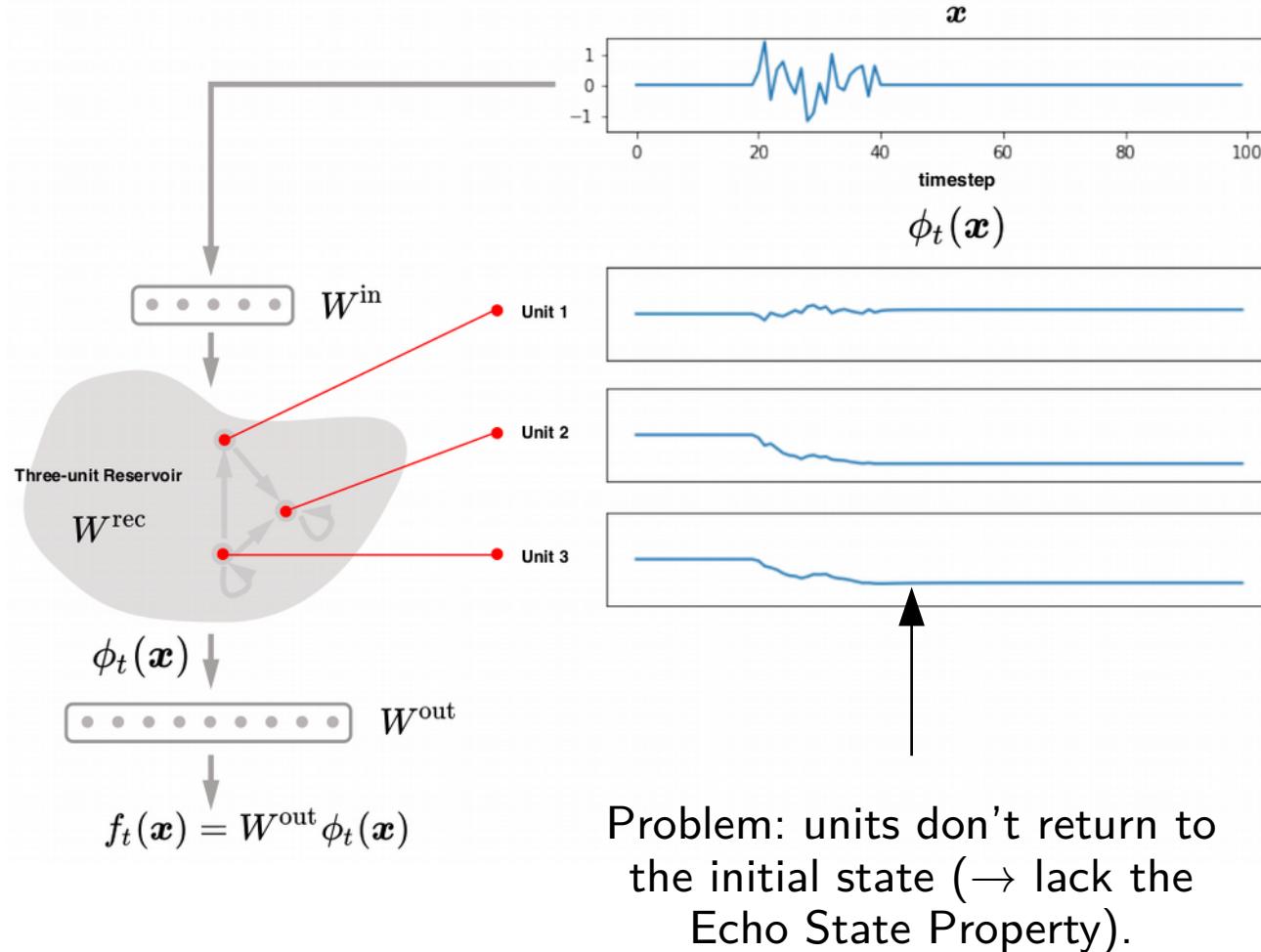
ESN: A Physical Analogy



ESN at Work: Example 1



ESN at Work: Example 2



Echo State Property (ESP)

Goal of ESP: make sure that the solution is not dependent on the initial conditions.

Conditions for ESP: For linear ESN, ESP is guaranteed when highest eigenvalue of the recurrent weights matrix (spectral radius) is smaller than 1. For nonlinear ESN, there is no formal guarantee, but the condition still works empirically.

Analogy to Kernels: ESP ensures the feature map is L2-integrable, and therefore enables the construction of a kernel:

$$k_{\text{ESN}}(x, x') = \langle \phi_{\text{ESN}}(x), \phi_{\text{ESN}}(x') \rangle$$

Insight: An echo state network is a kernel (a smart one).

ESNs: How to Choose the Parameters?

Equations:

Set it large to create
abrupt nonlinearities

Set it such that spectral
radius is <1.

$$\tilde{\mathbf{h}}_t = \tanh(W^{\text{in}} \mathbf{x}_t + W^{\text{rec}} \mathbf{h}_{t-1}),$$

$$\mathbf{h}_t = (1 - \alpha) \mathbf{h}_{t-1} + \alpha \tilde{\mathbf{h}}_t,$$

$$\hat{\mathbf{y}}_t = W^{\text{out}} \text{concat}(\mathbf{x}_t, \mathbf{h}_t),$$

where $W^{\text{rec}} \in \mathbb{R}^{M \times M}$ with M is the
number of reservoir units, and

$\alpha \in (0, 1]$ is the leaking rate.

As large as you can afford
(typically $M > 10000$).

Set it close to zero if the
time series evolves slowly.



Cf. tutorial paper: M. Lukosevicius: A Practical Guide to Applying Echo State Networks. Neural Networks: Tricks of the Trade (2nd ed.) 2012: 659-686

Advantages and Disadvantages of ESNs

Advantages:

Easy to train (convex)

Work well for low-dimensional input

Embedded in the kernel framework

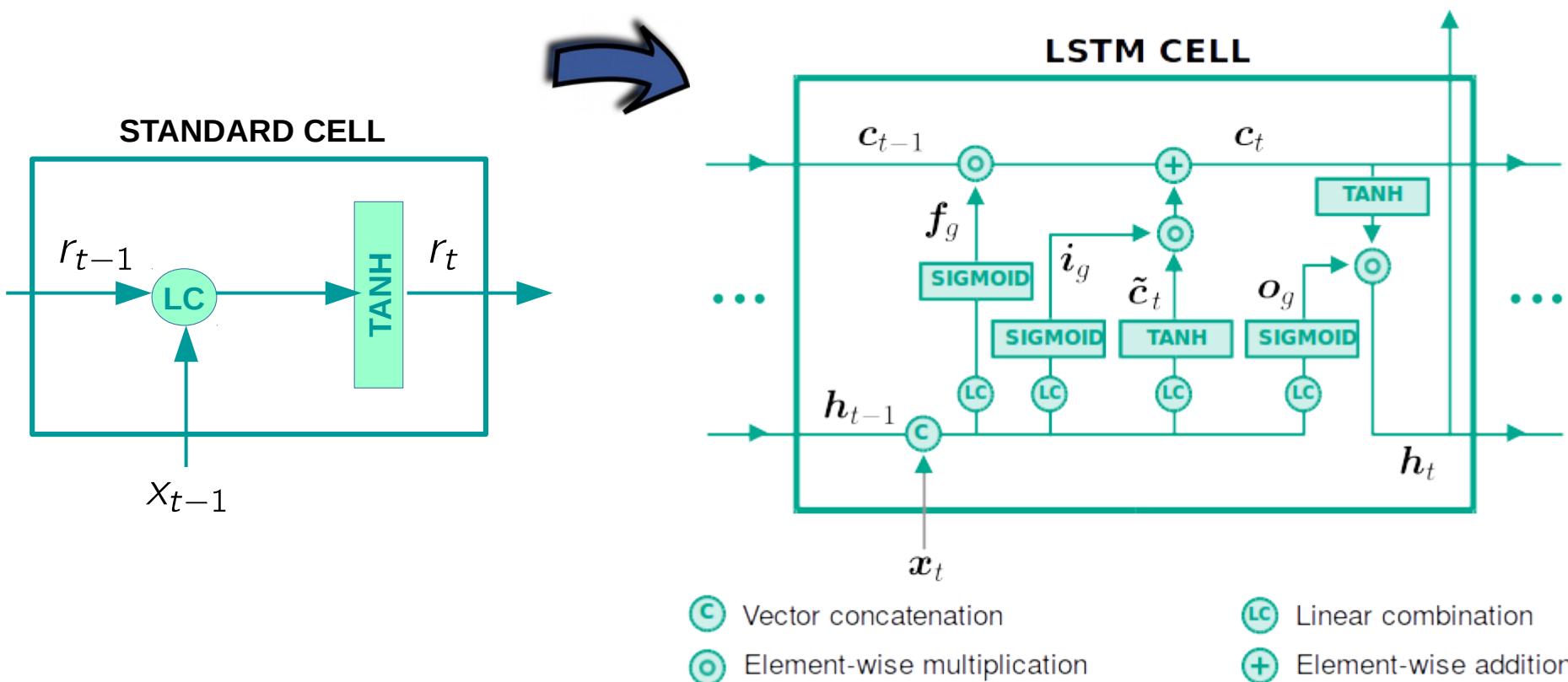
Disadvantages:

Can be tricky to design a reservoir with good properties

Require a big reservoir to solve complex problems

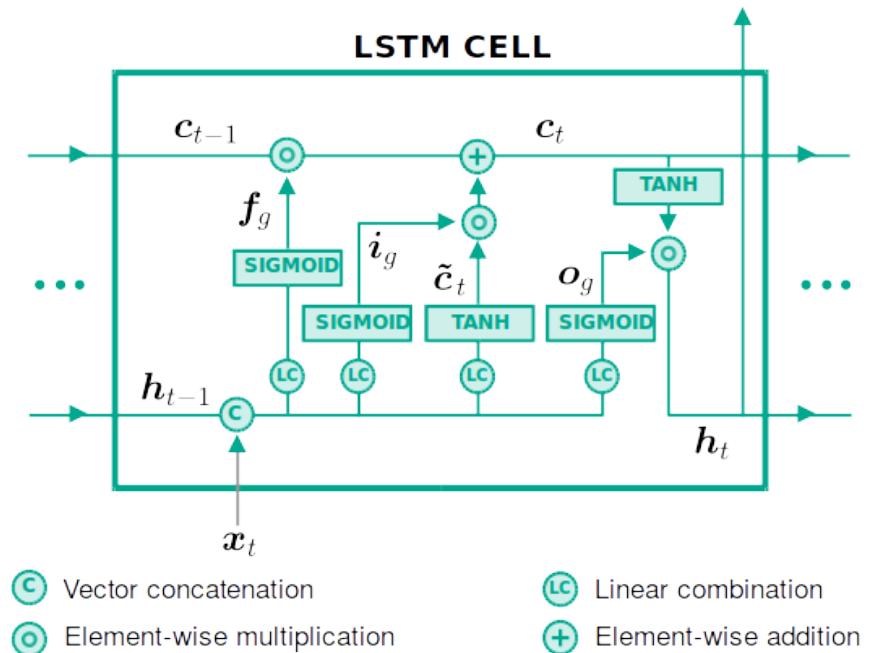
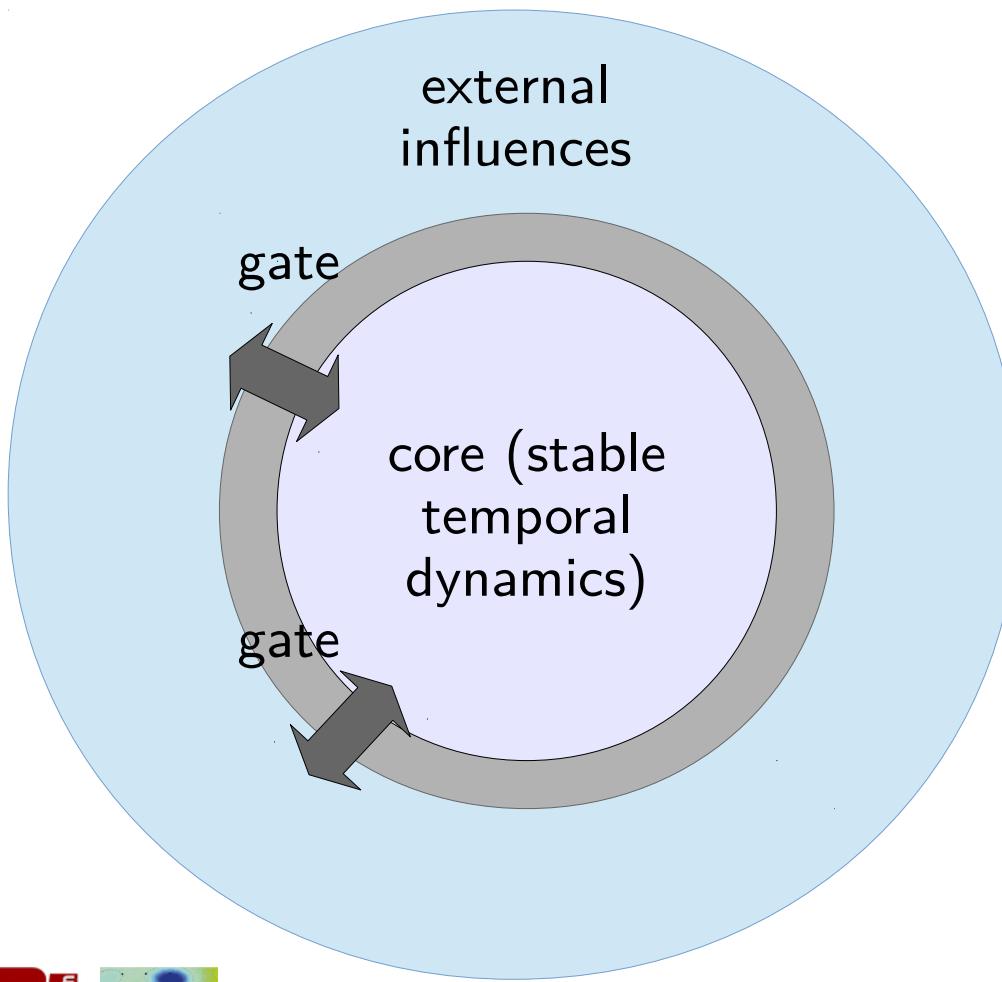
LSTMs [Hochreiter&Schmidhuber'97]

LSTM = RNN with long-short term memory



LSTMs

Key idea: protect the internal dynamics (core) from external influences (input and output feedback) through gates.



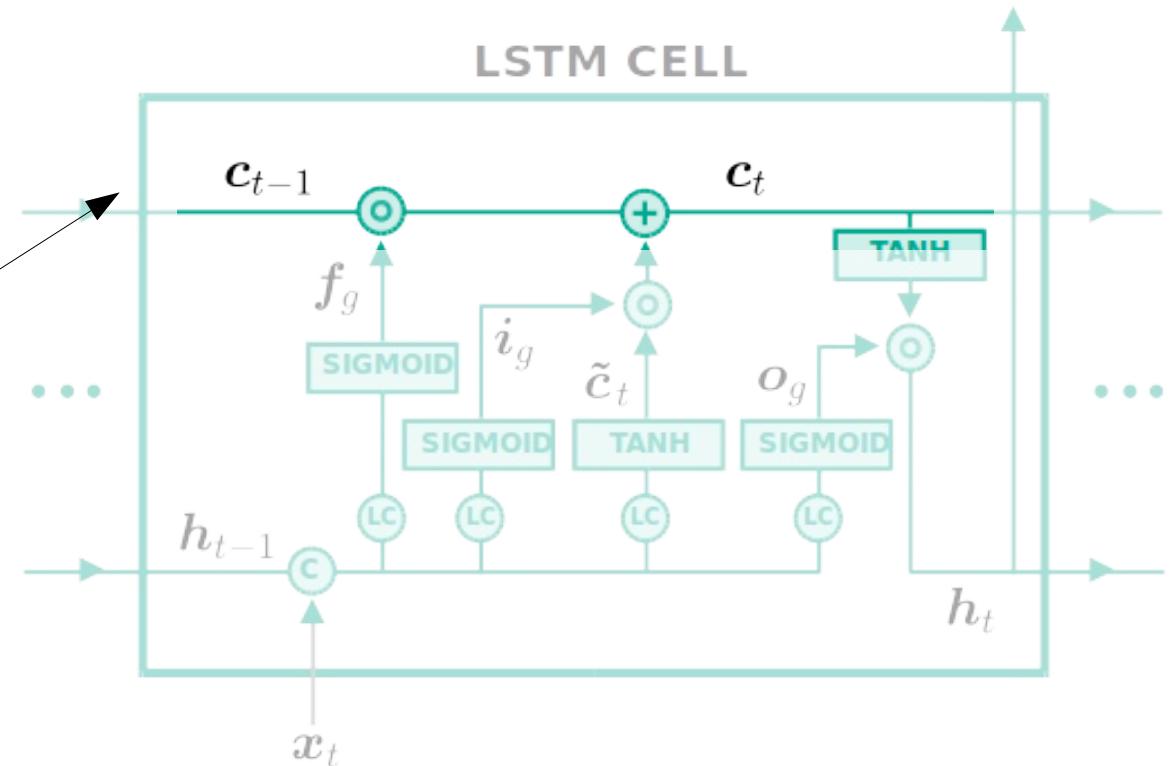
LSTMs

Key idea: protect the internal dynamics (core) from external influences (input and output feedback) through gates.

Internal dynamics

Stays stable over time
(only small increments).

No weight matrices and
squashing nonlinearities
on that path.



(C) Vector concatenation

(○) Element-wise multiplication

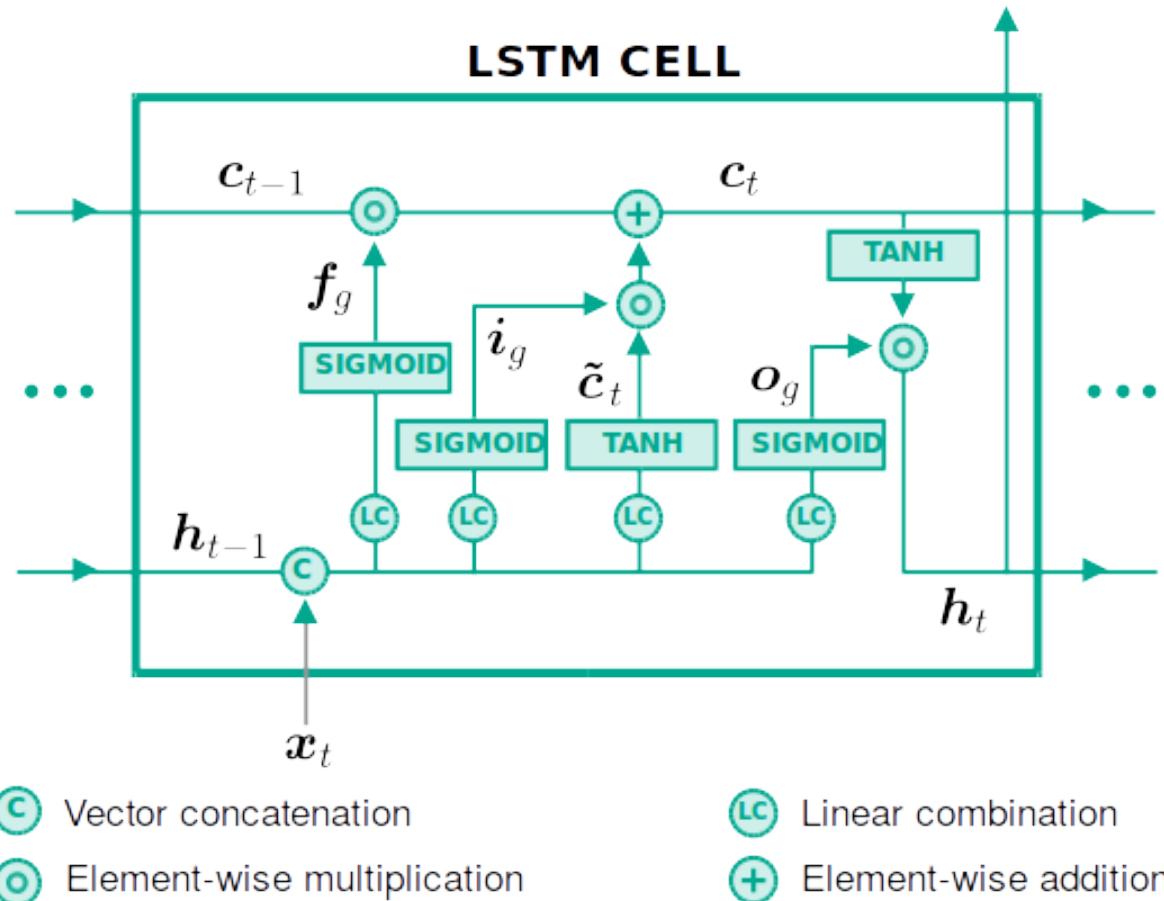
(LC) Linear combination

(+) Element-wise addition

LSTMs

Key idea: protect the internal dynamics (core) from external influences (input and output feedback) through gates.

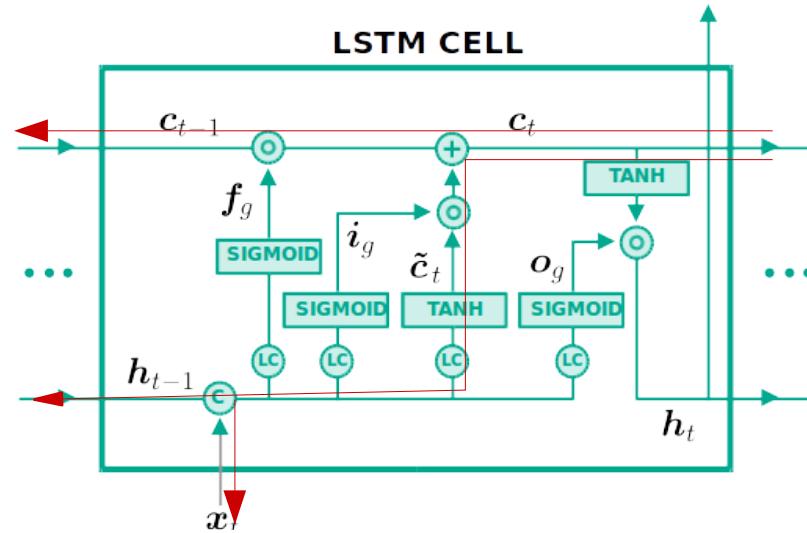
- Employ 3 gates:
 - Input gate i_g
 - Forget gate f_g
 - Output gate o_g



Explaining LSTMs

LSTM predicts well, but can we explain its predictions?

Yes (cf. Arras'17, Rieger & Chormai'18), special rules for propagating backward in that cell.



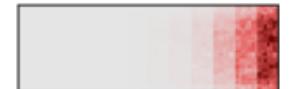
Explaining
majority:



Explaining negative sentiment in movie reviews:

1. do n't waste your money .
2. neither funny nor suspenseful nor particularly well-drawn .
3. it 's not horrible , just horribly mediocre .
4. ... too slow , too boring , and occasionally annoying .
5. it 's neither as romantic nor as thrilling as it should be .
6. the master of disaster - it 's a piece of dreck disguised as cc

Shallow



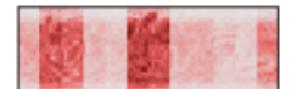
Deep



ConvDeep



R-LSTM



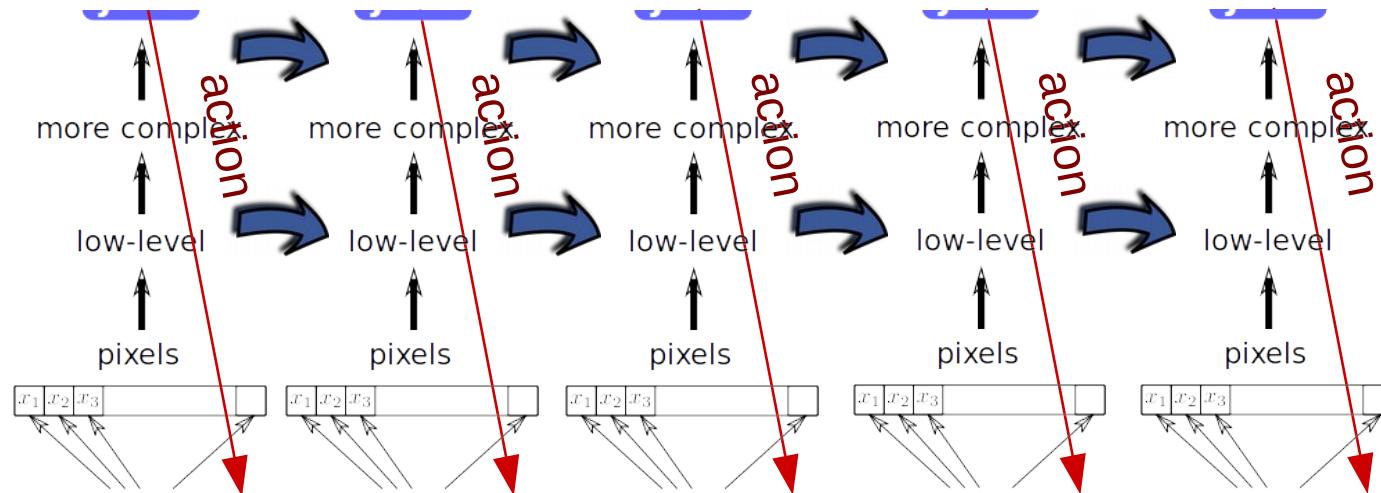
ConvR-LSTM



From Perception to Action

So far, we have built a state representation of the environment.

What about taking actions in an environment?



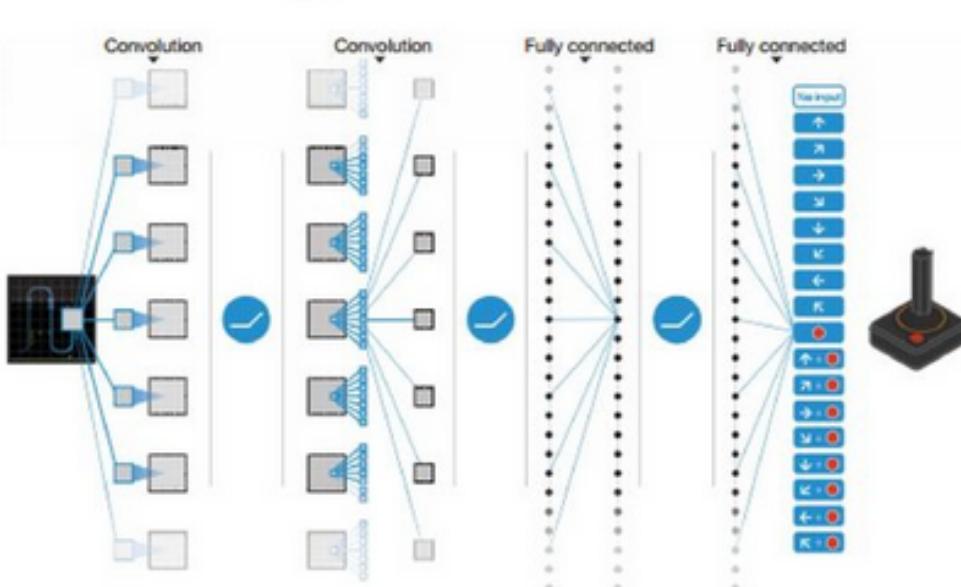
Neural Fitted Q Iteration

$$y_i = \mathbb{E}_{s' \sim \mathcal{E}} [r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) | s, a]$$

better Q values

$$L_i(\theta_i) = \mathbb{E}_{s, a \sim \rho(\cdot)} \left[(y_i - Q(s, a; \theta_i))^2 \right]$$

Deep Q Network



Source: Mnih'15
Human-level control through deep reinforcement learning

Summary

A large number of input data have a time component. Special models (e.g. recurrent neural networks) are needed for this.

Standard recurrent neural networks are hard to train (high curvature within and between layers, discontinuities in the error function).

Two approaches have been presented to make RNNs more applicable:

1. **Echo state networks** (train only the output layer)
2. **LSTMs** (adapt the structure to decouple the internal dynamics from the external influences)

RNNs is only one aspect. Reinforcement learning brings control capabilities. RNNs and RL work well in combination.