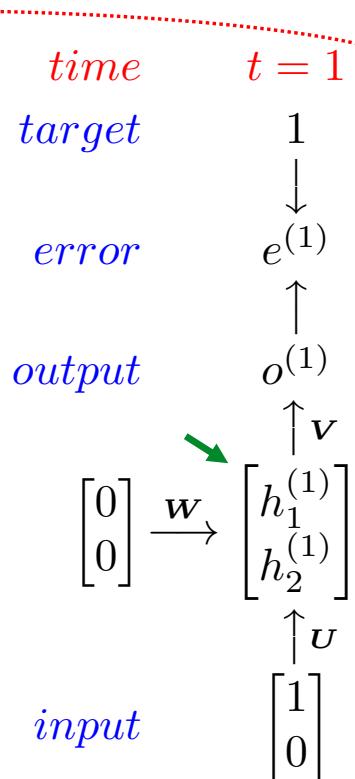
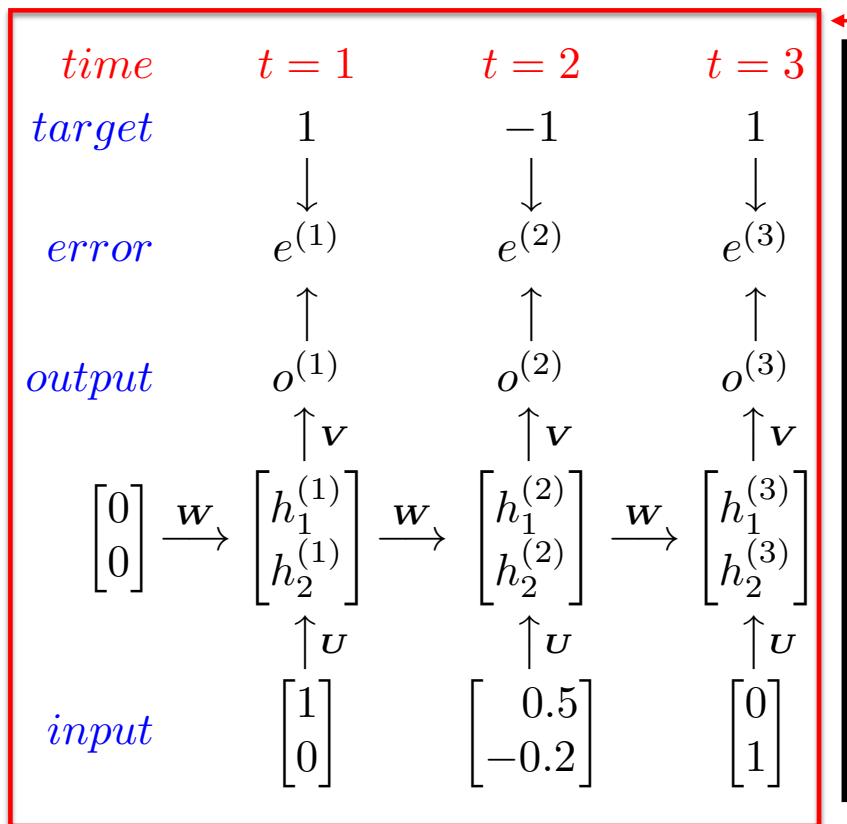


Real-Time Recurrent Learning (RTRL)



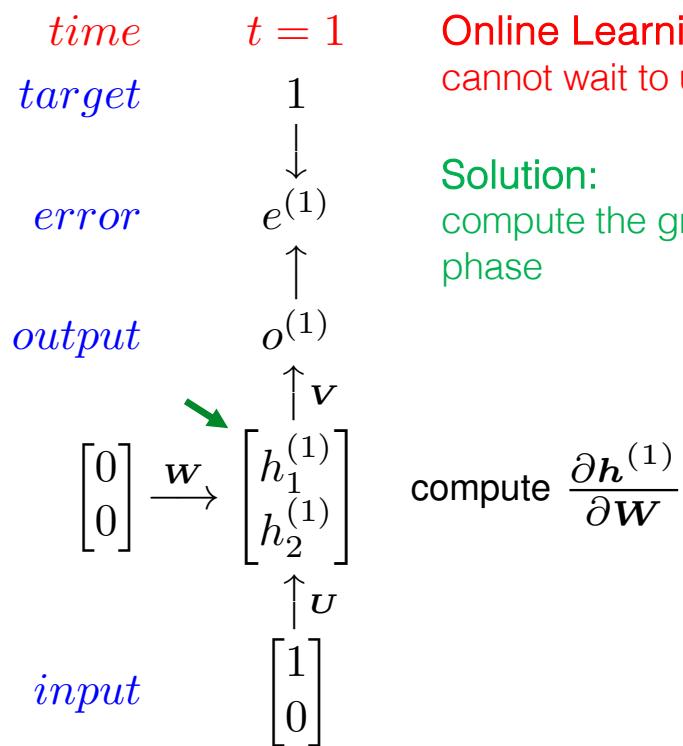
Online Learning:
cannot wait to unroll the recursive network

Solution:
compute the gradients during the forward phase

Real-Time Recurrent Learning (RTRL)

Equations

$$\begin{aligned}
 h^{(0)} &= \mathbf{0}, \\
 a_h^{(t)} &= \mathbf{U}x^{(t)} + \mathbf{W}h^{(t-1)} + \mathbf{b}, \\
 h^{(t)} &= \tanh(a_h^{(t)}), \\
 a_o^{(t)} &= \mathbf{V}h^{(t)} + c, \\
 o^{(t)} &= \tanh(a_o^{(t)}), \\
 L &= \sum_t e^{(t)} = \sum_t (y^{(t)} - o^{(t)})^2
 \end{aligned}$$



Online Learning:
cannot wait to unroll the recursive network

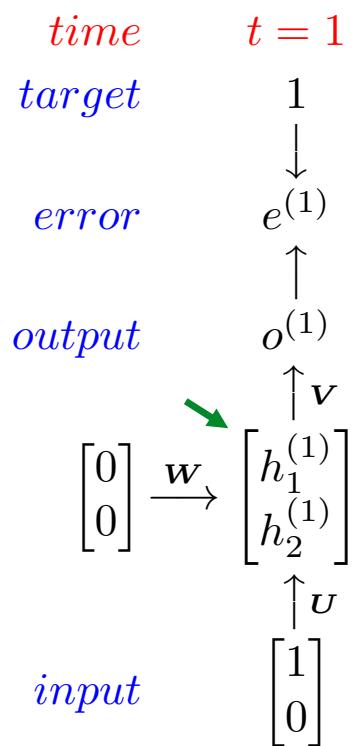
Solution:
compute the gradients during the forward phase

compute $\frac{\partial h^{(1)}}{\partial W}$

Real-Time Recurrent Learning (RTRL)

Equations

$$\begin{aligned} \mathbf{h}^{(0)} &= \mathbf{0}, \\ \mathbf{a}_h^{(t)} &= \mathbf{U}\mathbf{x}^{(t)} + \mathbf{W}\mathbf{h}^{(t-1)} + \mathbf{b}, \\ \mathbf{h}^{(t)} &= \tanh(\mathbf{a}_h^{(t)}), \\ \mathbf{a}_o^{(t)} &= \mathbf{V}\mathbf{h}^{(t)} + \mathbf{c}, \\ o^{(t)} &= \tanh(a_o^{(t)}), \\ L &= \sum_t e^{(t)} = \sum_t (y^{(t)} - o^{(t)})^2 \end{aligned}$$



Online Learning:

cannot wait to unroll the recursive network

Solution:

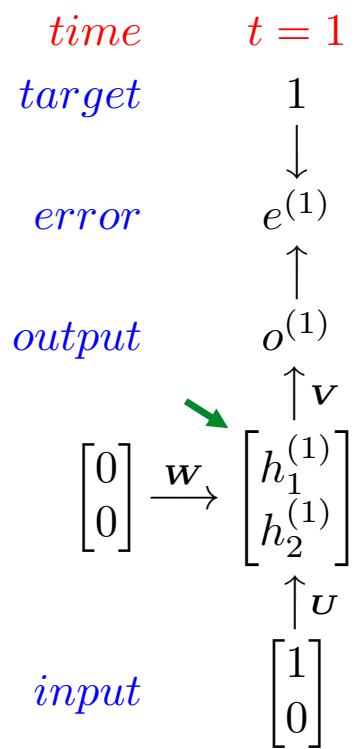
compute the gradients during the forward phase

here it is also possible to compute $\frac{\partial \mathbf{h}^{(1)}}{\partial \mathbf{U}}$

Real-Time Recurrent Learning (RTRL)

Equations

$$\begin{aligned}
 h^{(0)} &= \mathbf{0}, \\
 a_h^{(t)} &= \mathbf{U}x^{(t)} + \mathbf{W}h^{(t-1)} + \mathbf{b}, \\
 h^{(t)} &= \tanh(a_h^{(t)}), \\
 a_o^{(t)} &= \mathbf{V}h^{(t)} + c, \\
 o^{(t)} &= \tanh(a_o^{(t)}), \\
 L &= \sum_t e^{(t)} = \sum_t (y^{(t)} - o^{(t)})^2
 \end{aligned}$$



Online Learning:

cannot wait to unroll the recursive network

Solution:

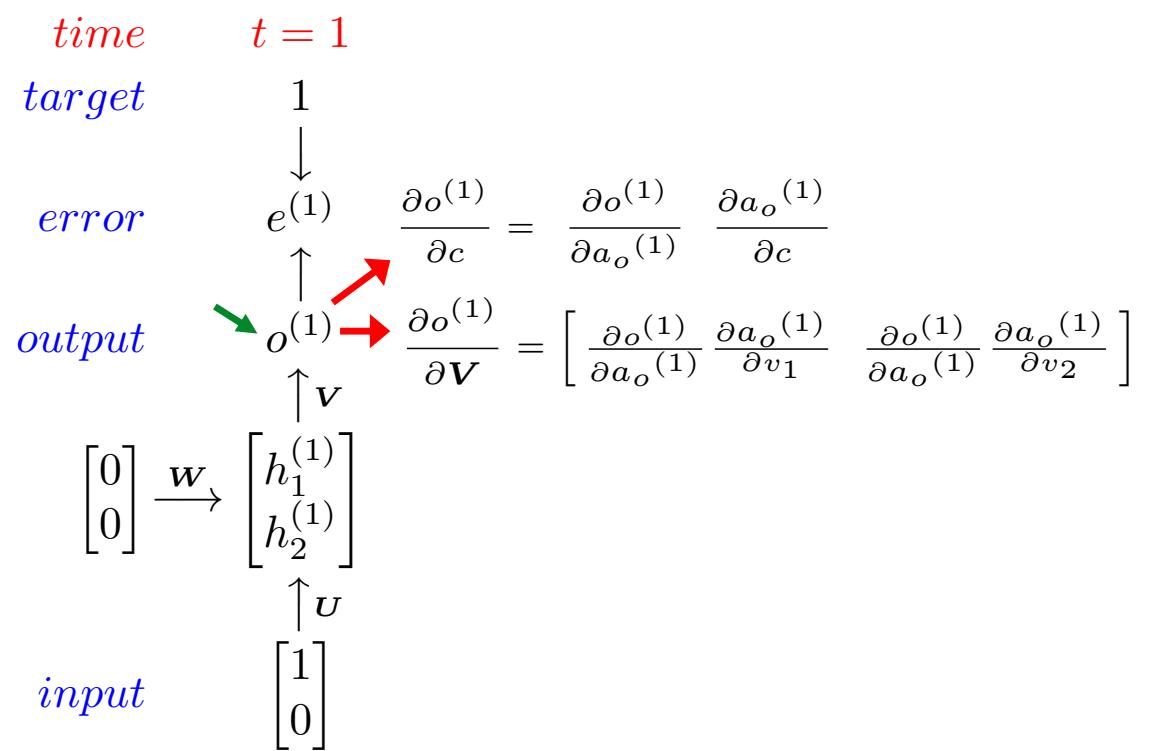
compute the gradients during the forward phase

... and $\frac{\partial \mathbf{h}^{(1)}}{\partial \mathbf{b}}$, which is a tensor in $\mathbb{R}^{2 \times 2}$

Real-Time Recurrent Learning (RTRL)

Equations

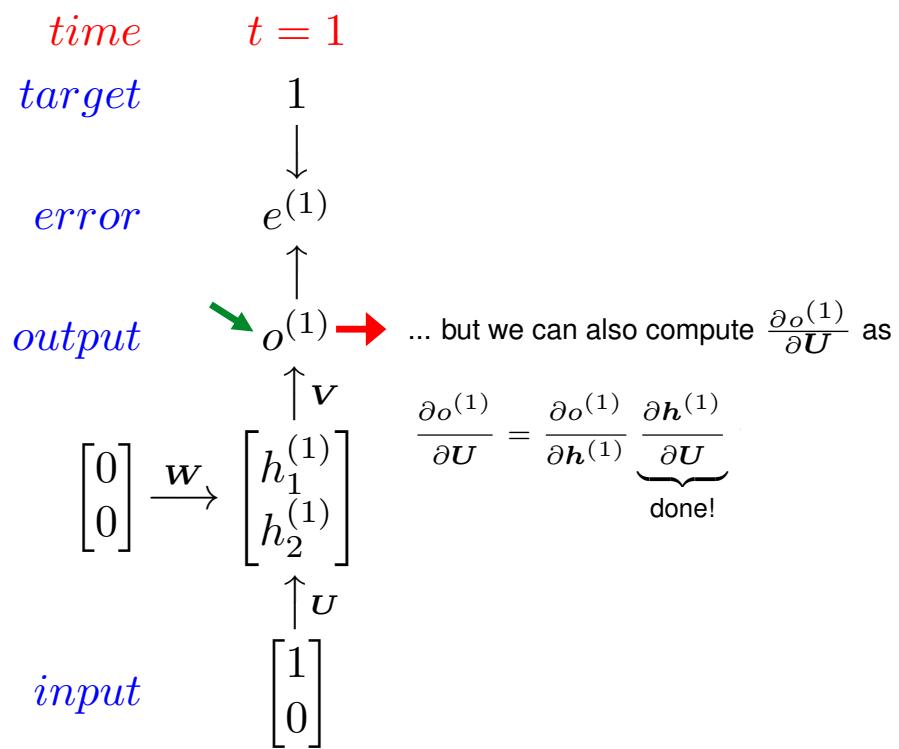
$$\begin{aligned}
 \mathbf{h}^{(0)} &= \mathbf{0}, \\
 \mathbf{a}_h^{(t)} &= \mathbf{U}\mathbf{x}^{(t)} + \mathbf{W}\mathbf{h}^{(t-1)} + \mathbf{b}, \\
 \mathbf{h}^{(t)} &= \tanh(\mathbf{a}_h^{(t)}), \\
 \mathbf{a}_o^{(t)} &= \mathbf{V}\mathbf{h}^{(t)} + \mathbf{c}, \\
 \mathbf{o}^{(t)} &= \tanh(\mathbf{a}_o^{(t)}), \\
 L &= \sum_t e^{(t)} = \sum_t (y^{(t)} - o^{(t)})^2
 \end{aligned}$$



Real-Time Recurrent Learning (RTRL)

Equations

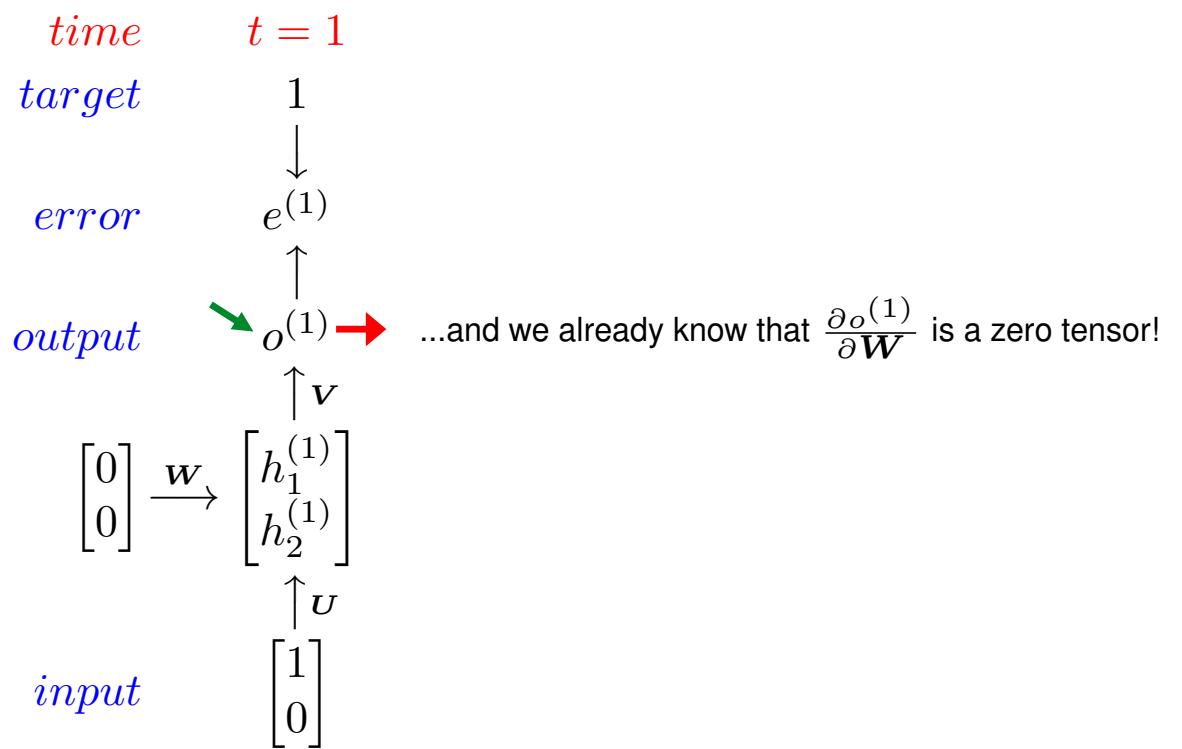
$$\begin{aligned}
 h^{(0)} &= \mathbf{0}, \\
 a_h^{(t)} &= \mathbf{U}x^{(t)} + \mathbf{W}h^{(t-1)} + \mathbf{b}, \\
 h^{(t)} &= \tanh(a_h^{(t)}), \\
 a_o^{(t)} &= \mathbf{V}h^{(t)} + c, \\
 o^{(t)} &= \tanh(a_o^{(t)}), \\
 L &= \sum_t e^{(t)} = \sum_t (y^{(t)} - o^{(t)})^2
 \end{aligned}$$



Real-Time Recurrent Learning (RTRL)

Equations

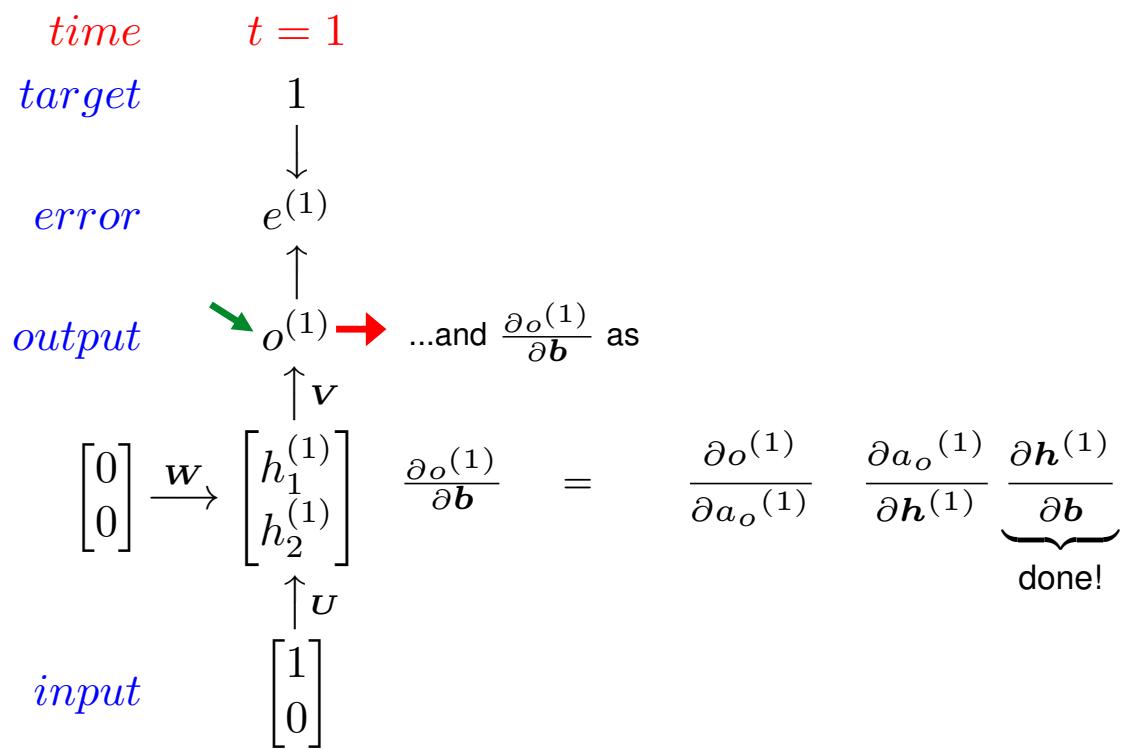
$$\begin{aligned} \mathbf{h}^{(0)} &= \mathbf{0}, \\ \mathbf{a}_h^{(t)} &= \mathbf{U}\mathbf{x}^{(t)} + \mathbf{W}\mathbf{h}^{(t-1)} + \mathbf{b}, \\ \mathbf{h}^{(t)} &= \tanh(\mathbf{a}_h^{(t)}), \\ \mathbf{a}_o^{(t)} &= \mathbf{V}\mathbf{h}^{(t)} + \mathbf{c}, \\ \mathbf{o}^{(t)} &= \tanh(\mathbf{a}_o^{(t)}), \\ L &= \sum_t e^{(t)} = \sum_t (y^{(t)} - o^{(t)})^2 \end{aligned}$$



Real-Time Recurrent Learning (RTRL)

Equations

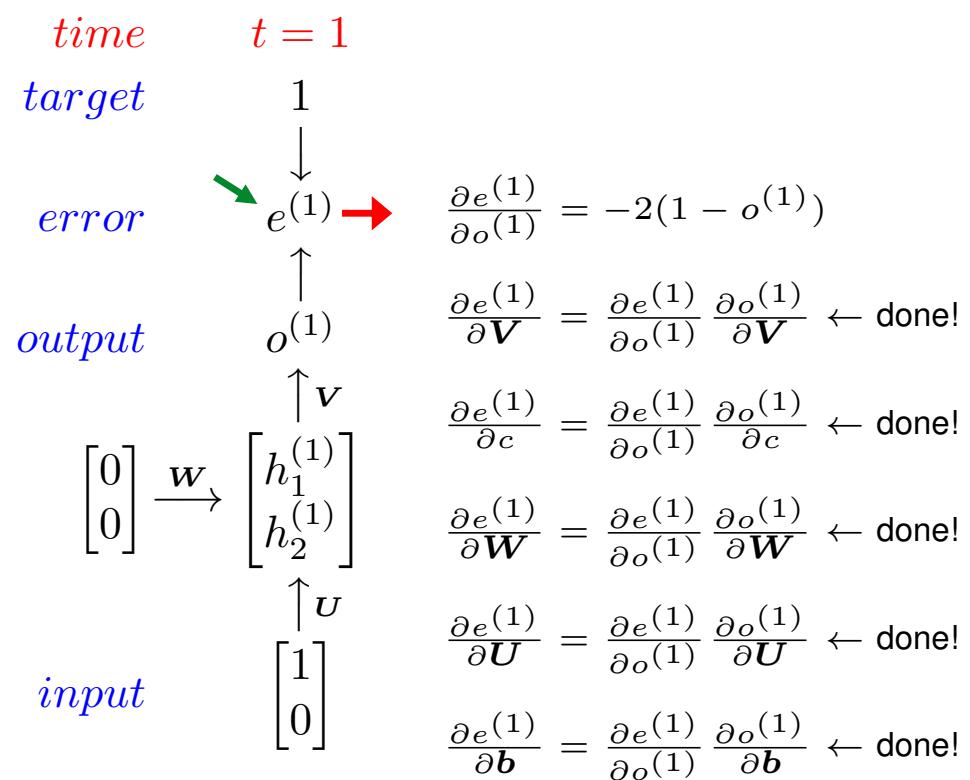
$$\begin{aligned}
 h^{(0)} &= 0, \\
 a_h^{(t)} &= \mathbf{U}\mathbf{x}^{(t)} + \mathbf{W}\mathbf{h}^{(t-1)} + \mathbf{b}, \\
 h^{(t)} &= \tanh(a_h^{(t)}), \\
 a_o^{(t)} &= \mathbf{V}\mathbf{h}^{(t)} + c, \\
 o^{(t)} &= \tanh(a_o^{(t)}), \\
 L &= \sum_t e^{(t)} = \sum_t (y^{(t)} - o^{(t)})^2
 \end{aligned}$$



Real-Time Recurrent Learning (RTRL)

Equations

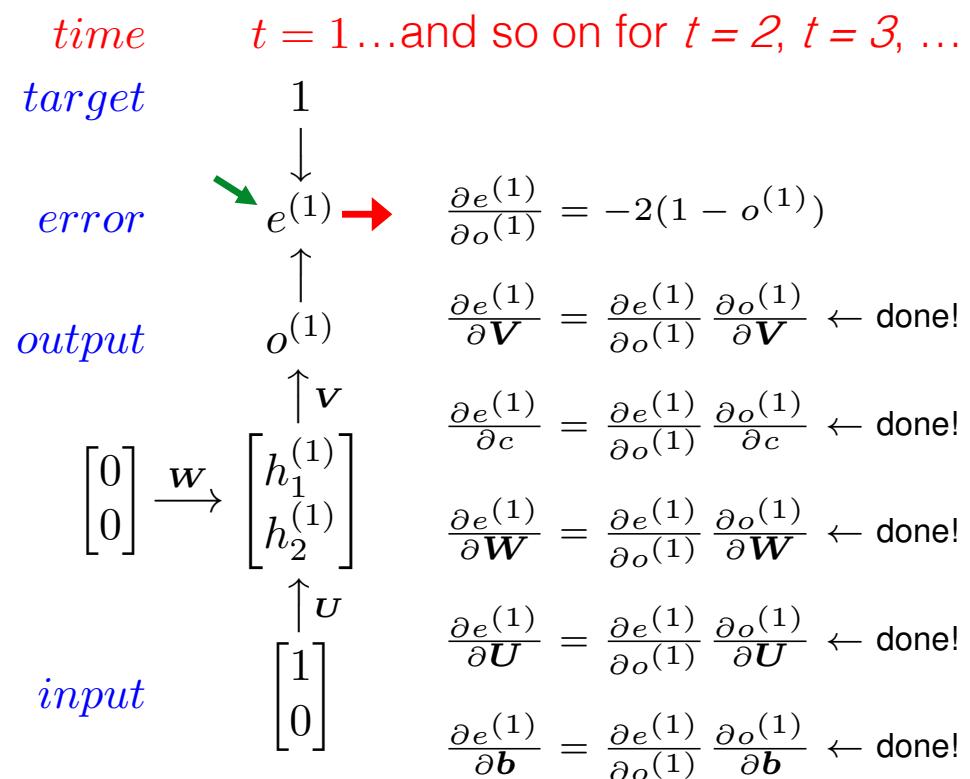
$$\begin{aligned}
 h^{(0)} &= \mathbf{0}, \\
 a_h^{(t)} &= \mathbf{U}x^{(t)} + \mathbf{W}h^{(t-1)} + \mathbf{b}, \\
 h^{(t)} &= \tanh(a_h^{(t)}), \\
 a_o^{(t)} &= \mathbf{V}h^{(t)} + c, \\
 o^{(t)} &= \tanh(a_o^{(t)}), \\
 L &= \sum_t e^{(t)} = \sum_t (y^{(t)} - o^{(t)})^2
 \end{aligned}$$



Real-Time Recurrent Learning (RTRL)

Equations

$$\begin{aligned}
 h^{(0)} &= 0, \\
 a_h^{(t)} &= \mathbf{U}\mathbf{x}^{(t)} + \mathbf{W}h^{(t-1)} + \mathbf{b}, \\
 h^{(t)} &= \tanh(a_h^{(t)}), \\
 a_o^{(t)} &= \mathbf{V}h^{(t)} + c, \\
 o^{(t)} &= \tanh(a_o^{(t)}), \\
 L &= \sum_t e^{(t)} = \sum_t (y^{(t)} - o^{(t)})^2
 \end{aligned}$$



BPTT vs RTRL

- Both **BPTT** and **RTRL** compute the same quantities (gradients), but in different ways
- Computational complexity

	<i>Space</i>	<i>Time</i>
<i>BPTT</i>	$O(NT)$	$O(N^2T)$
<i>RTRL</i>	$O(N^3)$	$O(N^4)$

T : time steps

N : number of units

- Whatever the computation is performed, there is a serious problem...**vanishing/exploding of gradient!**

Gradient Vanishing/Exploding Problem how to cope with it (part 1)

Learning Long-term Dependencies is difficult

What is a long-term dependency ?

- when the *desidered* output at time t depends on the input at time $t - \tau$, with $t > \tau \gg 1$

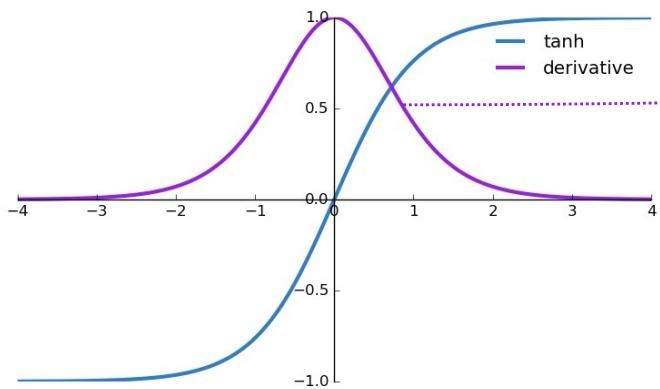
$$\mathbf{x}^{(t-\tau)} \rightarrow \mathbf{y}^{(t)}$$

e.g.

$$\mathbf{x}^{(t-100)} \rightarrow \mathbf{y}^{(t)}$$

- this means that, for the Recurrent Neural Network to output the correct *desidered* $\mathbf{y}^{(t)}$, it has:
 - to recognize* its dependency on $\mathbf{x}^{(t-\tau)}$, and
 - to use* $\mathbf{x}^{(t-\tau)}$ in the generation of $\mathbf{y}^{(t)}$

Learning Long-term Dependencies is difficult



$$t=3 \Rightarrow \frac{\partial h_1^{(3)}}{\partial \mathbf{w}} = \underbrace{(1 - (h_1^{(3)})^2)}_{w_{11}(1 - (h_1^{(2)})^2)} \left(\begin{bmatrix} h_1^{(2)} & h_2^{(2)} \\ 0 & 0 \end{bmatrix} + \underbrace{w_{12}(1 - (h_2^{(2)})^2)}_{w_{12}(1 - (h_2^{(1)})^2)} \begin{bmatrix} 0 & 0 \\ h_1^{(1)} & h_2^{(1)} \end{bmatrix} \right)$$

$$\frac{\partial h_2^{(3)}}{\partial \mathbf{w}} = \underbrace{(1 - (h_2^{(3)})^2)}_{w_{21}(1 - (h_1^{(2)})^2)} \left(\begin{bmatrix} 0 & 0 \\ h_1^{(2)} & h_2^{(2)} \end{bmatrix} + w_{22}(1 - (h_2^{(2)})^2) \begin{bmatrix} 0 & 0 \\ h_1^{(1)} & h_2^{(1)} \end{bmatrix} \right)$$

Learning Long-term Dependencies is difficult

- if $\|W\|$ is large enough, then $\|\frac{\partial h^{(t)}}{\partial h^{(t-\tau)}}\|$ will increase for increasing values of τ (*exploding gradient* \Rightarrow hard to decide the “right” descent step size)
- however, in order to robustly store past information, the dynamics of the network must exhibit attractors:

spectral radius $\rho(W) < 1$, i.e. W is *contractive*

- In the presence of attractors, *gradients vanish* going backward in time:

$$\lim_{t-\tau \rightarrow \infty} \frac{\partial h^{(t)}}{\partial h^{(\tau)}} = 0$$

\Rightarrow no learning with gradient descent...

Spectral Radius and Power Sequence of a Matrix

The *spectral radius* $\rho(\mathbf{W})$ of a matrix $\mathbf{W} \in \mathbb{C}^{n \times n}$ is defined as:

$$\rho(\mathbf{W}) = \max \{|\lambda_1|, \dots, |\lambda_n|\},$$

where $\lambda_1, \dots, \lambda_n$ are the (real or complex) eigenvalues of \mathbf{W} .

The *spectral radius* is closely related to the behaviour of the *convergence* of the power sequence of a matrix \mathbf{W}

Theorem

Let $\mathbf{W} \in \mathbb{C}^{n \times n}$ with spectral radius $\rho(\mathbf{W})$. Then $\rho(\mathbf{W}) < 1$ *if and only if*

$$\lim_{k \rightarrow \infty} \mathbf{W}^k = \mathbf{0}.$$

On the other hand, if $\rho(\mathbf{W}) > 1$,

$$\lim_{k \rightarrow \infty} \|\mathbf{W}^k\| = \infty.$$

The statement holds for any choice of matrix norm on $\mathbb{C}^{n \times n}$.

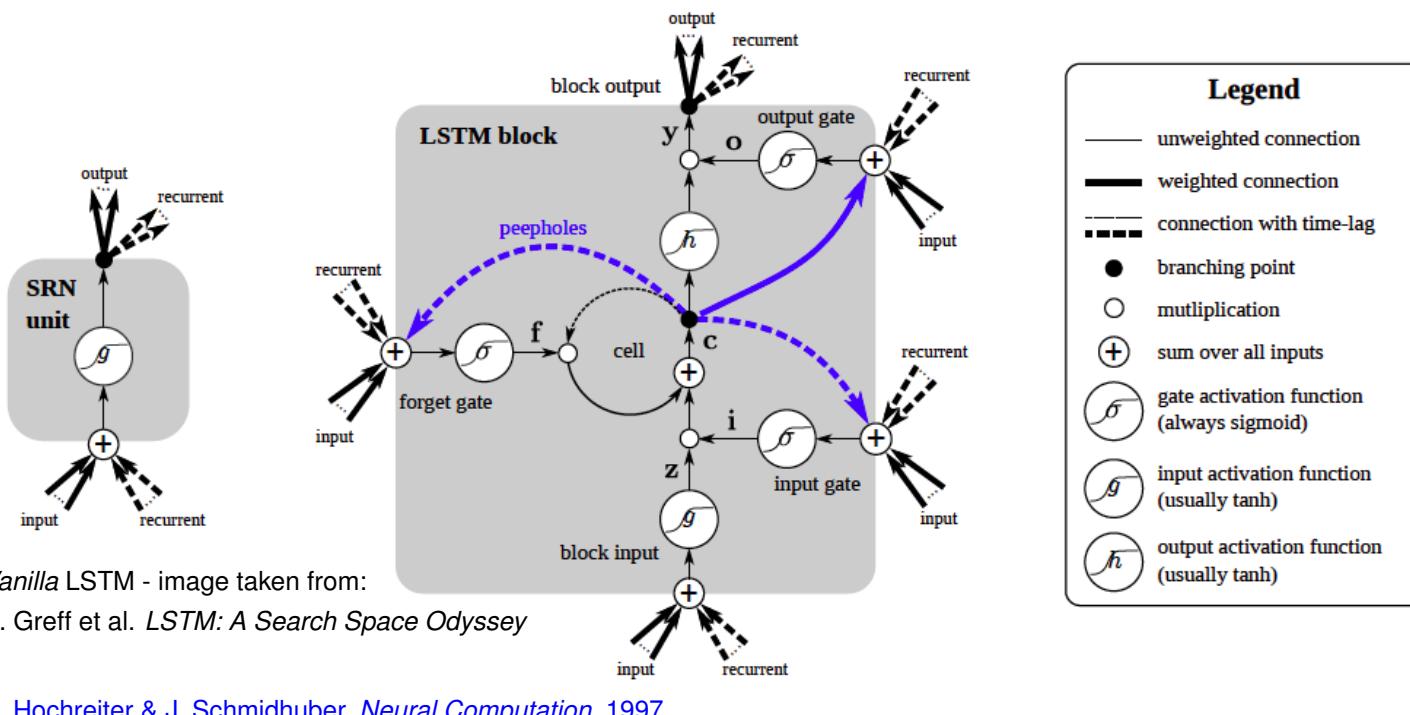
Vanishing/Exploding of Gradients: some remedies

Here are some approaches to try to *reduce* the vanishing/exploding gradients problem

- Architectural
 - Long Short-Term Memory or Gated Recurrent units
 - Reservoir Computing: Echo State Networks and Liquid State Machines
- Algorithmic
 - *Clipping gradients (avoids exploding gradients)*
 - *Hessian Free Optimization*
 - Smart Initialization: pre-training techniques

Long Short-Term Memory

Extension of RNN that can deal with long-term temporal dependencies
Mechanism that allows the networks to “remember” relevant information for a long period of time



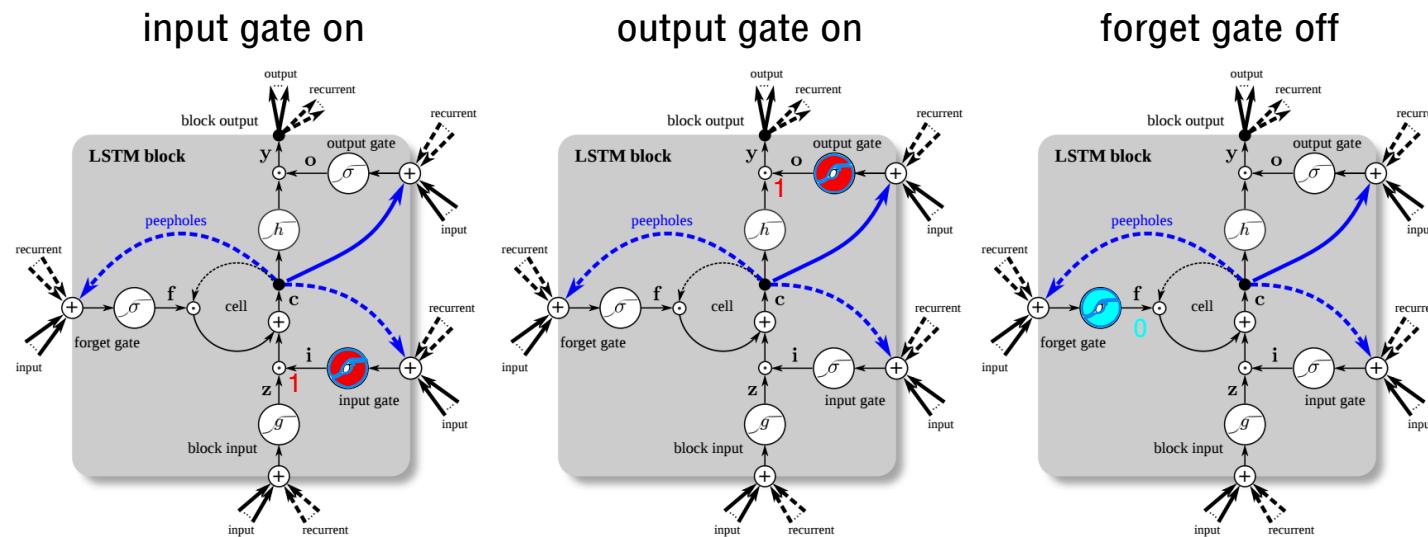
Vanilla LSTM - image taken from:

K. Greff et al. *LSTM: A Search Space Odyssey*

S. Hochreiter & J. Schmidhuber, *Neural Computation*, 1997

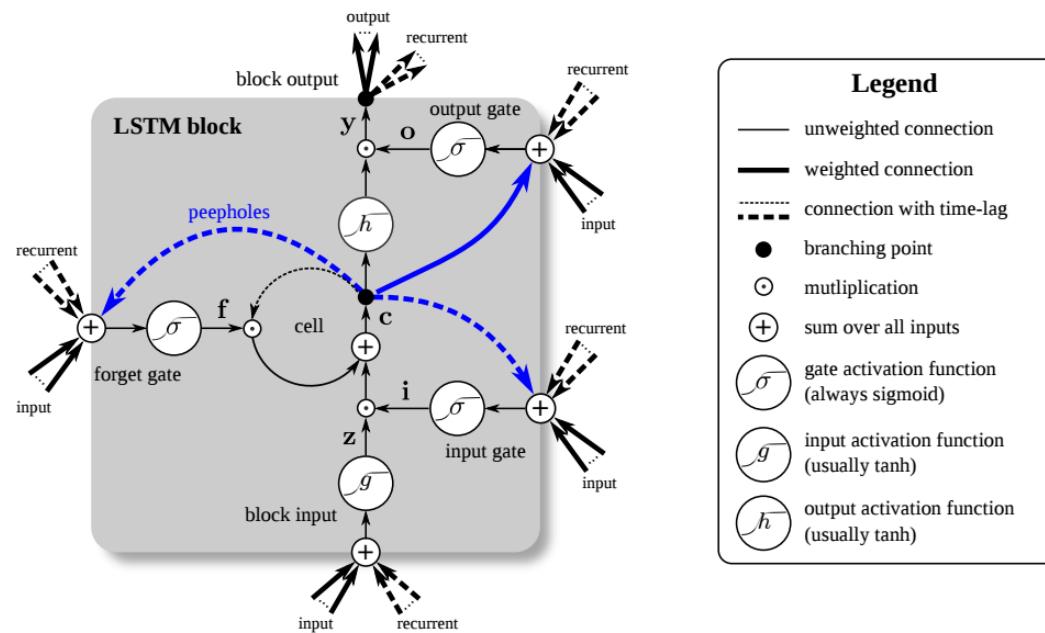
Long Short-Term Memory (*Vanilla* LSTM)

- Exploits a **linear** memory cell (*state*) that integrates input information through time
 - memory obtained by self-loop
 - gradient not down-sized by Jacobian of sigmoidal function \Rightarrow **no vanishing gradient**
- 3 gate units (with sigmoid: soft version of a 0/1 switch) control the information flow via multiplicative connections
 - input gate “on”: let input to flow in the memory cell
 - output gate “on”: let the current value stored in the memory cell to be read in output
 - forget gate “off”: let the current value stored in the memory cell to be reset to 0



Long Short-Term Memory (*Vanilla* LSTM)

- peepholes connections allow to directly control all gates to allow for easier learning of precise timings
- full back-propagation through time (BPTT) training introduced only in 2005



Long Short-Term Memory

Many variations of the *Vanilla* LSTM architecture have been studied, as

- No Input Gate (NIG)
- No Forget Gate (NFG)
- No Output Gate (NOG)
- No Input Activation Function (NIAF)
- No Output Activation Function (NOAF)
- No Peepholes (NP)

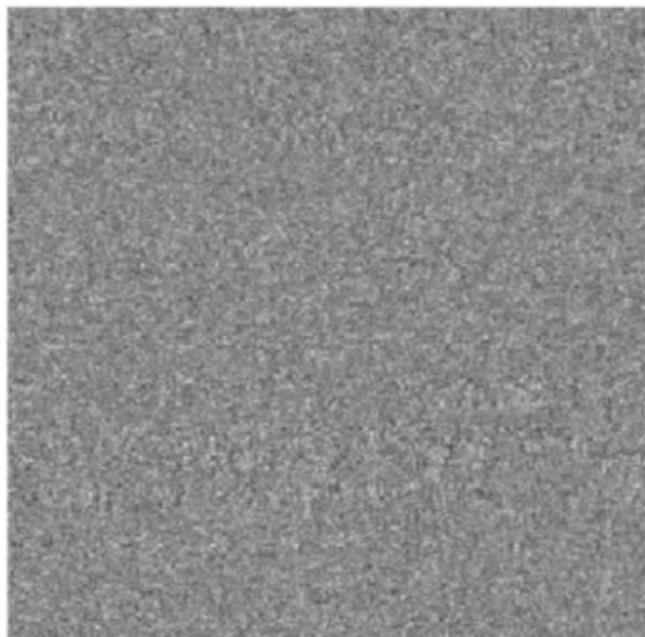
However:

- *Vanilla* LSTM performs reasonably well in general and variations do not significantly improve the performance
- the forget gate is crucial for LSTM performance
- if the cell state is unbounded the output activation function is needed

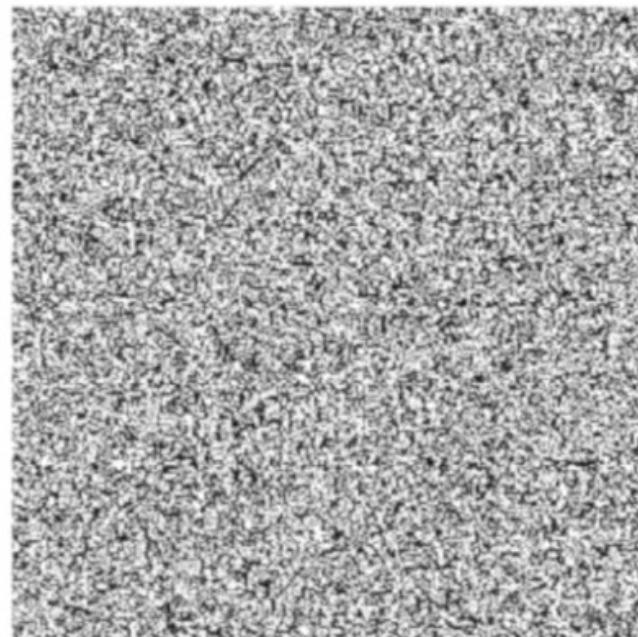
Long Short-Term Memory

gradient as an image

127



127



Simplifying LSTM: Gated Recurrent Units

Is it possible to simplify LSTM units ?

- Gated Recurrent Units (GRU) do that by using a single gating unit that simultaneously controls the forgetting factor and the decision to update the state unit

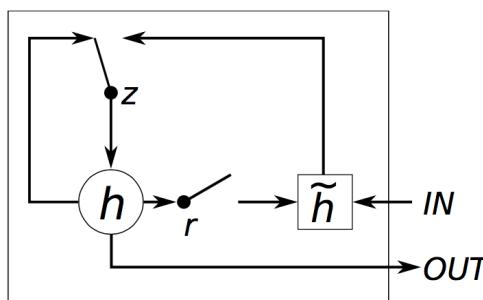


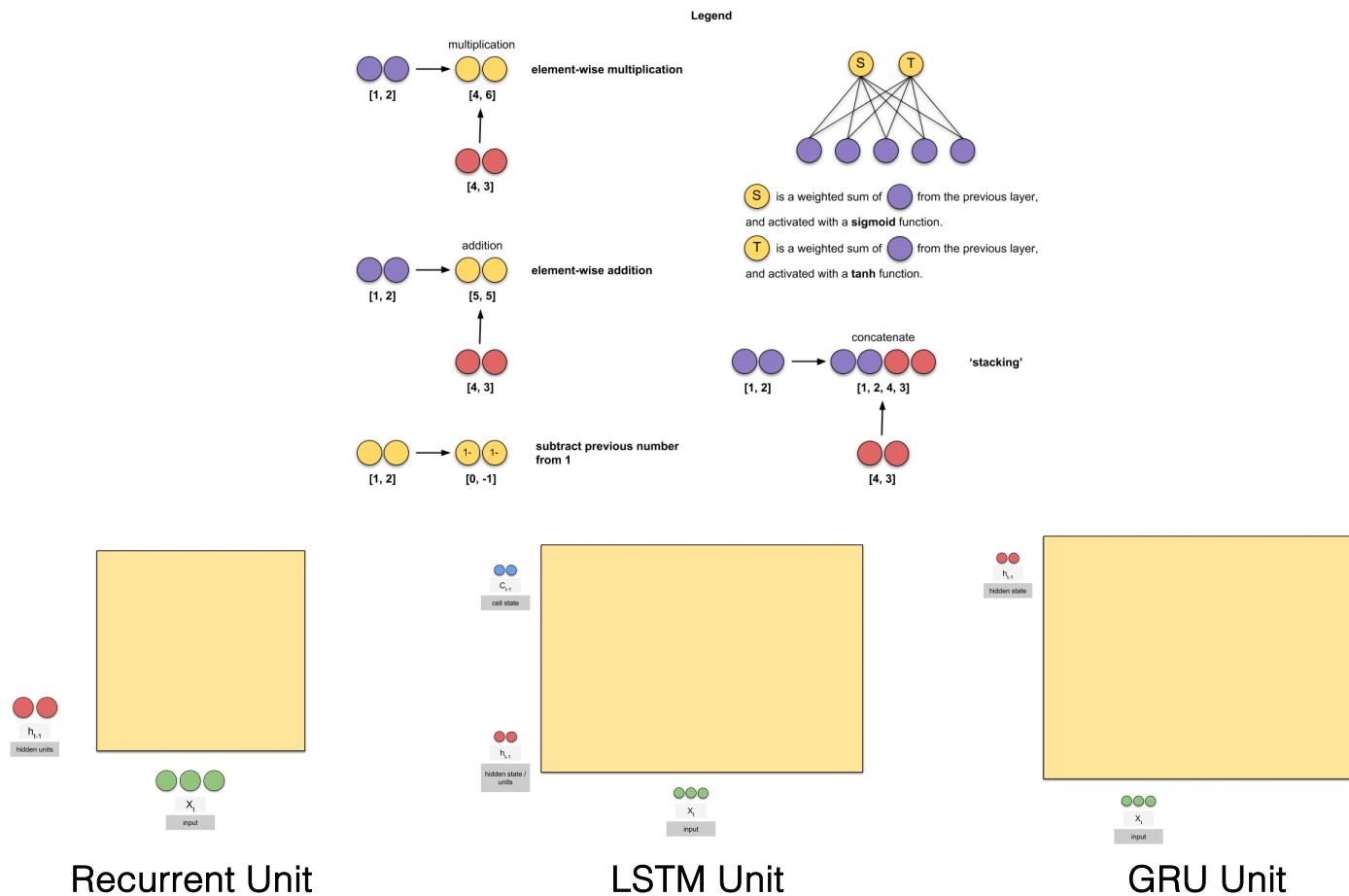
image taken from Chung et al., NIPS Workshop on Deep Learning, 2014

$$\mathbf{h}^{(t)} = \mathbf{z}^{(t)} \odot \mathbf{h}^{(t-1)} + \underbrace{(\mathbf{1} - \mathbf{z}^{(t)}) \odot \sigma(\mathbf{Ux}^{(t)} + \mathbf{W}(\mathbf{r}^{(t)} \odot \mathbf{h}^{(t-1)}))}_{\tilde{\mathbf{h}}^{(t)}}$$

- the *update* gate z selects whether the hidden state is to be updated with a new hidden state \tilde{h}
- the *reset* gate r decides whether the previous hidden state is ignored

Cho et al., EMNLP, 2014

Summing Up



Recap: Vanishing/Exploding of Gradients: some remedies

Here are some approaches to try to *reduce* the vanishing/exploding gradients problem

- Architectural

- Long Short-Term Memory or Gated Recurrent units
- Reservoir Computing: Echo State Networks and Liquid State Machines

Done!



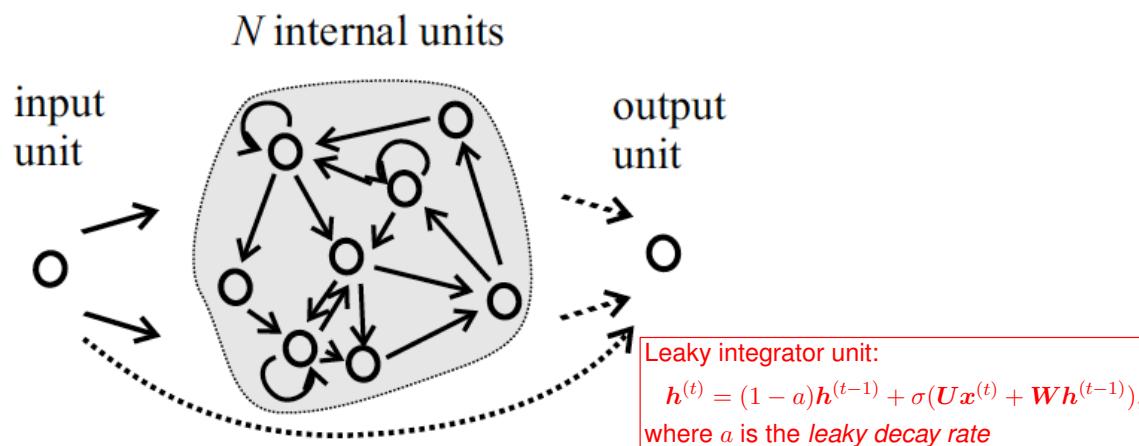
- Algorithmic

- *Clipping gradients (avoids exploding gradients)*
- *Hessian Free Optimization*
- Smart Initialization: pre-training techniques

Reservoir Computing: Echo State Nets and Liquid State Machines

Idea: fix the input-to-hidden and hidden-to-hidden connections at random values and only learn the output units connections

image taken from:
H. Jaeger, NIPS 2002



- Echo State Networks (ESN): standard recurrent neurons (+ leaky integrators)
Jaeger, Technical Report GMD Report 148, 2001
- Liquid State Machines (LSM): spiking integrate-and-fire neurons and dynamic synaptic connection models
Maass et al., *Neural Computation*, 2002

Reservoir Computing: additional details

- ➊ A recurrent neural network with no output (the *reservoir*) is randomly created and remains unchanged during training
 - it is passively excited by the input signal $x^{(t)}$
 - the hidden state $h^{(t)}$ maintains a nonlinear version of input history
 - in order to produce a “rich” set of dynamics, the *reservoir* should
 - be **big** (hundreds to thousands units)
 - be **sparsely** (hidden weight matrix W up to 20% possible connections) and **randomly** (uniform distribution symmetric around zero) connected
 - satisfy the ***echo state property***, i.e. the effect of the current state $h^{(t)}$ and the current input $x^{(t)}$ on a future state $h^{(t+\tau)}$ should vanish gradually as time passes ($\tau \rightarrow \infty$), in practice:
$$\text{spectral radius } \rho(W) < 1, \text{ i.e. } W \text{ is } \textcolor{red}{\text{contractive}}$$
- ➋ on the contrary, the input (U) and *optional* output feedback weight matrices are dense (still random with uniform distribution)
- ➌ The output is computed as a linear combination of the input-excited reservoir, and the linear combination is obtained by linear regression

Reservoir Computing: additional details

- alternative topologies for the reservoir, e.g. small-world, scale-free, and topologies inspired by spatial growth, were compared with no significant improvement over simple random networks
- on the other side, very simple topologies can be very effective

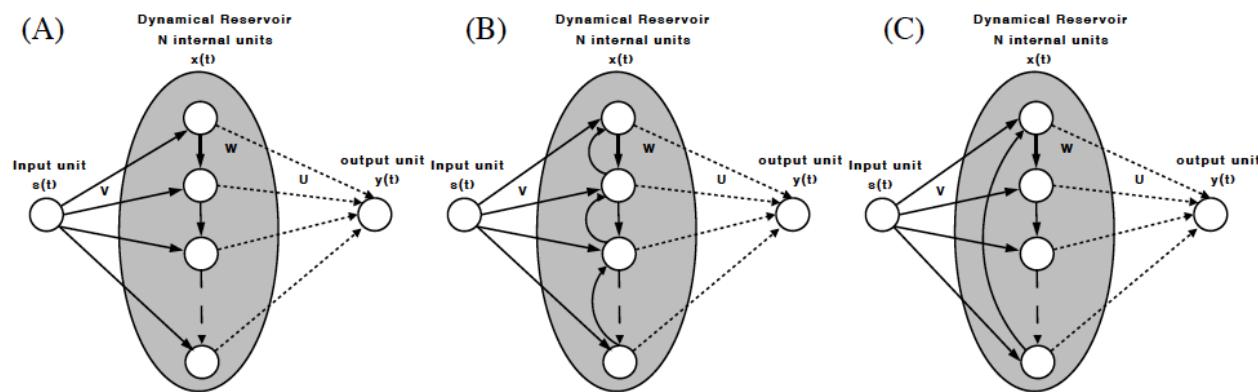


Fig. 2. (A) Delay Line Reservoir (DLR). (B) Delay Line Reservoir with feedback connections (DLRB). (C) Simple Cycle Reservoir (SCR).
from Rodan and Tino, *IEEE TNN*, 2011

- a simple cycle reservoir topology obtains performances comparable to those of ESN
- competitive reservoirs can be constructed deterministically
- the memory capacity of simple linear cyclic reservoirs can be made to be arbitrarily close to the proved optimal memory capacity value

Reservoir Computing: additional details

- *Intrinsic Plasticity* (IP) has attracted a wide attention in the reservoir computing community
- IP is a computationally efficient online learning rule to adjust threshold and gain of sigmoid reservoir neurons
 - it drives the neurons' output activities to approximate exponential distributions
 - the exponential distribution maximizes the entropy of a non-negative random variable with a fixed mean, thus enabling the neurons to transmit maximal information

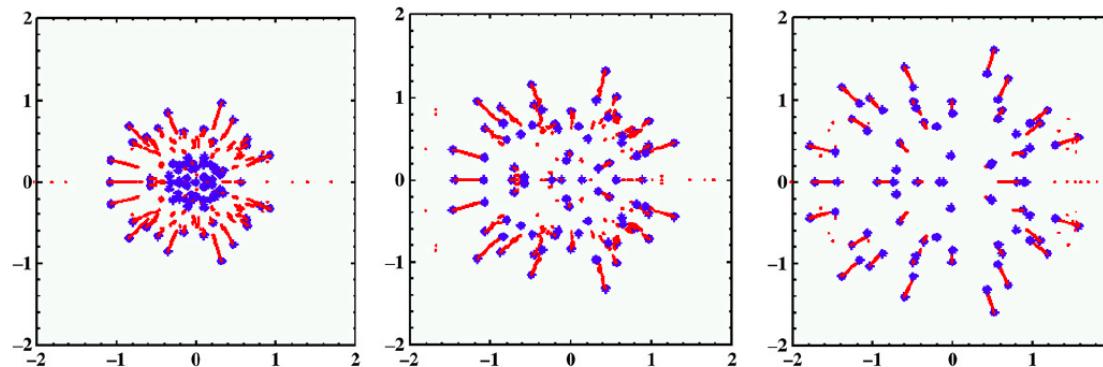


Fig. 6. Changes in the eigenvalue distribution introduced by IP for epochs 0–4, 4–9, 9–14 of a 100 neuron network learning a one-step-ahead prediction of the Mackey–Glass attractor (as detailed in Section 4). Larger dots denote the positions of weight matrix eigenvalues at the beginning and end of the 5-epoch intervals, small dotted lines indicate the path of their movements away from the center.

from Steil, *Neural Networks*, 2007

Deep Echo State Networks

Deep layered organization of RC models have been investigated in terms of

- occurrence of multiple time-scale
- increasing of richness of the dynamics

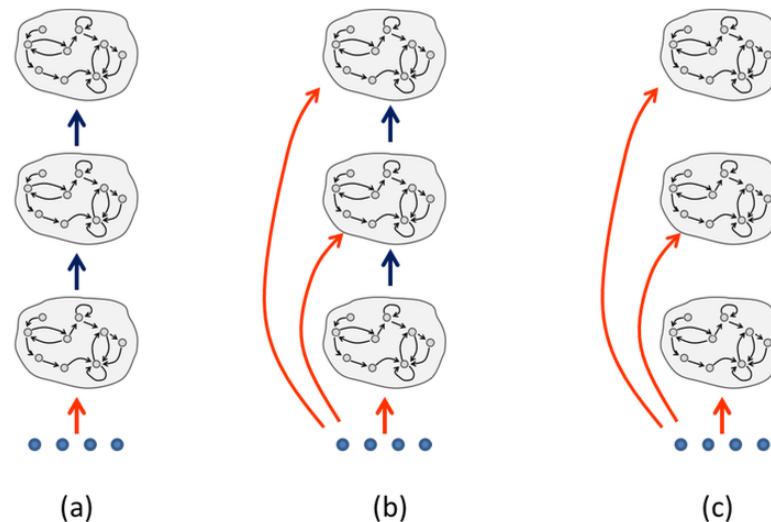


Figure 1: Deep RC architectures: (a) deepESN, (b) deepESN-IA, (c) groupedESN.
from Gallicchio and Micheli, *Neurocomputing*, 2017

Deep Echo State Networks: Richness of Dynamics

Entropy of output distribution $f_i(x)$ of reservoir unit i : $-\int f_i(x) \log f_i(x) dx$

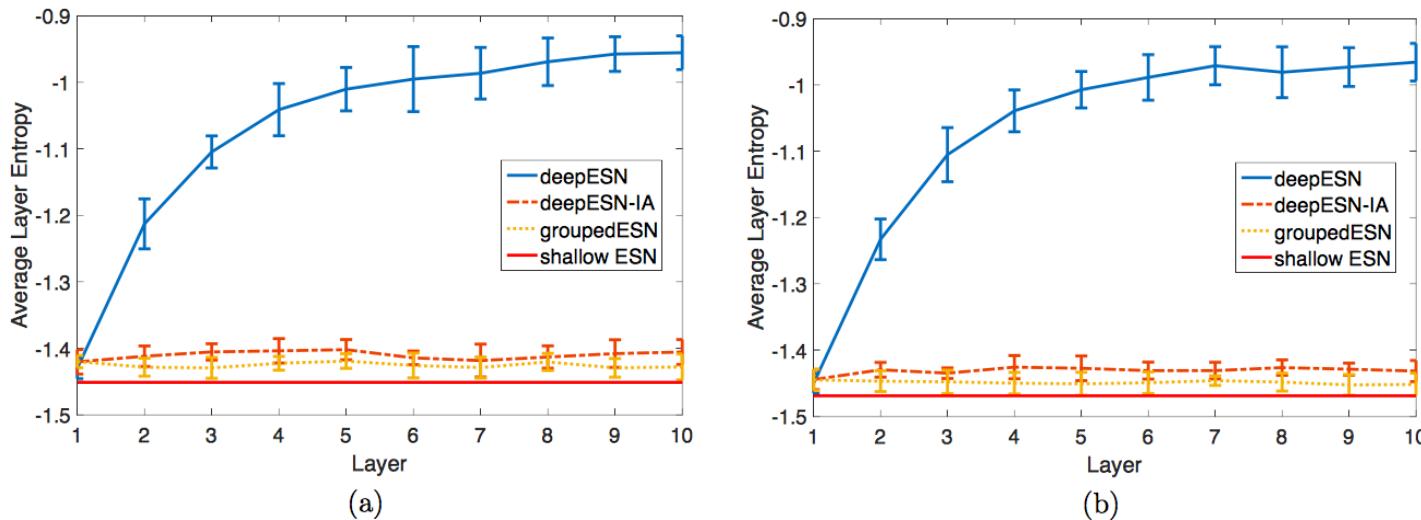


Figure 6: Layer-wise averaged entropy of reservoir states on the Artificial dataset (a) and on the Wikipedia dataset (b), computed for deepESN, groupedESN and deepESN-IA, with $a = 0.55$ and $\rho = 0.9$ for every layer (or sub-reservoir) and IP learning. For groupedESN the results refer to sub-reservoirs. The average entropy of the shallow ESN counterpart is reported as a continuous red line across each plot.

from Gallicchio and Micheli, *Neurocomputing*, 2017

Deep Echo State Networks: Memory Capacity

Task: reconstruct the input with increasing delay

- target $\mathbf{y}_k^{(t)} = \mathbf{x}^{(t-k)}$, $\forall k \in [0, \dots, \infty]$
- Memory Capacity (MC):

$$\sum_{k=0}^{\infty} r^2(\mathbf{x}^{(t-k)}, \mathbf{o}_k^{(t)})$$

where $r^2(\mathbf{x}^{(t-k)}, \mathbf{o}_k^{(t)})$ is the squared correlation coefficient between:

- the input $\mathbf{x}^{(t-k)}$ with delay k
- the corresponding output $\mathbf{o}_k^{(t)}$ generated by the network at time t for delay k

Deep Echo State Networks

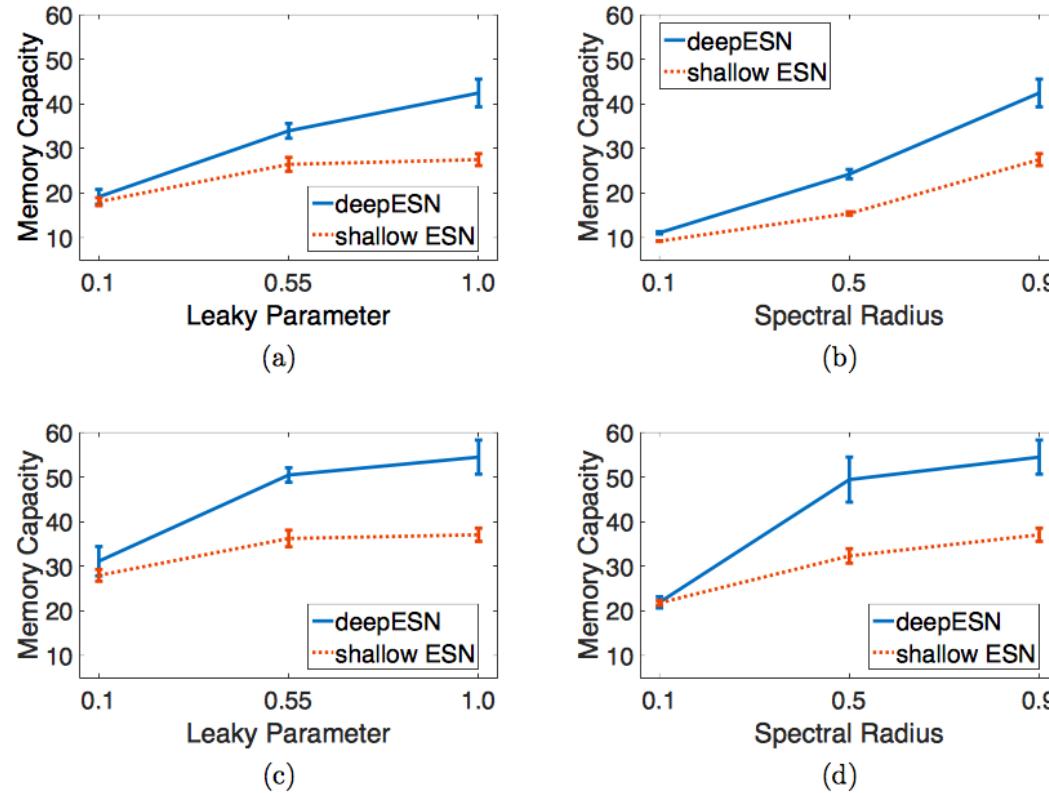


Figure 7: MC results of deepESN and shallow ESN for different values of the spectral radius ρ and of the leaky parameter a . (a): different values of a , without IP, (b): different values of ρ , without IP, (c): different values of a , with IP, (d): different values of ρ , with IP. For deepESN each reservoir layer has the same hyper-parametrization.

from Gallicchio and Micheli, *Neurocomputing*, 2017

Deep Recurrent Neural Networks: Examples in Sequential Domains

How to Construct Deep Recurrent Neural Networks

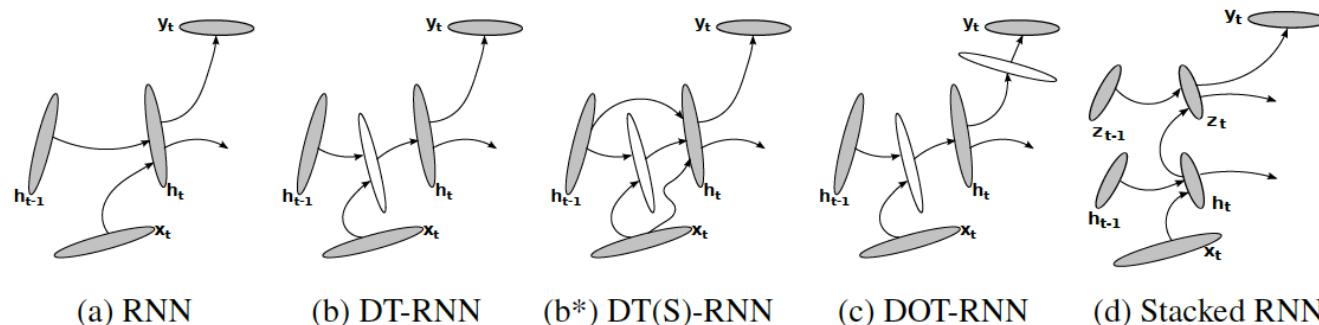


Figure 2: Illustrations of four different recurrent neural networks (RNN). (a) A conventional RNN. (b) Deep Transition (DT) RNN. (b*) DT(S)-RNN with shortcut connections (c) Deep Transition, Deep Output (DOT) RNN. (d) Stacked RNN
image taken from: Pascanu et al., arXiv:1312.6026v5

Empirical experiments on polyphonic music prediction and language modeling showed that RNN benefits from having a deeper architecture

Pascanu et al., arXiv:1312.6026v5, 2014

Examples of Deep Architectures

- Novel deep RNN architectures/mechanisms have been designed for specific applications
- Examples of architectures using
 - Encoder-Decoder
 - Alignment/Attention
- image to text applications

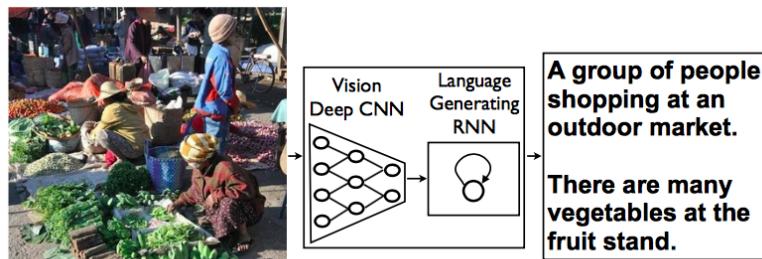


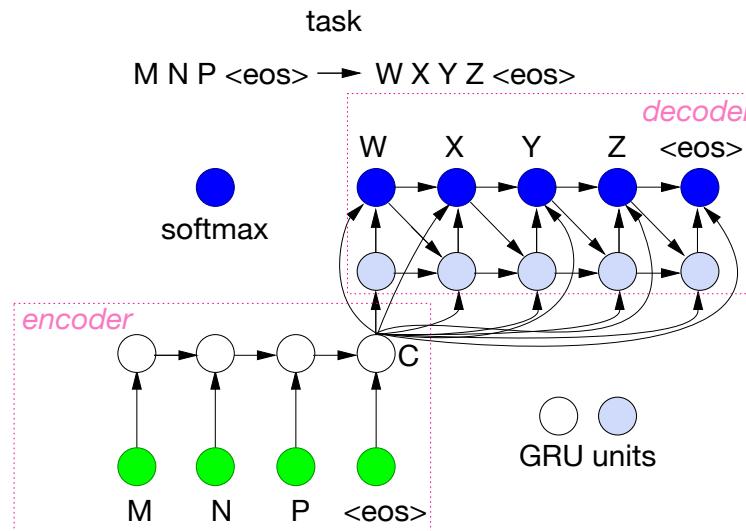
image taken from: Vinyals et al., *IEEE TPAMI*, 2016

Other interesting architectures that we do not have time to look at:

- Adaptive Computation Time, “External” memory, Pointers

Encoding-Decoding Networks

$$\text{softmax}(\mathbf{z})_i = \frac{\exp(z_i)}{\sum_j \exp(z_j)}$$



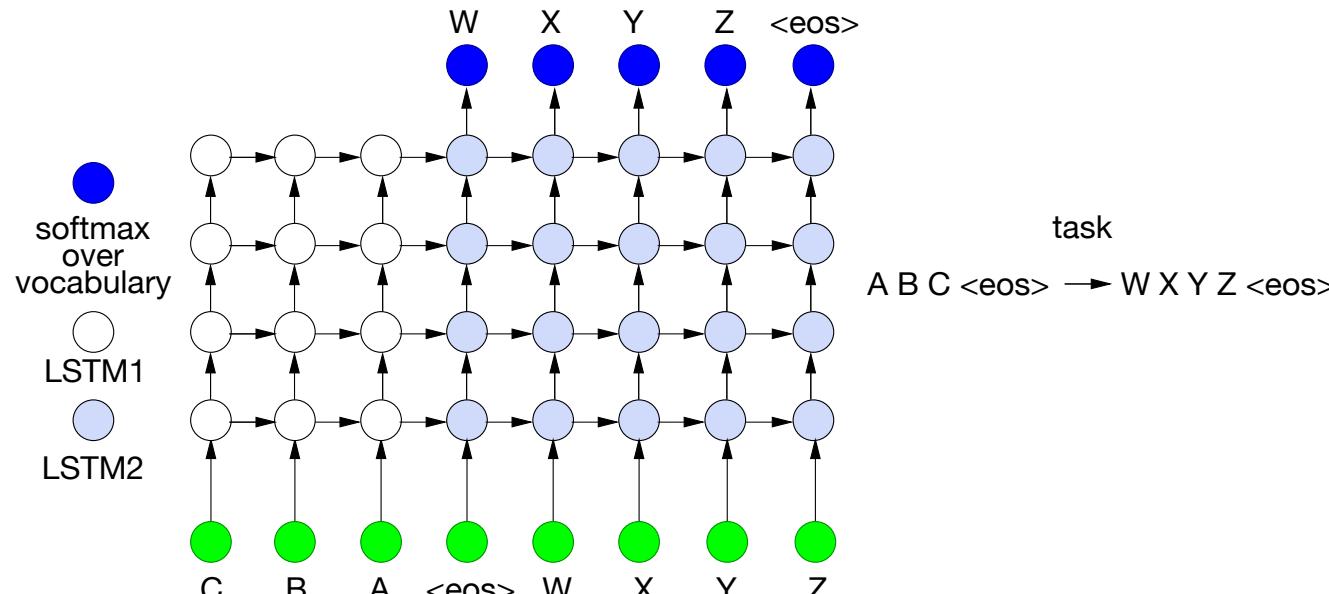
End-to-end learning of $(\max_{\theta} \frac{1}{N} \sum_{n=1}^N \log P_{\theta}(\mathbf{s}_n^d | \mathbf{s}_n^i))$

$$P(\mathbf{o}_t | \mathbf{o}_{t-1}, \mathbf{o}_{t-2}, \dots, \mathbf{o}_1, \mathbf{c}) = g(\mathbf{h}_t, \mathbf{o}_{t-1}, \mathbf{c})$$

where $\mathbf{h}_t = f(\mathbf{h}_{t-1}, \mathbf{o}_{t-1}, \mathbf{c})$

Cho et al., EMNLP 2014

Sequence to Sequence Networks



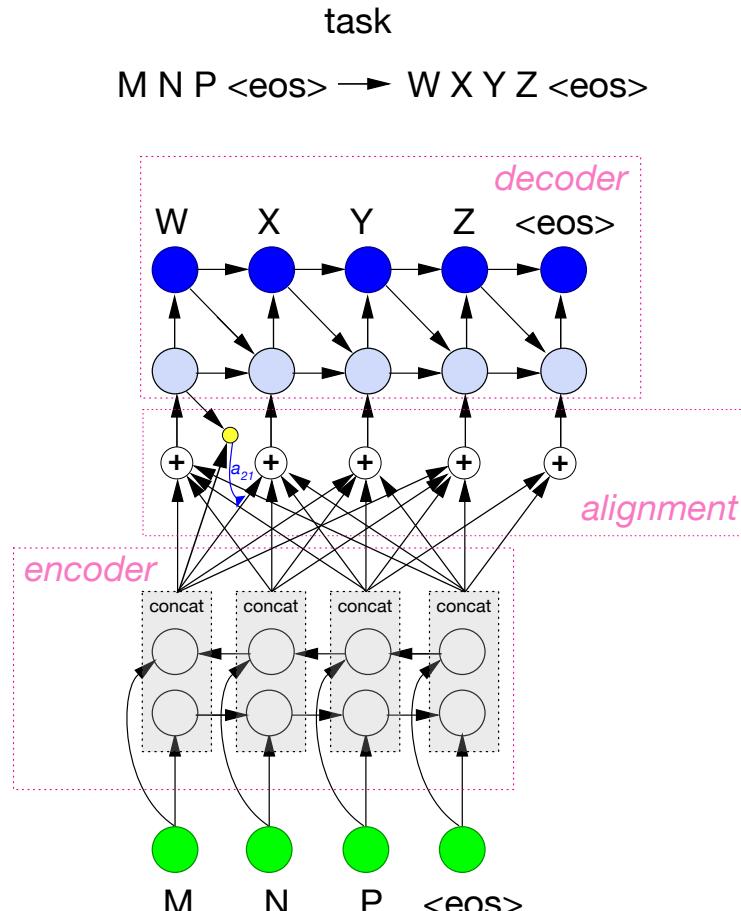
End-to-end learning of

$$P(\mathbf{o}_1, \mathbf{o}_2, \dots, \mathbf{o}_{T'} | \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T) = \prod_{t=1}^{T'} P(\mathbf{o}_t | \mathbf{v}, \mathbf{o}_{t-1}, \mathbf{o}_{t-2}, \dots, \mathbf{o}_1)$$

where $\mathbf{v} = \mathbf{h}_{LSTM1}(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T)$

Sutskever et al., NIPS 2014

Encoding-Decoding Networks with Alignment



Bahdanau et al., ICLR 2015

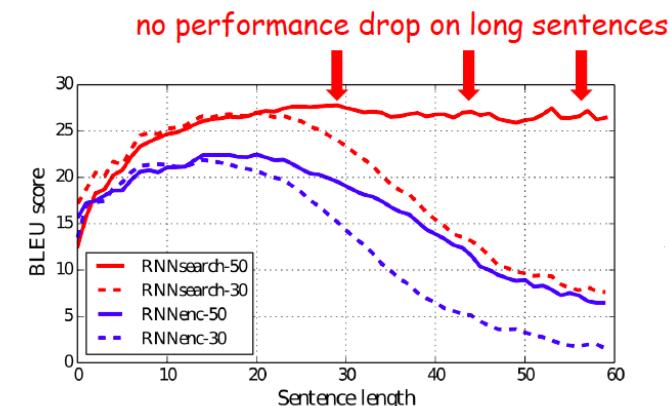
$$P(\mathbf{o}_t | \mathbf{o}_{t-1}, \dots, \mathbf{o}_1, \mathbf{s}^i) = g(\mathbf{o}_{t-1}, \mathbf{h}_t^{dec}, \mathbf{c}_t)$$

$$\mathbf{h}_t = f(\mathbf{h}_{t-1}^{dec}, \mathbf{o}_{t-1}, \mathbf{c}_t)$$

$$\mathbf{c}_t = \sum_{j=1}^T \alpha_{ij} \mathbf{h}_j^{enc}$$

$$\alpha_{ij} = \text{softmax}_j(e_{ij})$$

$$e_{ij} = a(\mathbf{h}_{i-1}^{dec}, \mathbf{h}_j^{enc})$$



One example of Image-Text Alignment (Karpathy & Fei-Fei, IEEE TPAMI, 2017)

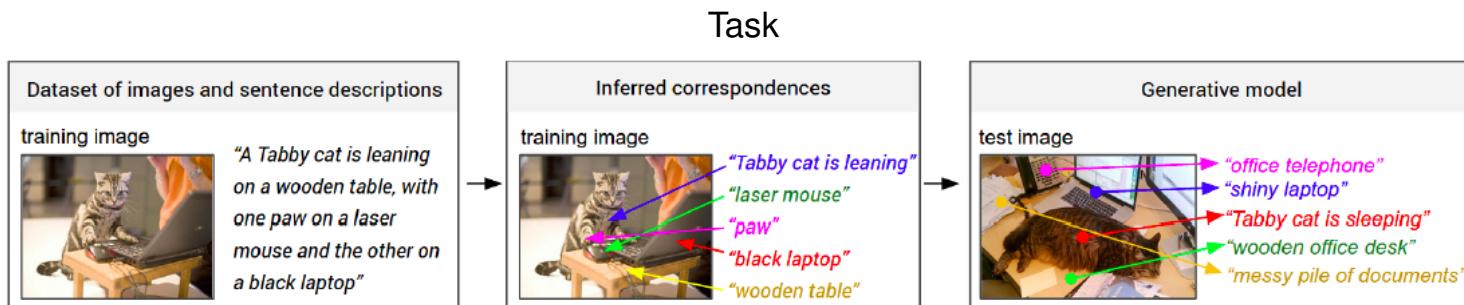


Figure 2. Overview of our approach. A dataset of images and their sentence descriptions is the input to our model (left). Our model first infers the correspondences (middle, Section 3.1) and then learns to generate novel descriptions (right, Section 3.2).

