

# Recurrent Neural Networks

Deep Learning Summer School 2017

June 27, 2017

Yoshua Bengio

Montreal Institute for Learning Algorithms

Université de Montréal

Université   
de Montréal

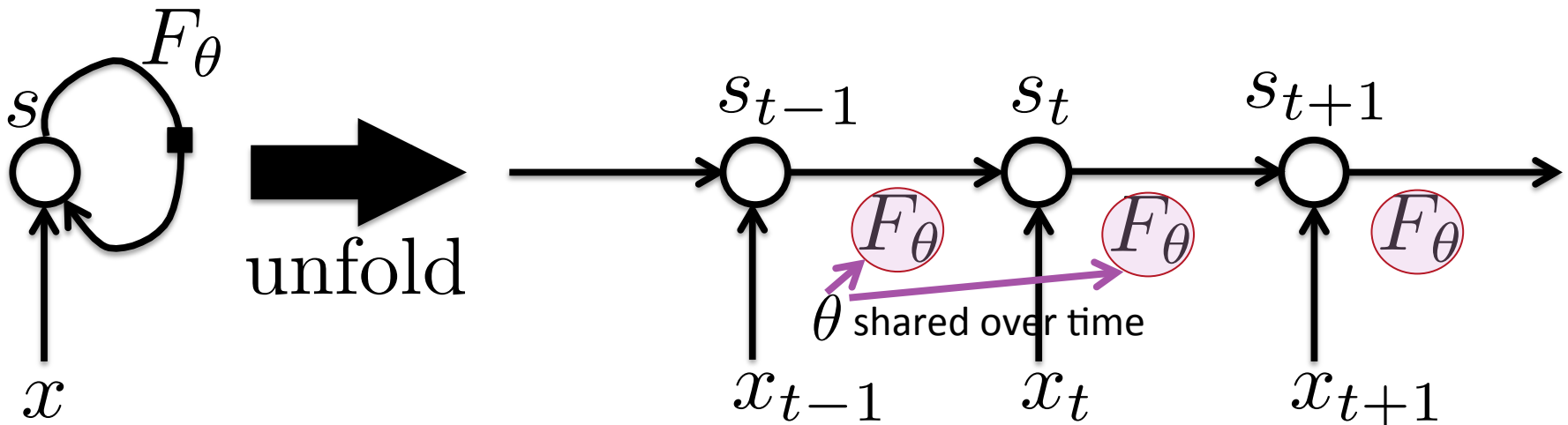
**CIFAR | ICRA**



# Recurrent Neural Networks

- Selectively summarize an input sequence in a fixed-size state vector via a recursive update

$$s_t = F_\theta(s_{t-1}, x_t)$$

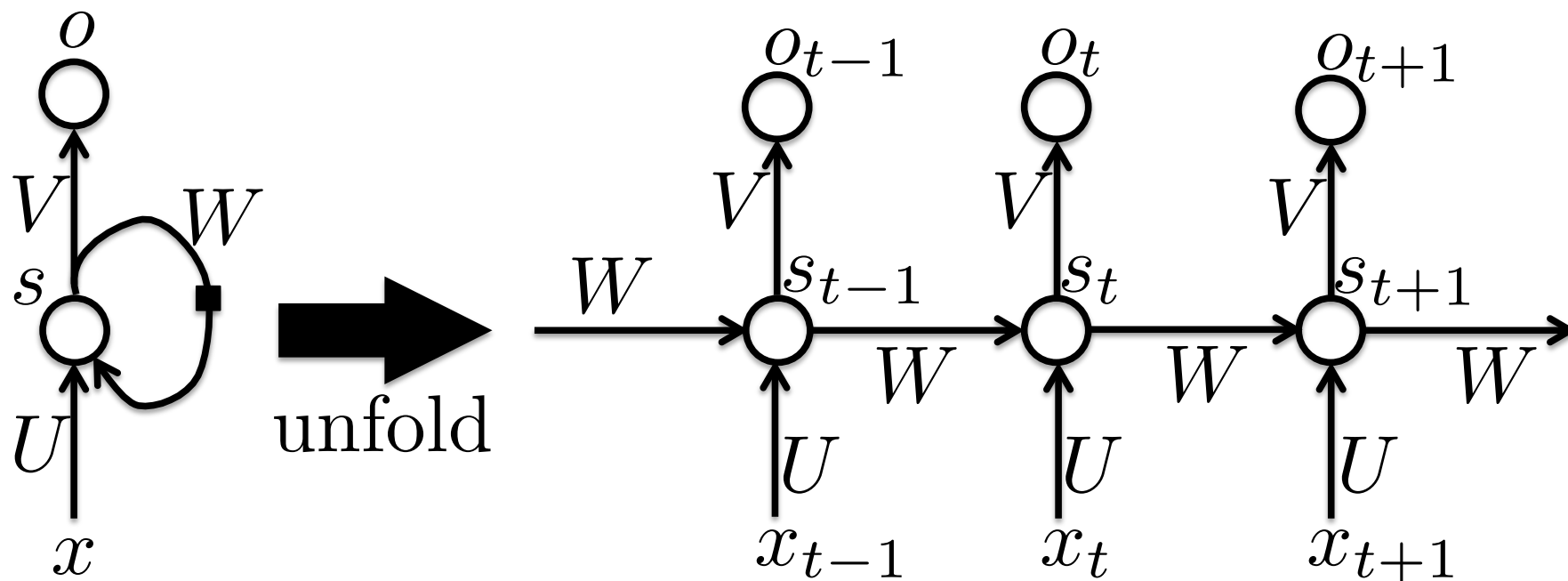


$$s_t = G_t(x_t, x_{t-1}, x_{t-2}, \dots, x_2, x_1)$$

➔ Generalizes naturally to new lengths not seen during training

# Recurrent Neural Networks

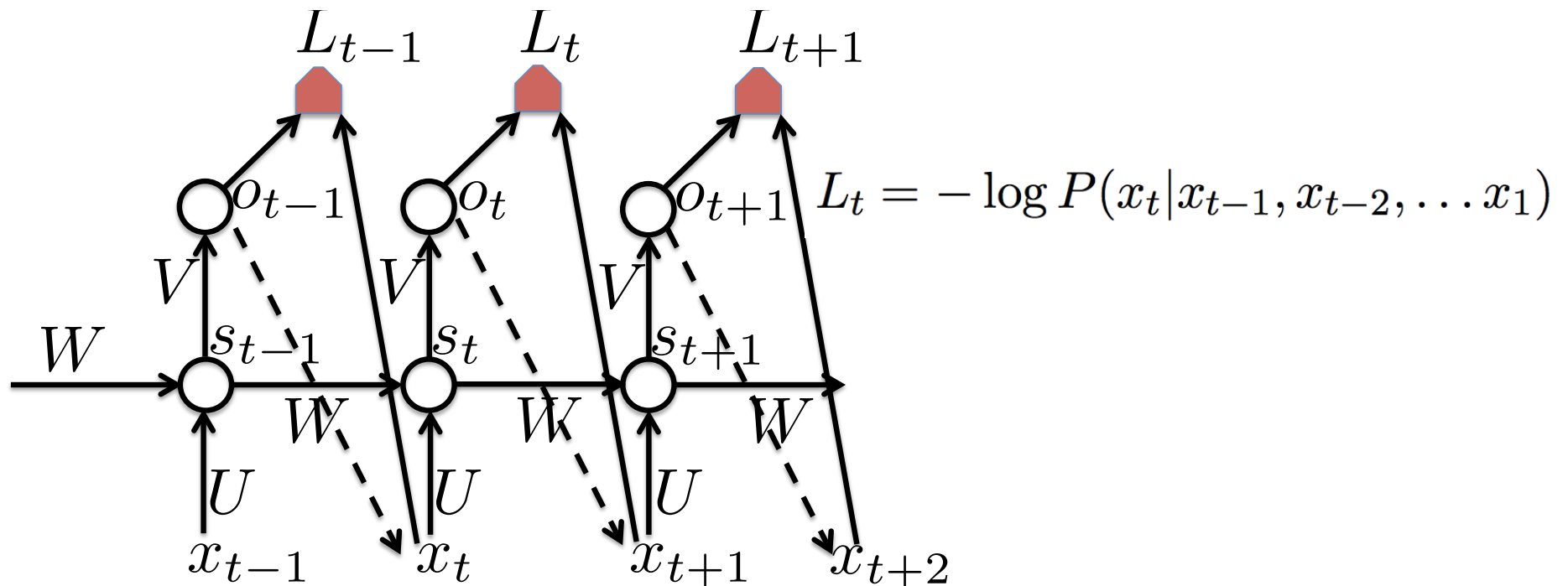
- Can produce an output at each time step: unfolding the graph tells us how to back-prop through time.



# Generative RNNs

- An RNN can represent a fully-connected **directed generative model**: every variable predicted from all previous ones.

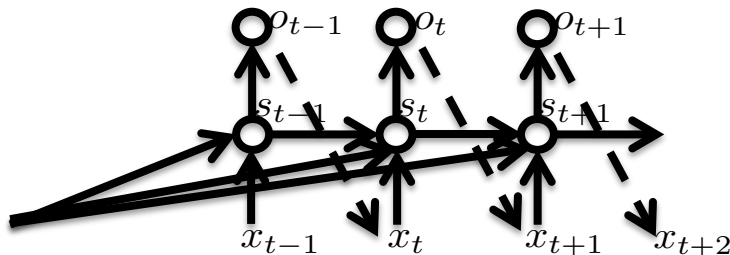
$$P(\mathbf{x}) = P(x_1, \dots, x_T) = \prod_{t=1}^T P(x_t | x_{t-1}, x_{t-2}, \dots, x_1)$$



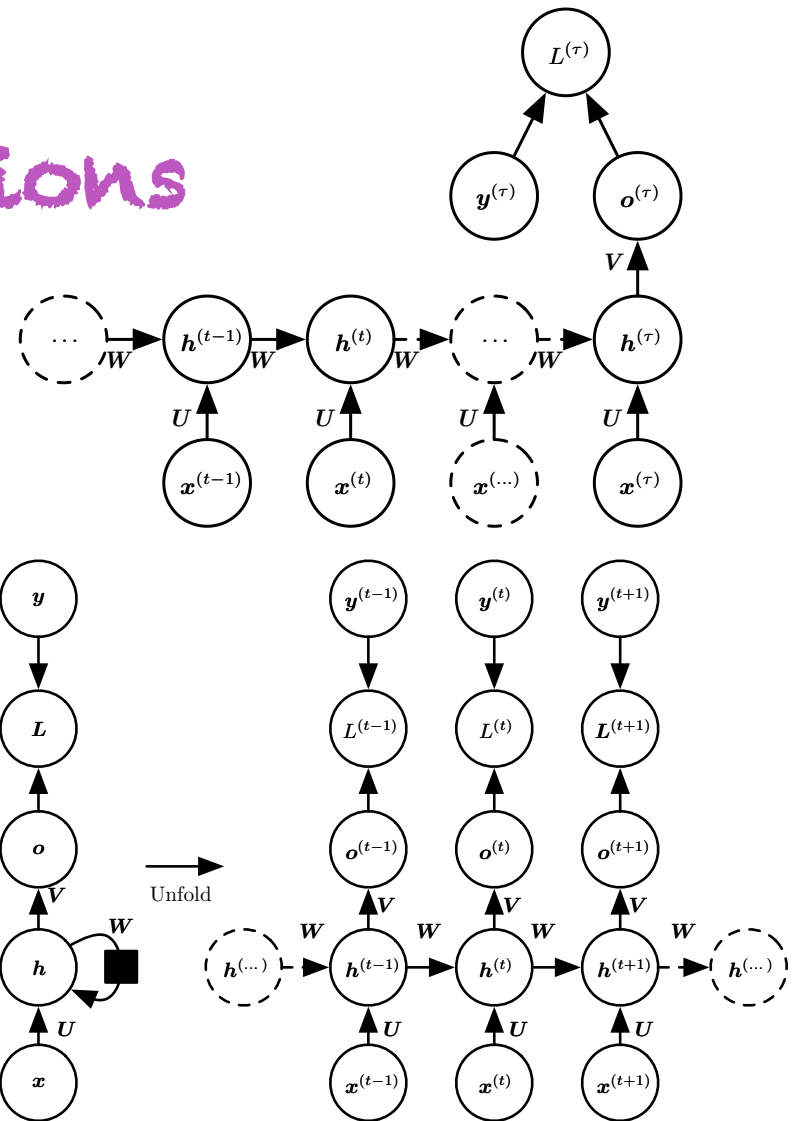
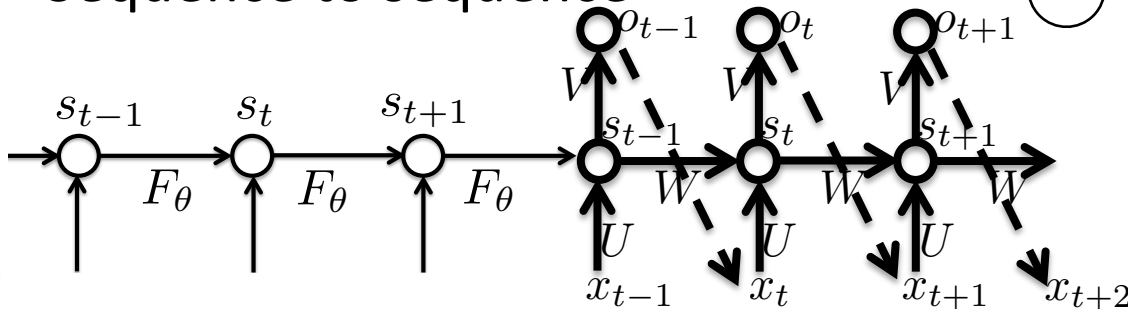


# Conditional Distributions

- Sequence to vector
- Sequence to sequence of the same length, aligned
- Vector to sequence



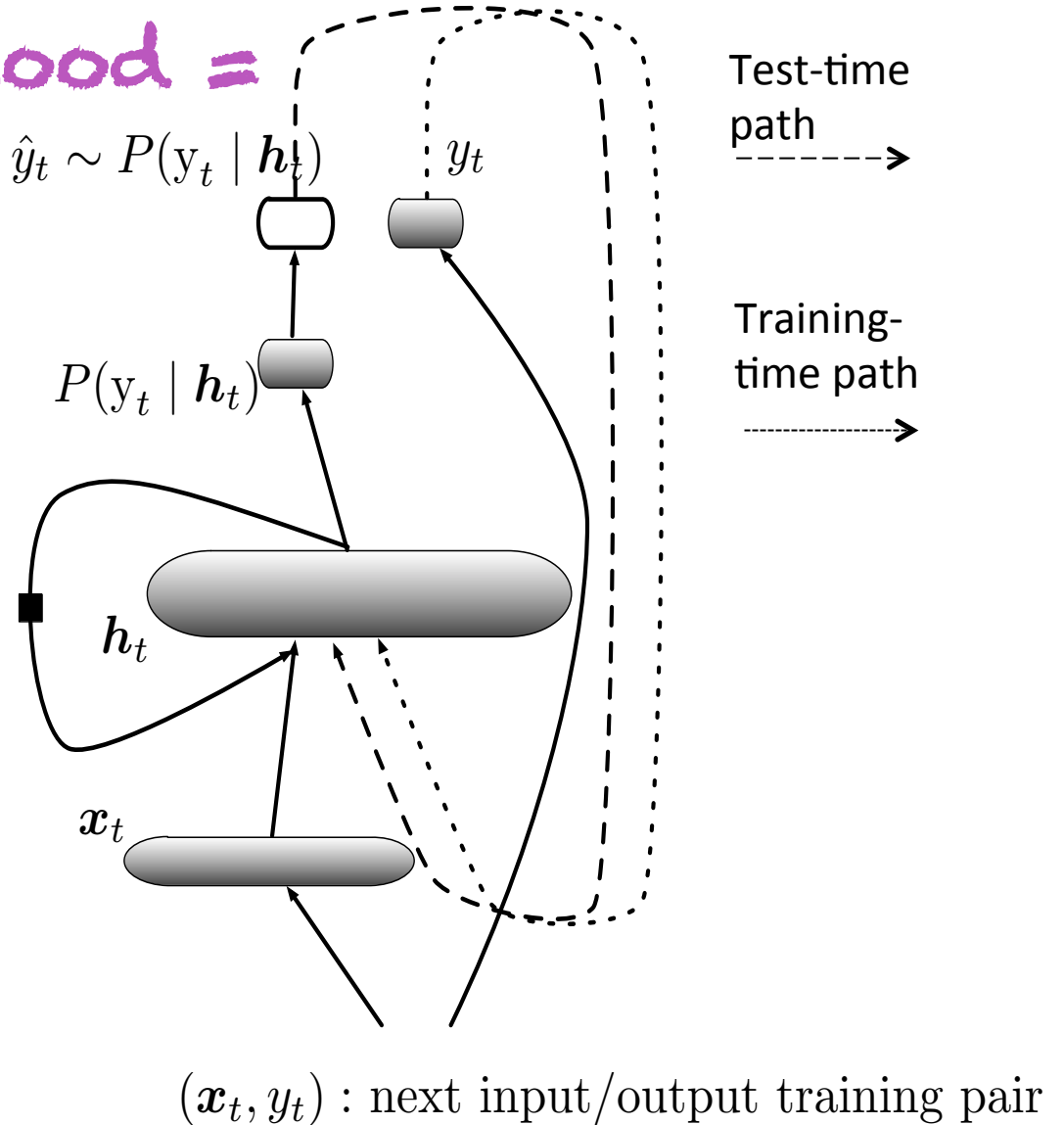
- Sequence to sequence



# Maximum Likelihood = Teacher Forcing

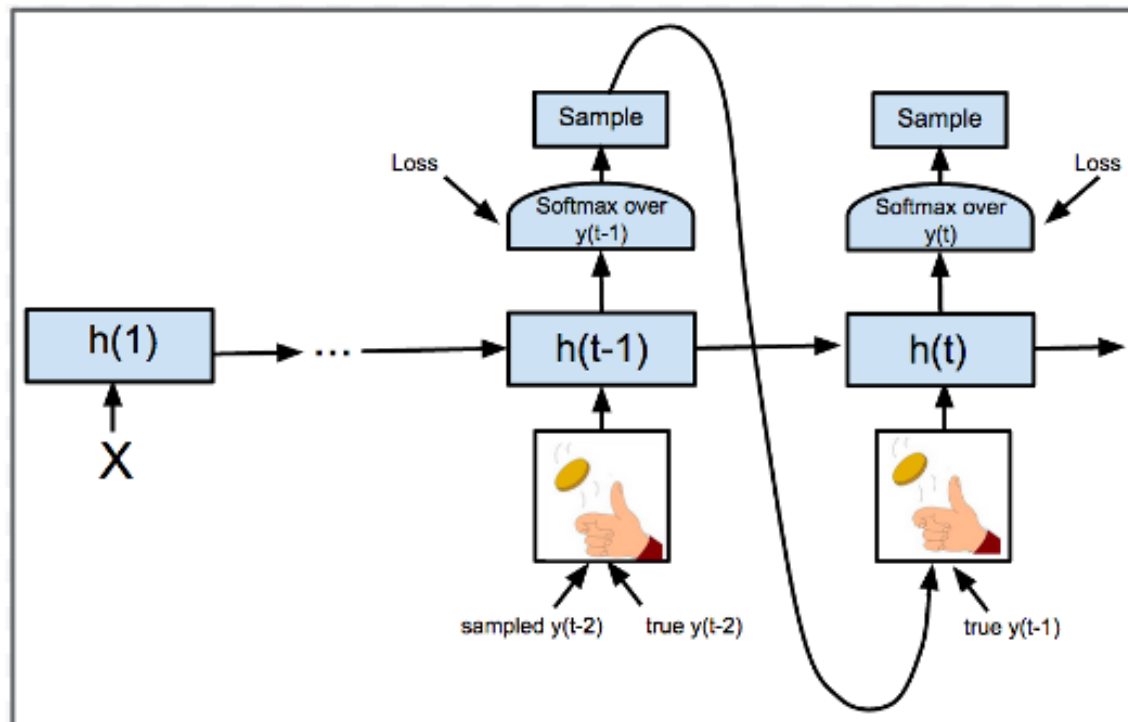
$$\hat{y}_t \sim P(y_t | h_t)$$

- During training, past  $y$  in input is from training data
- At generation time, past  $y$  in input is generated
- Mismatch can cause "compounding error"



# Ideas to reduce the train/generate mismatch in teacher forcing

- Scheduled sampling (*S. Bengio et al, NIPS 2015*)



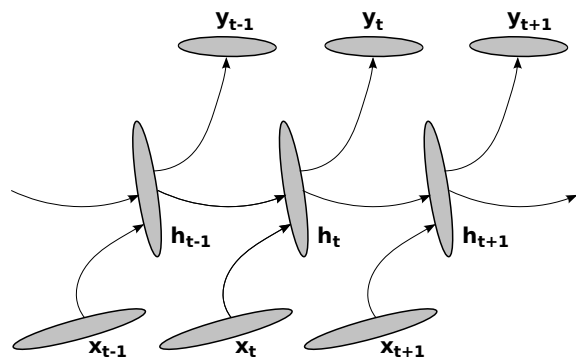
Related to  
SEARN (Daumé et al 2009)  
DAGGER (Ross et al 2010)

Gradually increase the probability of using the model's samples vs the ground truth as input.

- Backprop through open-loop sampling recurrence & minimize long-term cost (but which one? GAN would be most natural → Professor Forcing)

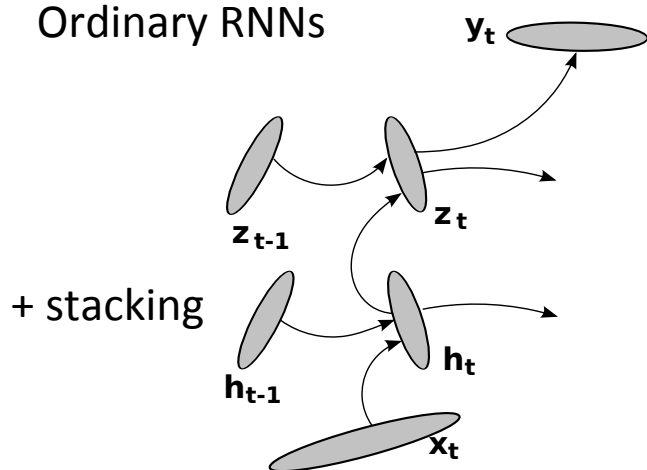
# Increasing the Expressive Power of RNNs with more Depth

- ICLR 2014, *How to construct deep recurrent neural networks*

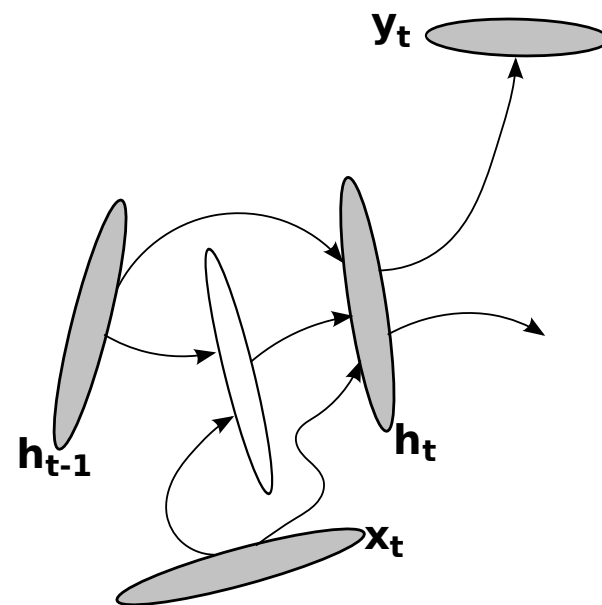
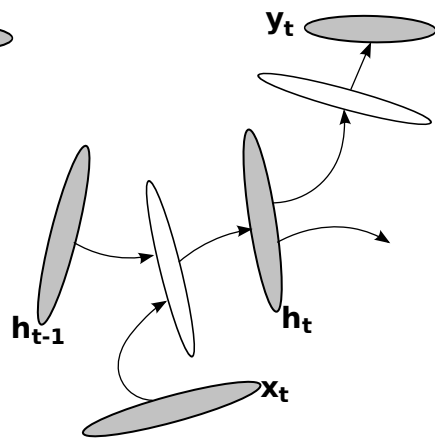


Ordinary RNNs

+ deep hid-to-out  
+ deep hid-to-hid  
+ deep in-to-hid



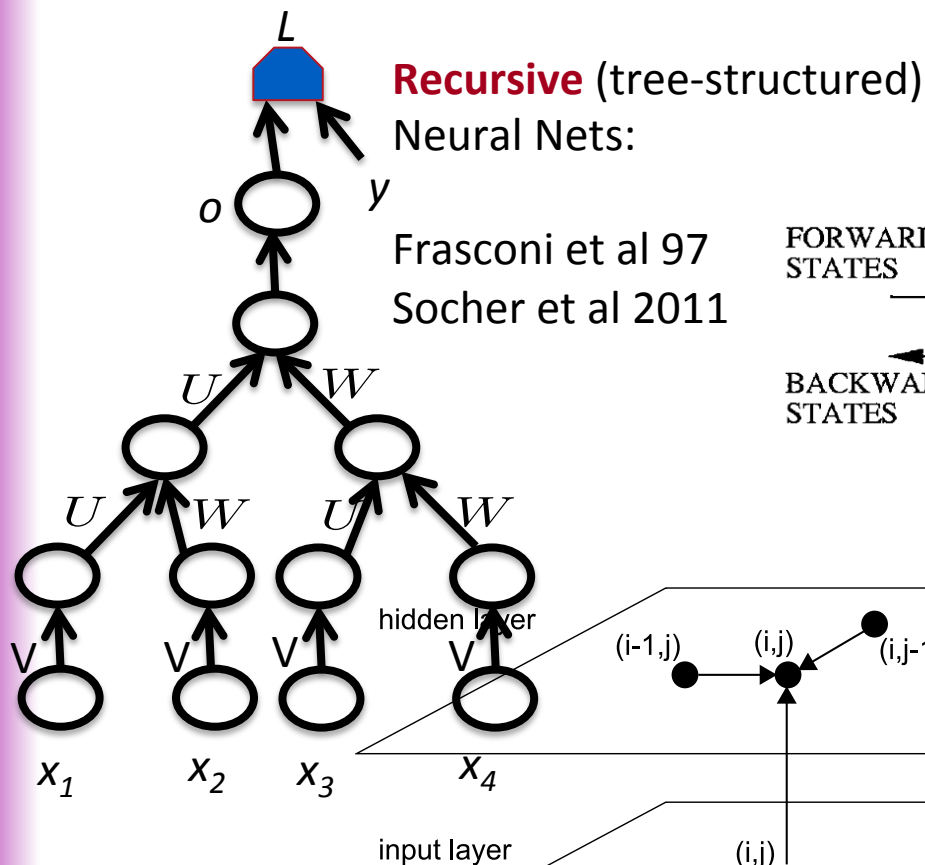
+ stacking



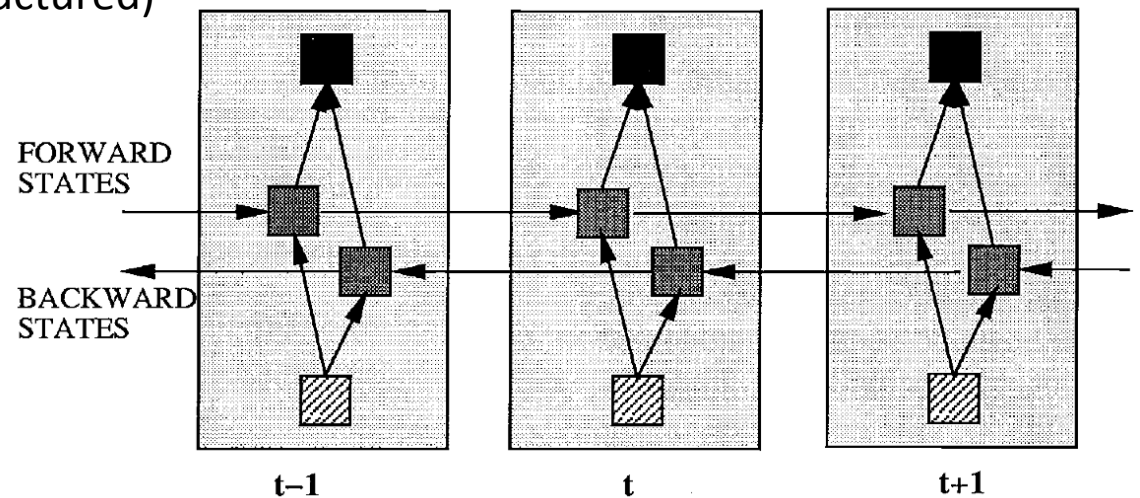
+ skip connections for creating shorter paths

# Bidirectional RNNs, Recursive Nets, Multidimensional RNNs, etc.

- The unfolded architecture needs not be a straight chain

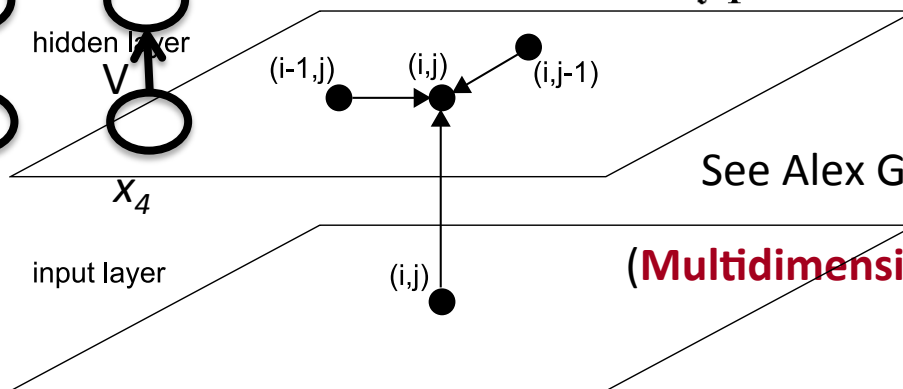


## Bidirectional RNNs (Schuster and Paliwal, 1997)



See Alex Graves's work, e.g., 2012

## (Multidimensional RNNs, Graves et al 2007)



# Multiplicative Interactions

(Wu et al, 2016, arXiv:1606.06630)

- Multiplicative Integration RNNs:

- Replace

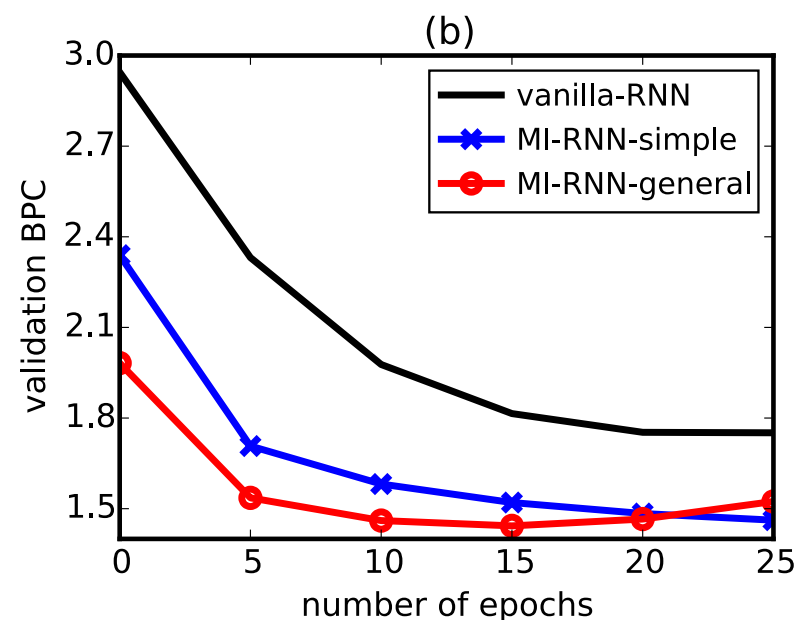
$$\phi(\mathbf{W}\mathbf{x} + \mathbf{U}\mathbf{z} + \mathbf{b})$$

- By

$$\phi(\mathbf{W}\mathbf{x} \odot \mathbf{U}\mathbf{z} + \mathbf{b})$$

- Or more general:

$$\phi(\alpha \odot \mathbf{W}\mathbf{x} \odot \mathbf{U}\mathbf{z} + \beta_1 \odot \mathbf{U}\mathbf{z} + \beta_2 \odot \mathbf{W}\mathbf{x} + \mathbf{b})$$



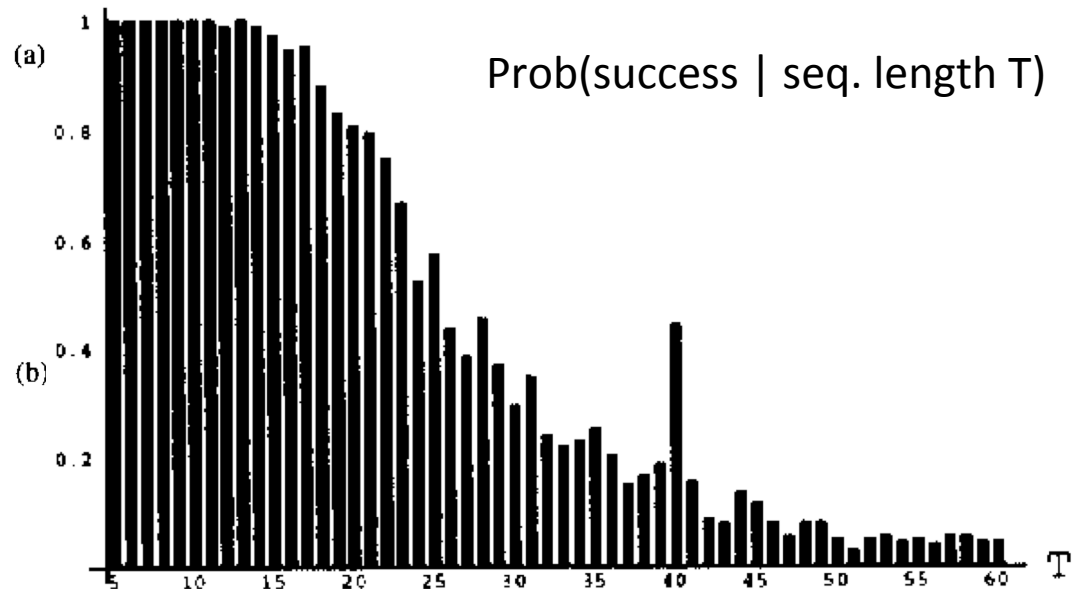
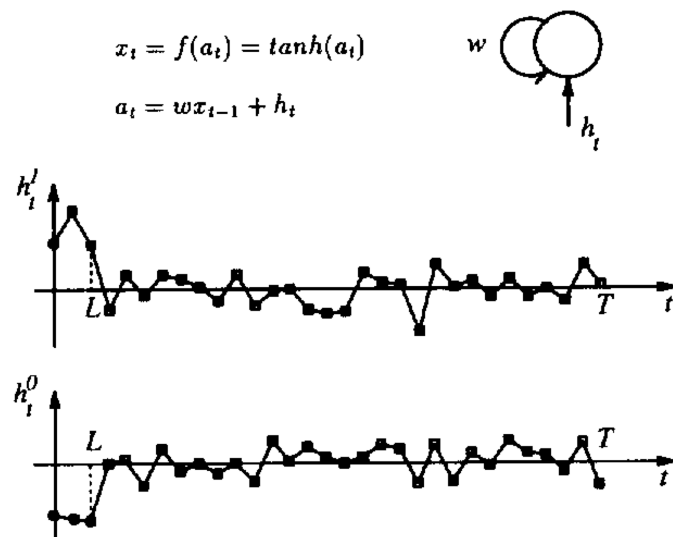
# Learning Long-Term Dependencies with Gradient Descent is Difficult



Y. Bengio, P. Simard & P. Frasconi, IEEE Trans. Neural Nets, **1994**

# Simple Experiments from 1991 while I was at MIT

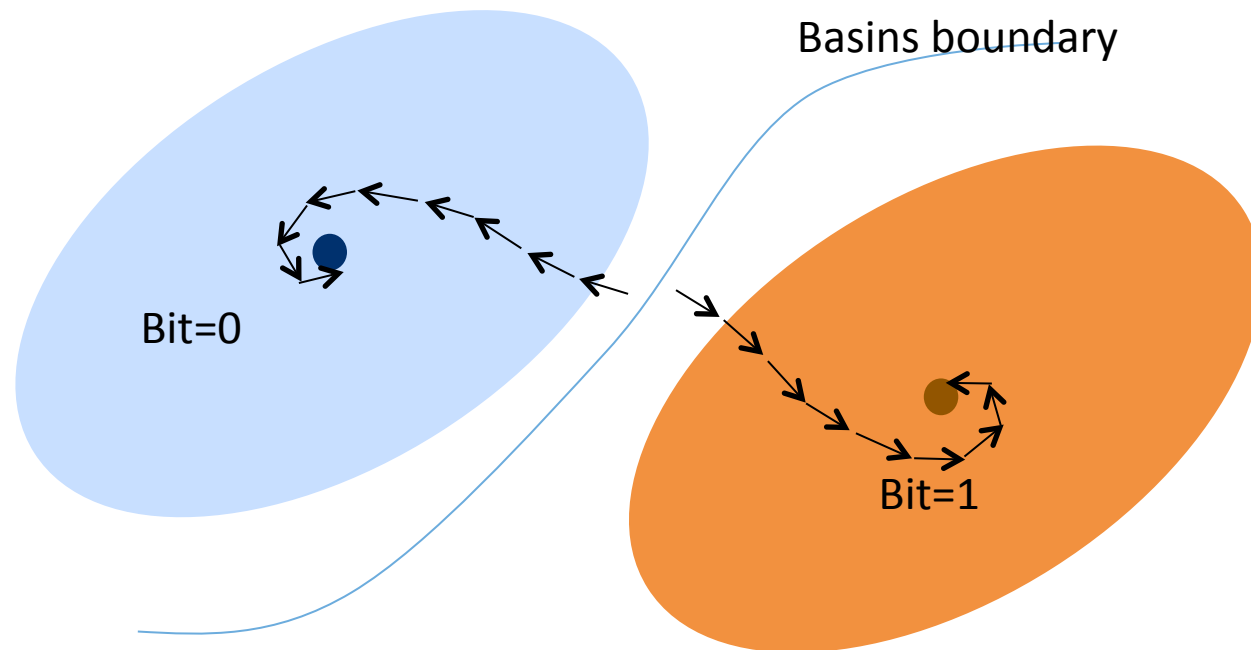
- 2 categories of sequences
- Can the single tanh unit learn to store for  $T$  time steps 1 bit of information given by the sign of initial input?





# How to store 1 bit? Dynamics with multiple basins of attraction in some dimensions

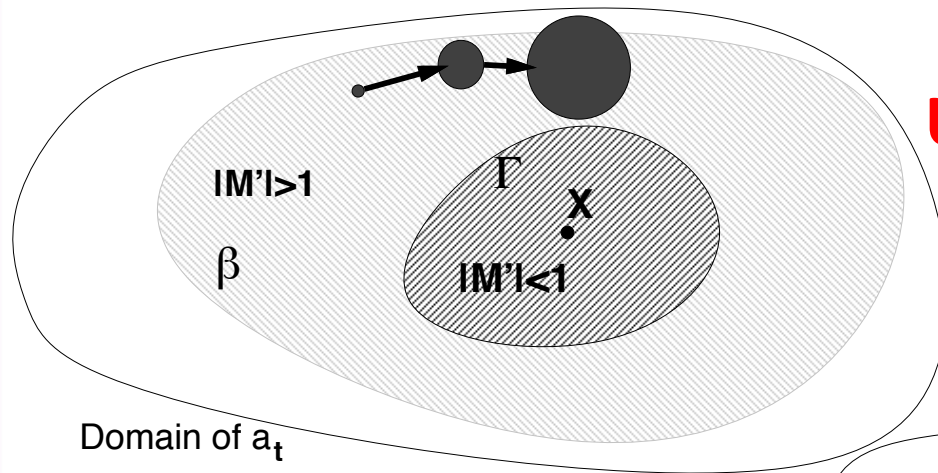
- Some subspace of the state can store 1 or more bits of information if the dynamical system has multiple basins of attraction in some dimensions



Note: gradients MUST be high near the boundary

# Robustly storing 1 bit in the presence of bounded noise

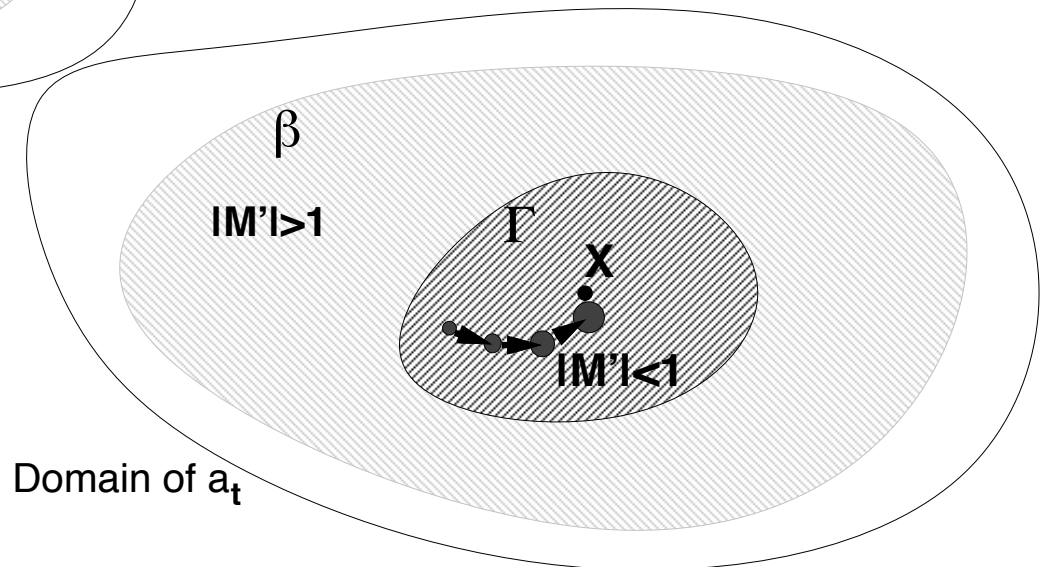
- With spectral radius  $> 1$ , noise can kick state out of attractor



**UNSTABLE**

- Not so with radius  $< 1$

**CONTRACTIVE**  
**→ STABLE**



## Storing Reliably $\rightarrow$ Vanishing gradients

- Reliably storing bits of information requires spectral radius  $< 1$
- The product of  $T$  matrices whose spectral radius is  $< 1$  is a matrix whose spectral radius converges to 0 at exponential rate in  $T$

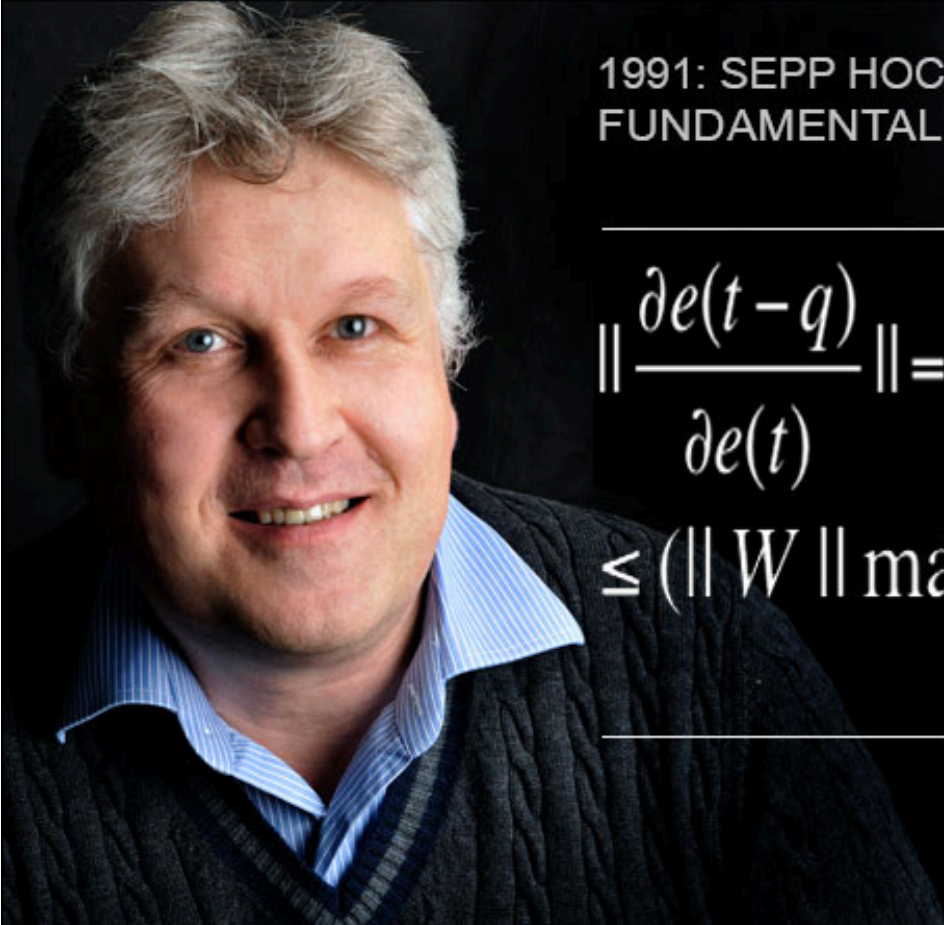
$$L = L(s_T(s_{T-1}(\dots s_{t+1}(s_t, \dots))))$$

$$\frac{\partial L}{\partial s_t} = \frac{\partial L}{\partial s_T} \frac{\partial s_T}{\partial s_{T-1}} \dots \frac{\partial s_{t+1}}{\partial s_t}$$

- If spectral radius of Jacobian is  $< 1 \rightarrow$  propagated gradients vanish

# Vanishing or Exploding Gradients

- Hochreiter's 1991 MSc thesis (in German) had independently discovered that backpropagated gradients in RNNs tend to either vanish or explode as sequence length increases

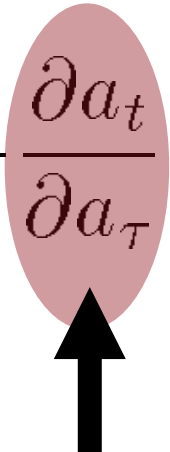


1991: SEPP HOCHREITER'S ANALYSIS OF THE FUNDAMENTAL DEEP LEARNING PROBLEM

$$\left\| \frac{\partial e(t-q)}{\partial e(t)} \right\| = \left\| \prod_{m=1}^q W F'(Net(t-m)) \right\|$$
$$\leq (\|W\| \max_{Net} \{ \|F'(Net)\| \})^q$$

## Why it hurts gradient-based Learning

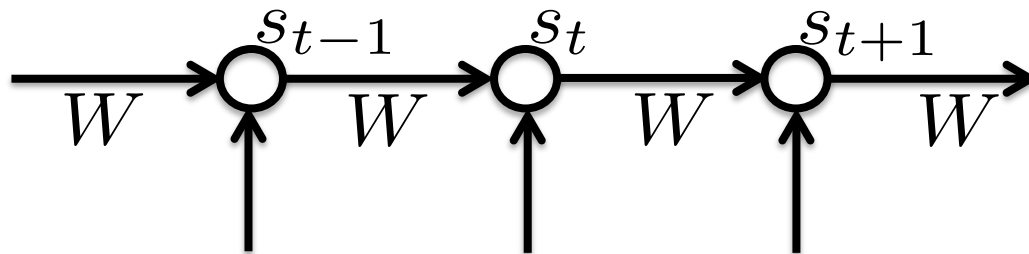
- Long-term dependencies get a weight that is exponentially smaller (in T) compared to short-term dependencies

$$\frac{\partial C_t}{\partial W} = \sum_{\tau \leq t} \frac{\partial C_t}{\partial a_\tau} \frac{\partial a_\tau}{\partial W} = \sum_{\tau \leq t} \frac{\partial C_t}{\partial a_t} \frac{\partial a_t}{\partial a_\tau} \frac{\partial a_\tau}{\partial W}$$


Becomes exponentially smaller  
for longer time differences,  
when spectral radius < 1

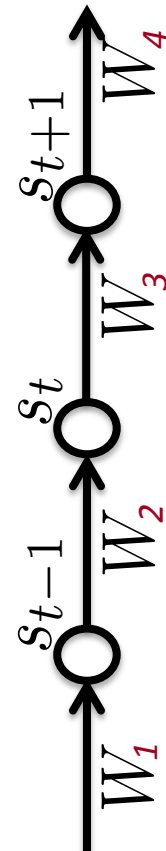
# Vanishing Gradients in Deep Nets are Different from the Case in RNNs

- If it was just a case of vanishing gradients in deep nets, we could just rescale the per-layer learning rate, but that does not really fix the training difficulties.



- Can't do that with RNNs because the weights are shared, & total true gradient = sum over different "depths"

$$\frac{\partial C_t}{\partial W} = \sum_{\tau \leq t} \frac{\partial C_t}{\partial a_\tau} \frac{\partial a_\tau}{\partial W} = \sum_{\tau \leq t} \frac{\partial C_t}{\partial a_t} \frac{\partial a_t}{\partial a_\tau} \frac{\partial a_\tau}{\partial W}$$



# To store information robustly the dynamics must be contractive

- The RNN gradient is a product of Jacobian matrices, each associated with a step in the forward computation. To store information robustly in a finite-dimensional state, the dynamics must be contractive [Bengio et al 1994].

$$L = L(s_T(s_{T-1}(\dots s_{t+1}(s_t, \dots))))$$
$$\frac{\partial L}{\partial s_t} = \frac{\partial L}{\partial s_T} \frac{\partial s_T}{\partial s_{T-1}} \dots \frac{\partial s_{t+1}}{\partial s_t}$$

Storing bits robustly requires e-values < 1

- Problems:
    - e-values of Jacobians > 1 → *gradients explode*
    - or e-values < 1 → *gradients shrink & vanish*
    - or random → *variance grows exponentially*
- Gradient clipping

# RNN Tricks

(Pascanu, Mikolov, Bengio, ICML 2013; Bengio, Boulanger & Pascanu, ICASSP 2013)

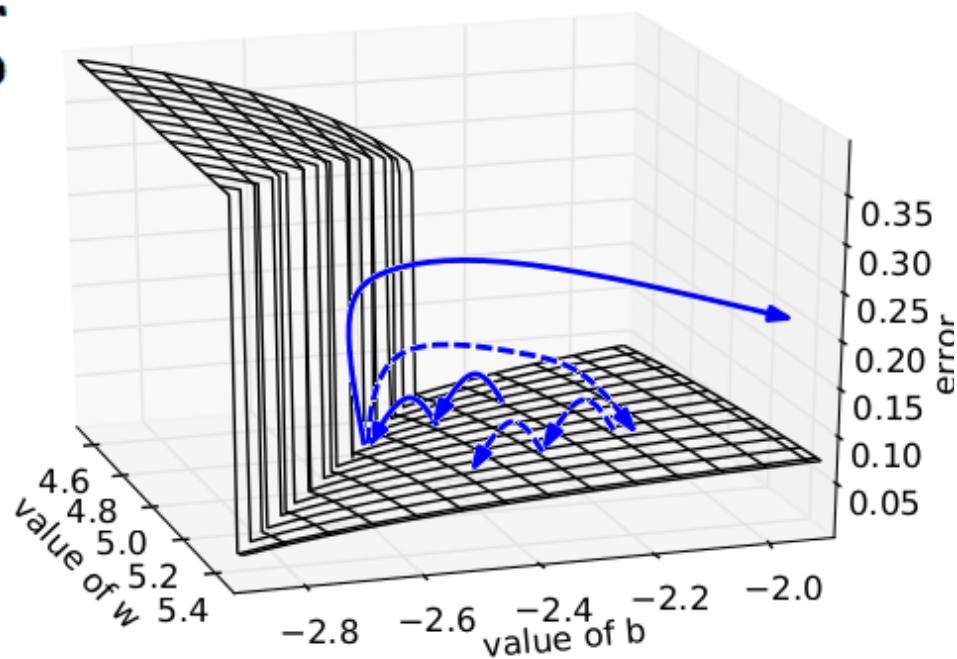
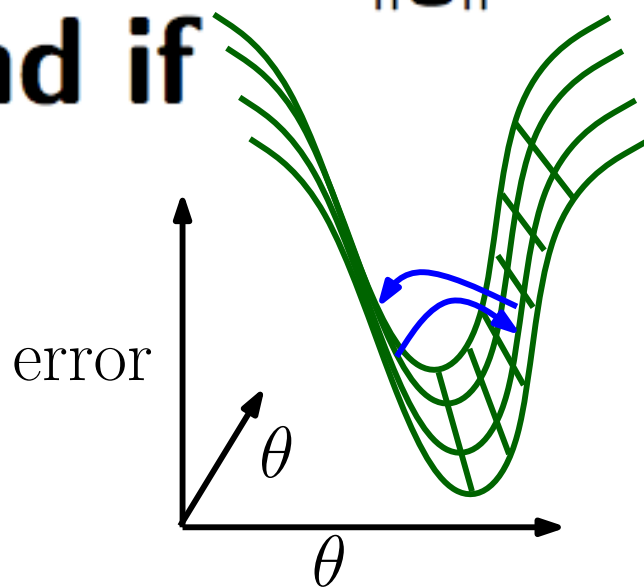
- Clipping gradients (avoid exploding gradients)
- Leaky integration (propagate long-term dependencies)
- Momentum (cheap 2<sup>nd</sup> order)
- Initialization (start in right ballpark avoids exploding/vanishing)
- Sparse Gradients (symmetry breaking)
- Gradient propagation regularizer (avoid vanishing gradient)
- Gated self-loops (LSTM & GRU, reduces vanishing gradient)



# Dealing with Gradient Explosion by Gradient Norm Clipping

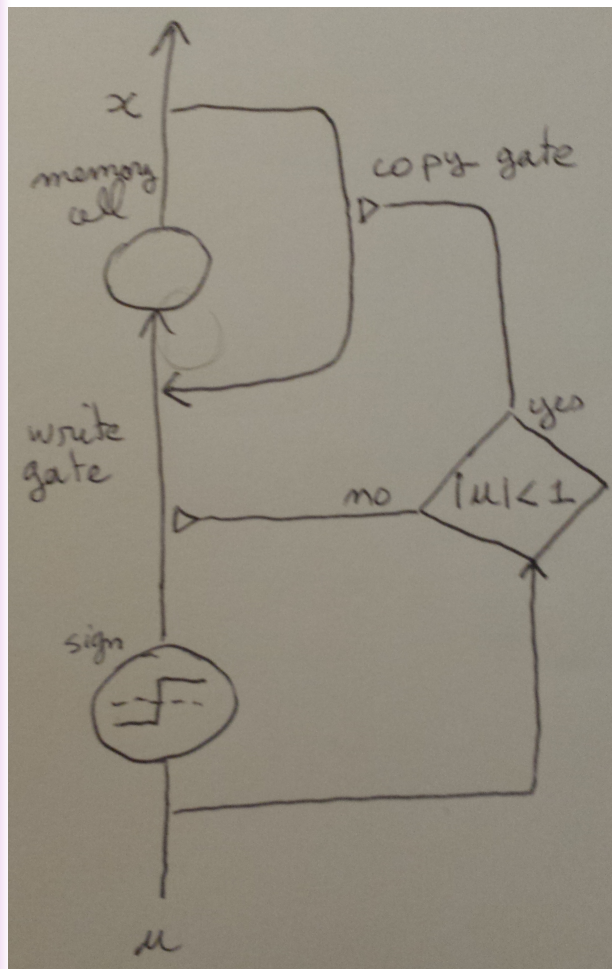
(Mikolov thesis 2012;  
Pascanu, Mikolov, Bengio, ICML 2013)

$\hat{\mathbf{g}} \leftarrow \frac{\partial \text{error}}{\partial \theta}$   
**if**  $\|\hat{\mathbf{g}}\| \geq \text{threshold}$  **then**  
     $\hat{\mathbf{g}} \leftarrow \frac{\text{threshold}}{\|\hat{\mathbf{g}}\|} \hat{\mathbf{g}}$   
**end if**



# Conference version (1993) of the 1994 paper by the same authors had a predecessor of GRU and targetprop

*(The problem of learning long-term dependencies in recurrent networks, Bengio, Frasconi & Simard ICNN'1993)*



## IV. A TRAINABLE FLIP-FLOP

- Flip-flop unit to store 1 bit, with gating signal to control when to write

$$x_{t+1} = f(x_t, u_t)$$

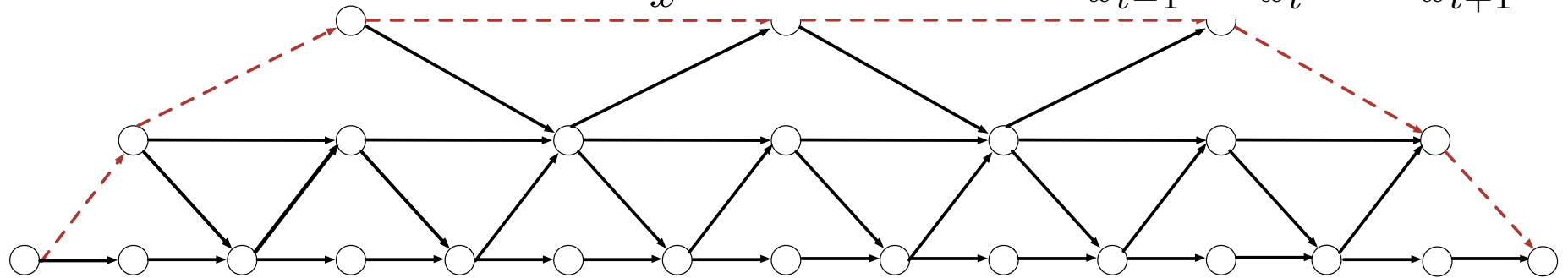
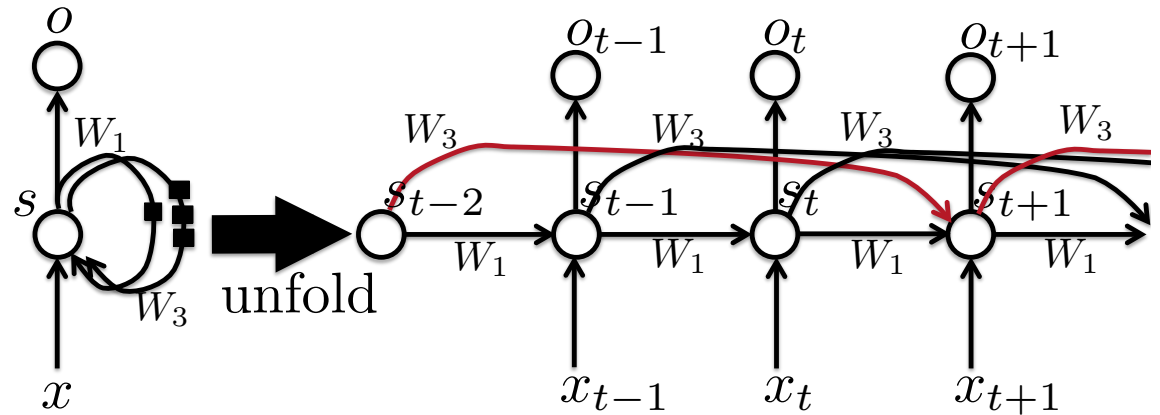
$$f(x, u) = \begin{cases} 1 & \text{if } |u| < 1 \text{ and } x \geq 0 \\ & \text{or if } u \geq 1 \\ -1 & \text{otherwise} \end{cases} \quad (8)$$

- Pseudo-backprop through it by a form of targetprop

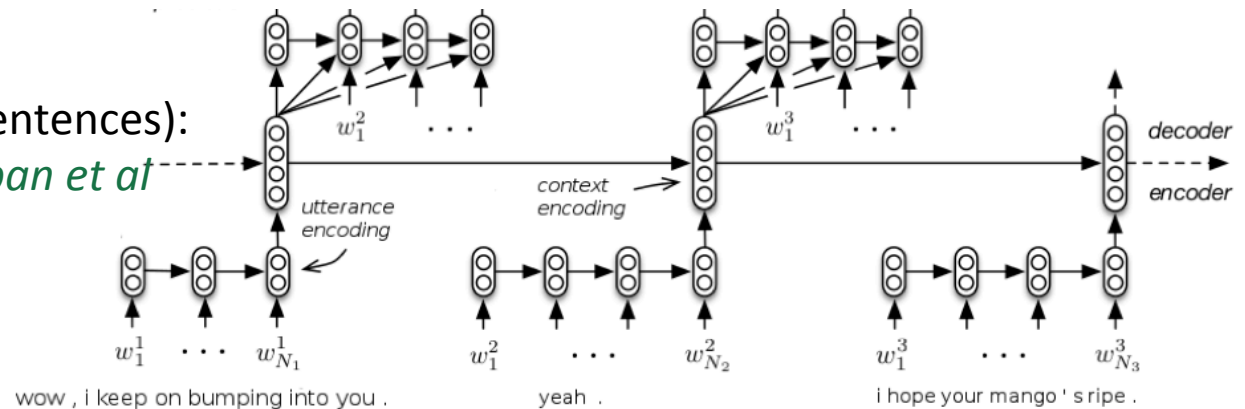
$$\Delta x(\Delta f, u) = \begin{cases} \Delta f & \text{if } |u| < 1 \\ 0 & \text{otherwise} \end{cases} \quad (11)$$

# Delays & Hierarchies to Reach Farther

- Delays and multiple time scales, *Elhihi & Bengio NIPS 1995*, *Koutnik et al ICML 2014*
- *How to do this right?*
- *How to automatically and adaptively do it?*



Hierarchical RNNs (words / sentences):  
*Sordani et al CIKM 2015, Serban et al AACL 2016*

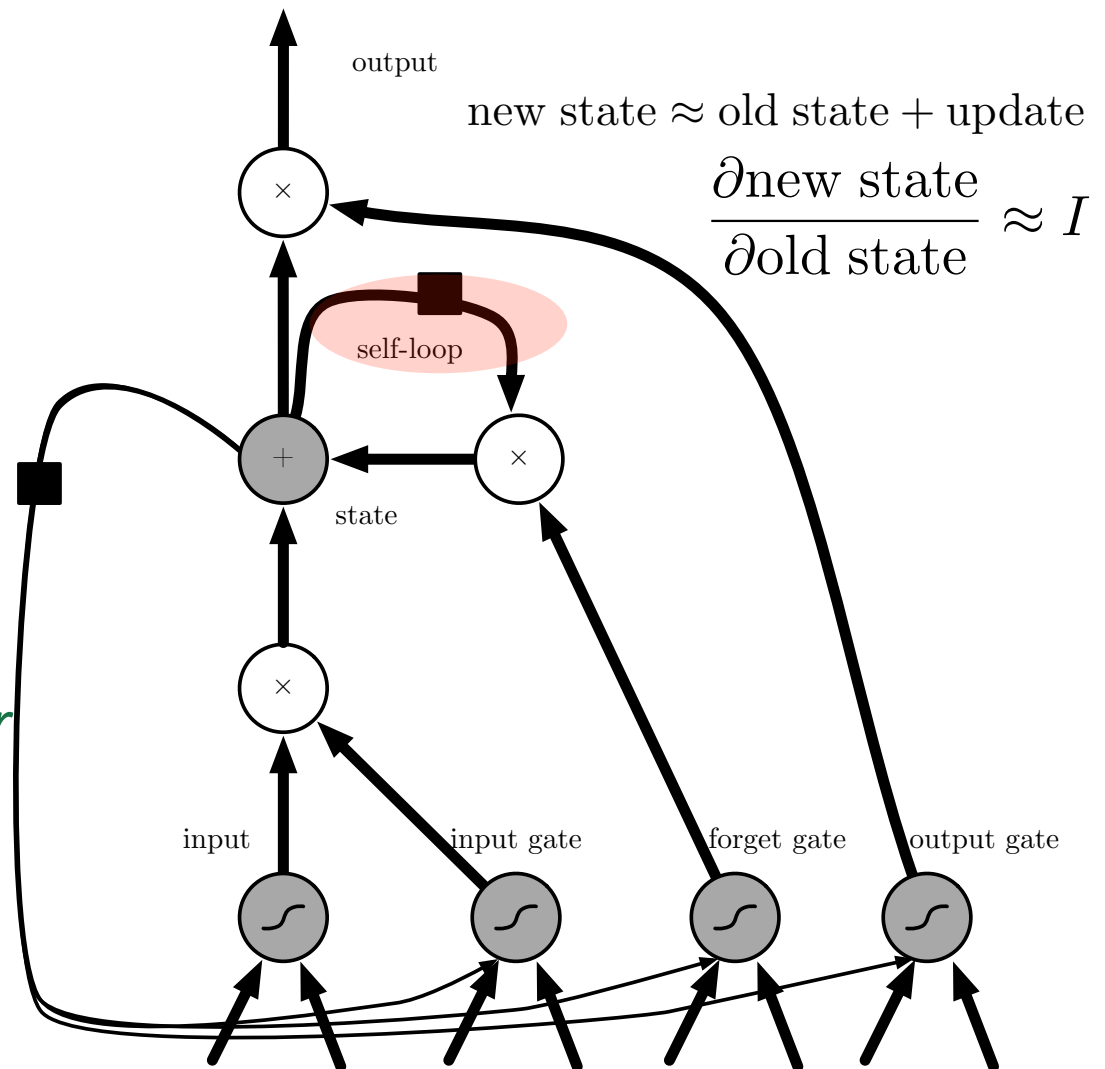


# Fighting the vanishing gradient: LSTM & GRU

(Hochreiter 1991); first version of the LSTM, called Neural Long-Term Storage with self-loop

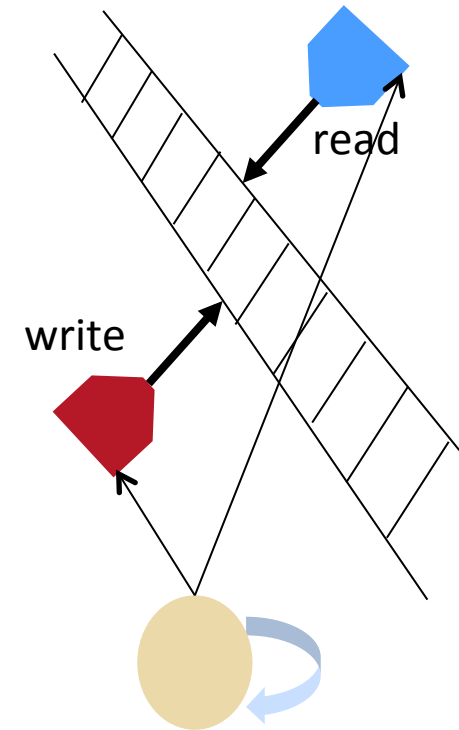
- Create a path where gradients can flow for longer with a self-loop
- Corresponds to an eigenvalue of Jacobian slightly less than 1
- LSTM is now **heavily used** (Hochreiter & Schmidhuber 1997)
- GRU light-weight version (Cho et al 2014)

LSTM: (Hochreiter & Schmidhuber 1997)



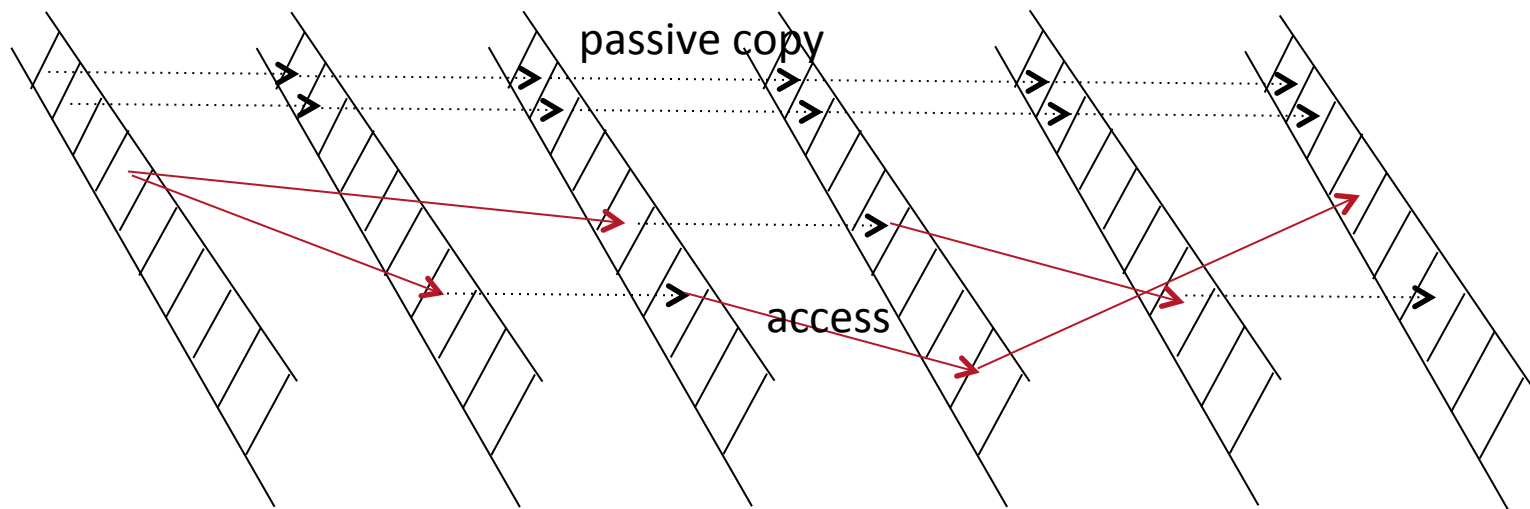
# Fast Forward 20 years: Attention Mechanisms for Memory Access

- Neural Turing Machines (*Graves et al 2014*)
- and Memory Networks (*Weston et al 2014*)
- Use a content-based attention mechanism (*Bahdanau et al 2014*) to control the read and write access into a memory
- The attention mechanism outputs a softmax over memory locations



# Large Memory Networks: Sparse Access Memory for Long-Term Dependencies

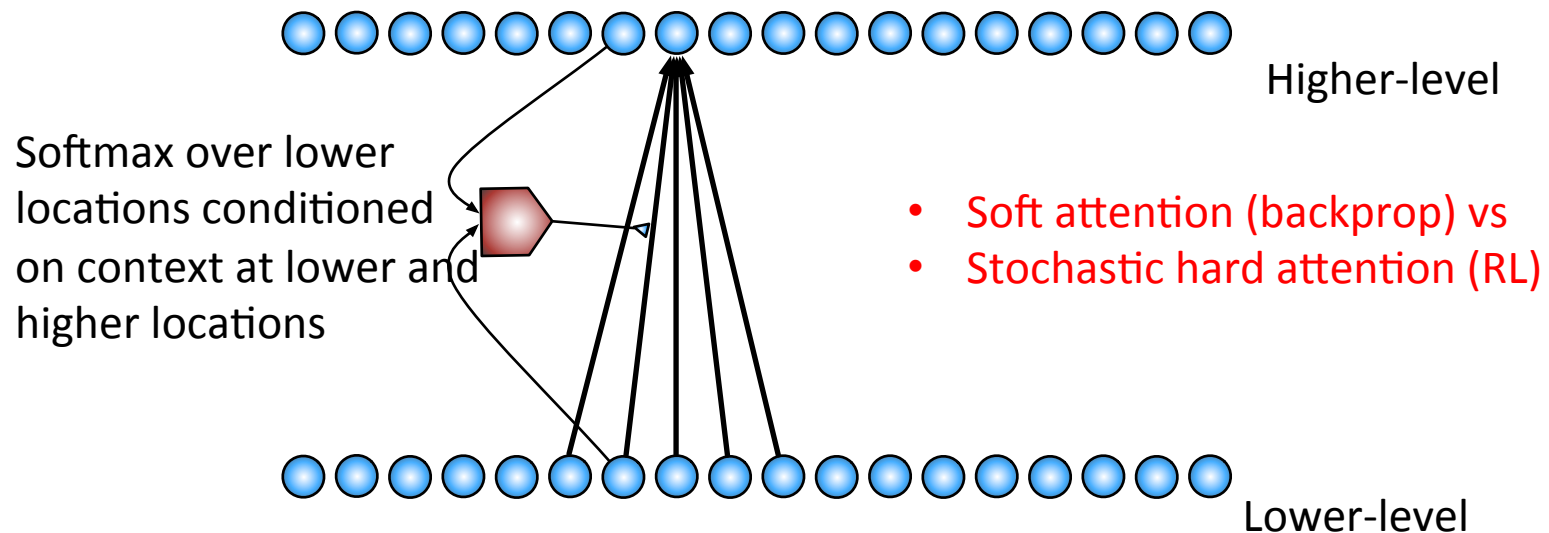
- Memory = part of the state
- Memory-based networks are special RNNs
- A mental state stored in an external memory can stay for arbitrarily long durations, until it is overwritten (partially or not)
- Forgetting = vanishing gradient.
- Memory = **higher-dimensional state**, avoiding or reducing the need for forgetting/vanishing



# Attention Mechanism for Deep Learning

(Bahdanau, Cho & Bengio, ICLR 2015; Jean et al ACL 2015; Jean et al WMT 2015; Xu et al ICML 2015; Chorowski et al NIPS 2015; Firat, Cho & Bengio 2016)

- Consider an input (or intermediate) sequence or image
- Consider an upper level representation, which can choose « where to look », by assigning a weight or probability to each input position, as produced by an MLP, applied at each position



# End-to-End Machine Translation with Recurrent Nets and Attention Mechanism

(Bahdanau et al ICLR 2015, Jean et al ACL 2015, Gulcehre et al 2015, Firat et al 2016)

- Reached the state-of-the-art in one year, from scratch

(a) English→French (WMT-14)

	<b>NMT(A)</b>	Google	P-SMT
NMT	32.68	30.6*	<b>37.03°</b>
+Cand	33.28	–	
+UNK	33.99	32.7°	
+Ens	<b>36.71</b>	<b>36.9°</b>	

(b) English→German (WMT-15)

Model	Note
<b>24.8</b>	Neural MT
24.0	U.Edinburgh, Syntactic SMT
23.6	LIMSI/KIT
22.8	U.Edinburgh, Phrase SMT
22.7	KIT, Phrase SMT

(c) English→Czech (WMT-15)

Model	Note
<b>18.3</b>	Neural MT
18.2	JHU, SMT+LM+OSM+Sparse
17.6	CU, Phrase SMT
17.4	U.Edinburgh, Phrase SMT
16.1	U.Edinburgh, Syntactic SMT



# Google-Scale NMT Success

(Wu et al & Dean, Nature, 2016)

- After beating the classical phrase-based MT on the academic benchmarks, there remained the question: will it work on the very large scale datasets like used for Google Translate?
- Distributed training, very large model ensemble
- Not only does it work in terms of BLEU but it makes a killing in terms of human evaluation on Google Translate data

Table 10: Side-by-side scores on production data

	PBMT	GNMT	Human	Relative Improvement
English → Spanish	3.594±1.58	5.031±1.09	5.140±1.04	93%
English → French	3.518±1.70	5.032±1.22	5.215±1.03	89%
English → Portuguese	3.675±1.64	4.856±1.29	4.973±1.17	91%
English → Chinese	2.457±1.48	4.154±1.42	4.580±1.26	80%
Spanish → English	3.410±1.65	4.921±1.16	4.930±1.12	99%
French → English	3.639±1.63	5.000±1.07	5.016±1.09	99%
Portuguese → English	3.471±1.74	5.029±1.05	5.040±1.03	99%
Chinese → English	1.994±1.47	3.884±1.37	4.334±1.20	81%

# Pointing the Unknown Words

Gulcehre, Ahn, Nallapati, Zhou & Bengio ACL 2016

Based on 'Pointer Networks', Vinyals et al 2015

The next word generated can either come from vocabulary or is copied from the input sequence.

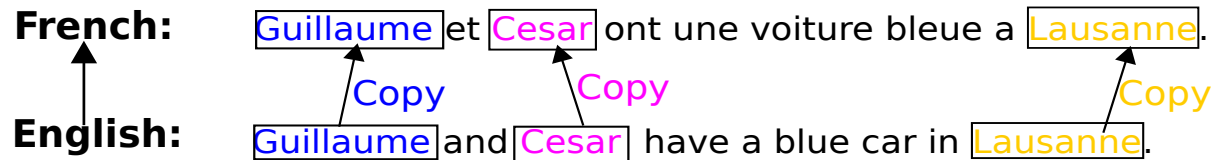


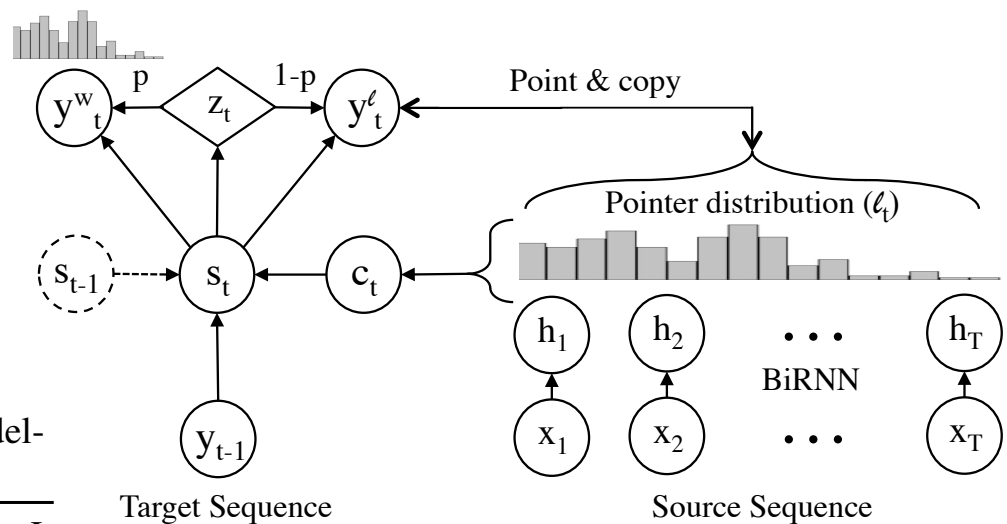
Table 5: Europarl Dataset (EN-FR)

Machine Translation	BLEU-4	
	NMT	20.19
NMT + PS	<b>23.76</b>	

Table 3: Results on Gigaword Corpus for modeling UNK's with pointers in terms of recall.

	Rouge-1	Rouge-2	Rouge-L
NMT + lvt	36.45	17.41	33.90
NMT + lvt + PS	<b>37.29</b>	<b>17.75</b>	<b>34.70</b>

Vocabulary softmax

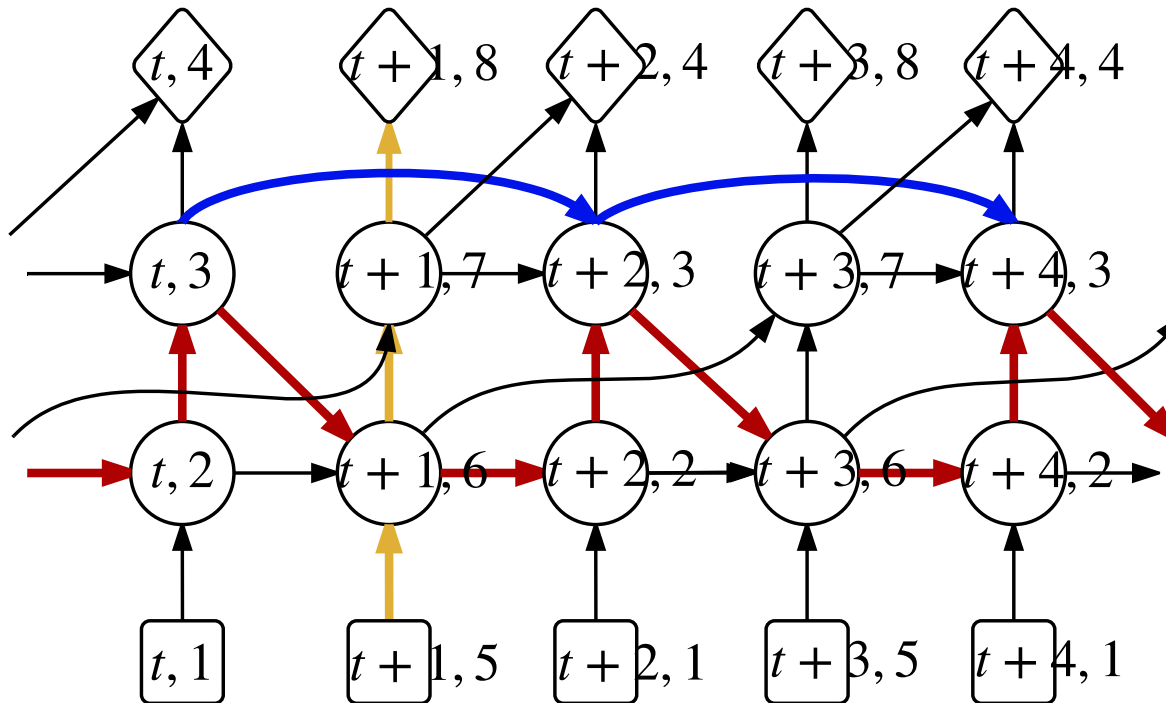


Text summarization

# Designing the RNN Architecture

(Architectural Complexity Measures of Recurrent Neural Networks  
Zhang et al 2016, arXiv:1602.08210)

- **Recurrent depth**: max path length divided by sequence length
- **Feedforward depth**: max length from input to nearest output
- **Skip coefficient**: shortest path length divided sequence length



# It makes a difference

- Impact of change in recurrent depth

DATASET	MODELS\ARCHS	<i>sh</i>	<i>st</i>	<i>bu</i>	<i>td</i>
<i>PennTreebank</i>	<i>tanh</i> RNN	1.54	1.59	1.54	<b>1.49</b>
<i>text8</i>	<i>tanh</i> RNN-SMALL	1.80	1.82	1.80	<b>1.77</b>
	<i>tanh</i> RNN-LARGE	1.69	1.67	1.64	<b>1.59</b>
	LSTM-SMALL	1.65	1.66	1.65	<b>1.63</b>
	LSTM-LARGE	1.52	1.53	1.52	<b>1.49</b>

- Impact of change in skip coefficient

RNN( <i>tanh</i> )	<i>s</i> = 1	<i>s</i> = 5	<i>s</i> = 9	<i>s</i> = 13	<i>s</i> = 21
MNIST	34.9	46.9	74.9	85.4	<b>87.8</b>
<i>p</i> MNIST	49.8	79.1	84.3	<b>88.9</b>	88.0

LSTM	<i>s</i> = 1	<i>s</i> = 3	<i>s</i> = 5	<i>s</i> = 7	<i>s</i> = 9
MNIST	56.2	<b>87.2</b>	86.4	86.4	84.8
<i>p</i> MNIST	28.5	25.0	60.8	62.2	<b>65.9</b>

Model	MNIST	<i>p</i> MNIST
<i>i</i> RNN[25]	97.0	≈82.0
<i>u</i> RNN[24]	95.1	91.4
LSTM[24]	<b>98.2</b>	88.0
RNN( <i>tanh</i> )[25]	≈35.0	≈35.0
<i>stanh</i> ( <i>s</i> = 21, 11)	98.1	<b>94.0</b>

Architecture, <i>s</i>	(1), 1	(2), 1	(3), $\frac{k}{2}$	(4), <i>k</i>
MNIST <i>k</i> = 17	39.5	39.4	54.2	<b>77.8</b>
<i>k</i> = 21	39.5	39.9	69.6	<b>71.8</b>
<i>p</i> MNIST <i>k</i> = 5	55.5	66.6	74.7	<b>81.2</b>
<i>k</i> = 9	55.5	71.1	78.6	<b>86.9</b>

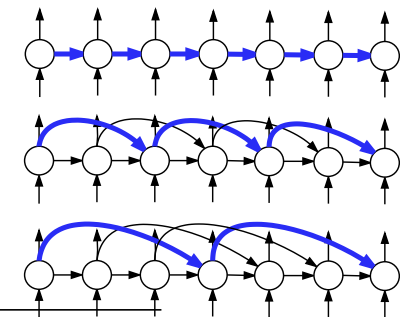
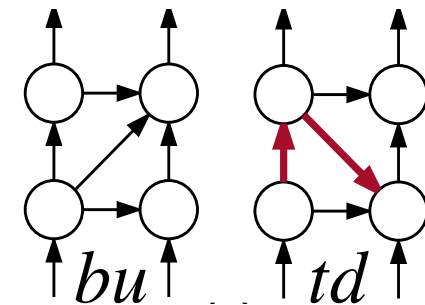


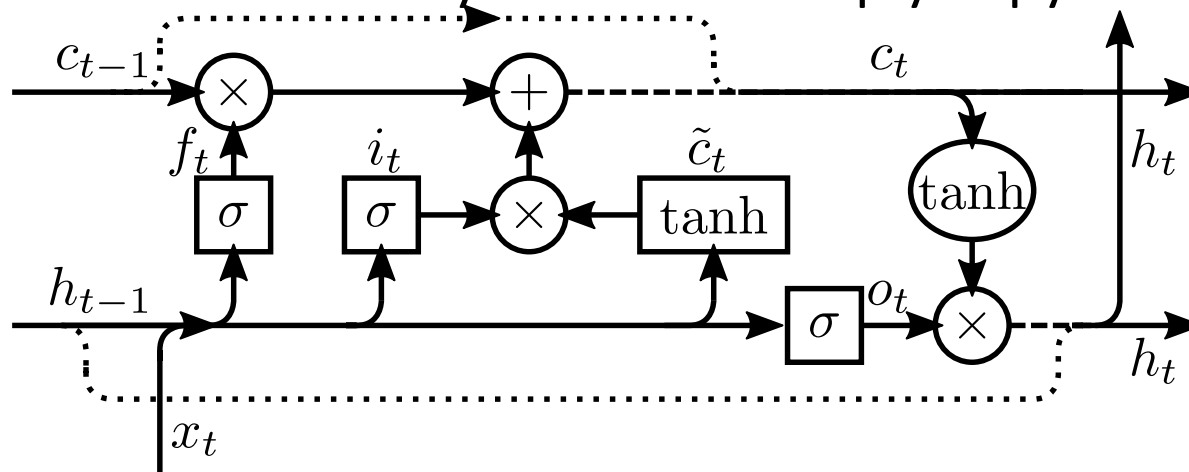
Table 2: Results for MNIST/*p*MNIST. **Top-left**: test accuracies with different *s* for *tanh* RNN. **Top-right**: test accuracies with different *s* for LSTM. **Bottom**: compared to previous results. **Bottom-right**: test accuracies for architectures (1), (2), (3) and (4) for *tanh* RNN.

# Near-Orthogonality to Help Information Propagation

- Initialization to orthogonal recurrent  $W$  *(Saxe et al 2013, ICLR2014)*
- Unitary matrices: all e-values of matrix are 1 *(Arjowski, Amar & Bengio ICML 2016)*

$$W = D_3 R_2 \mathcal{F}^{-1} D_2 \Pi R_1 \mathcal{F} D_1$$

- Zoneout: randomly choose to simply copy the state unchanged



*(Krueger et al 2016, submitted)*

# Variational Generative RNNs

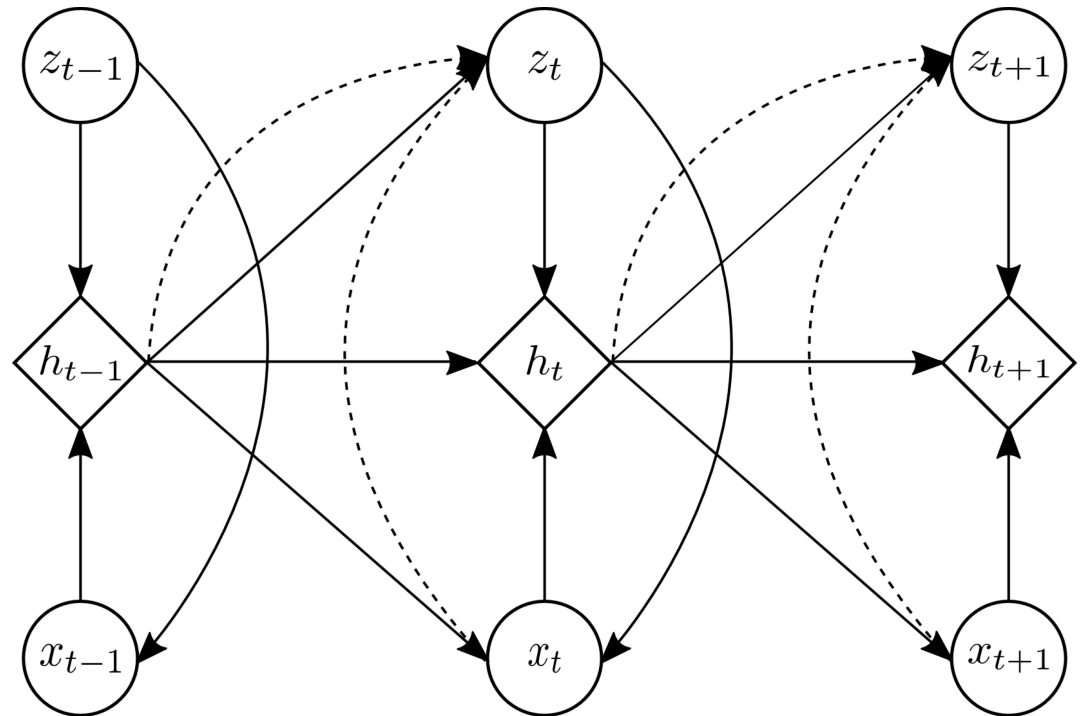
➔ Injecting higher-level variations / latent variables in RNNs

- (Chung et al, NIPS'2015)
- Regular RNNs have noise injected only in input space
- VRNNs also allow noise (latent variable) injected in top hidden layer; more « high-level » variability

Handwritten notes in Russian: "Вводятся шум в скрытый слой, что приводит к более высокой вариативности" (Noise is introduced into the hidden layer, which leads to higher variability).

Handwritten notes in Russian: "Вводятся шум в скрытый слой" (Noise is introduced into the hidden layer).

Handwritten notes in Russian: "Вводятся шум в скрытый слой, что приводит к более высокой вариативности" (Noise is introduced into the hidden layer, which leads to higher variability).



# Variational Hierarchical RNNs for Dialogue Generation (Serban et al 2016)

- Lower level = words of an utterance (turn of speech)
- Upper level = state of the dialogue
- Inject high-level choices

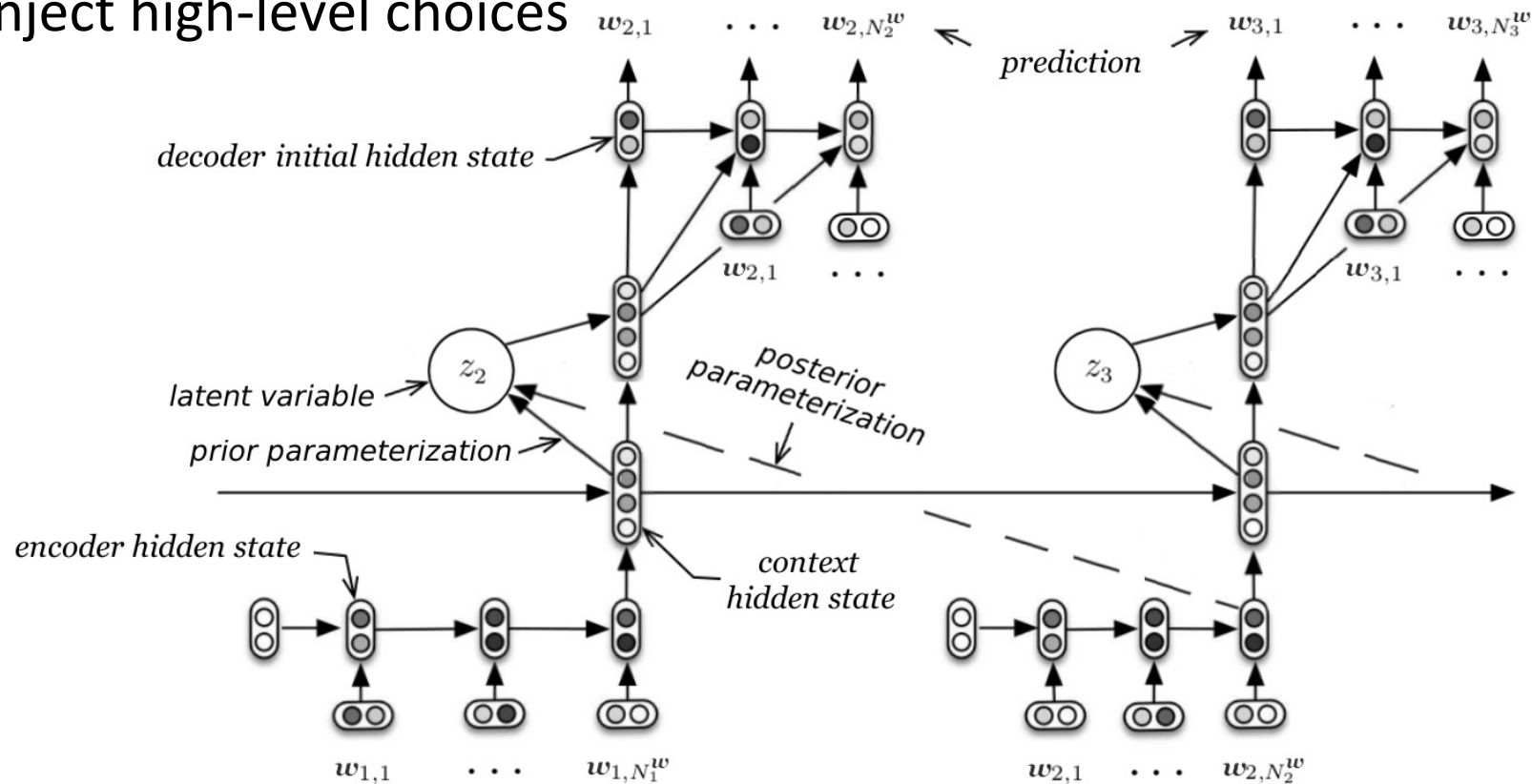


Table 1: Wins, losses and ties (in %) of the VHRED model against the baselines based on the human study on Twitter (mean preferences  $\pm$  90% confidence intervals)

Opponent	Short Contexts			Long Contexts		
	Wins	Losses	Ties	Wins	Losses	Ties
VHRED vs LSTM	32.3 $\pm$ 2.4	<b>42.5 <math>\pm</math> 2.6</b>	25.2 $\pm$ 2.3	<b>41.9 <math>\pm</math> 2.2</b>	36.8 $\pm$ 2.2	21.3 $\pm$ 1.9
VHRED vs HRED	<b>42.0 <math>\pm</math> 2.8</b>	31.9 $\pm$ 2.6	26.2 $\pm$ 2.5	<b>41.5 <math>\pm</math> 2.8</b>	29.4 $\pm$ 2.6	29.1 $\pm$ 2.6
VHRED vs TF-IDF	<b>51.6 <math>\pm</math> 3.3</b>	17.9 $\pm$ 2.5	30.4 $\pm$ 3.0	<b>47.9 <math>\pm</math> 3.4</b>	11.7 $\pm$ 2.2	40.3 $\pm$ 3.4

Table 2: Twitter examples for the neural network models. The  $\rightarrow$  token indicates a change of turn.

Context	Response
jamaica blood!! <url> $\rightarrow$ lol. yes [name]!. [name]! $\rightarrow$ what's poppin in toronto homie	<b>VHRED:</b> nothing much. just chillin. what's good with you? <b>LSTM:</b> something!!!!!!! <b>HRED:</b> something!!!!!!!
whhhhhhyyyy is it suddenly cold in tallahassee today?? this is florida! how can i dream of the beach when my fingers are icy cold?? $\rightarrow$ sittin on deck in sun looking at lake travis in austin! enjoying day b4 birthday	<b>VHRED:</b> the sun is shining!!! <b>LSTM:</b> i know!!!!!!! <b>HRED:</b> something!!!!!!!
reaches for [name]! $\rightarrow$ - approaches him warily - hey there buddy. - swings him into the air - how you doin? $\rightarrow$ tilts his head and your wariness and reaches for your face [name].	<b>VHRED:</b> - looks up at you and smiles - you're so cute. - kisses your cheek - <b>LSTM:</b> - smiles - i'm fine thank you. how are you? <b>HRED:</b> - grins, looking at you as i look at you -
i think i want an ipad2. but i already have an ibook. want vs. need. $\rightarrow$ if you have an ibook. then you deserve an ipad!	<b>VHRED:</b> that's what i'm thinking. i'm going to buy a new one. <b>LSTM:</b> i want one!!!! <b>HRED:</b> i don't know what to do with it. i don't know what to do.
ah echt? dan anrtwoord ik volgende keer via twitter als je me smst. $\rightarrow$ haha doe maar p. goed bezig hoor myrthe d $\rightarrow$ goed bezig? die snap ik even niet p	<b>VHRED:</b> haha, ik heb er geen zin in de something, maar ik heb er geen something meer something <b>LSTM:</b> haha something something something <b>HRED:</b> something something something something something something something something something
i love katie waissel more than like, anything in life. $\rightarrow$ i wanna c the show again! i need <number> meet her again! n i want a b-day tweet nxt week doubt that will happen x	<b>VHRED:</b> me too. i'm going to have to wait for her to come back. i'm so excited. <b>LSTM:</b> i'm sure she will. i'm sure she'll be fine. <b>HRED:</b> i'm sure she'll be fine, i'm sure she'll be fine.
my math teacher over here talkn bout nickels nd i'm tryna make big money i need to be in calculus $\rightarrow$ lol	<b>VHRED:</b> i need to go to the gym and get my hair done lol <b>LSTM:</b> lol <b>HRED:</b> how u been

VHRNN  
Results

-  
Twitter  
Dialogues

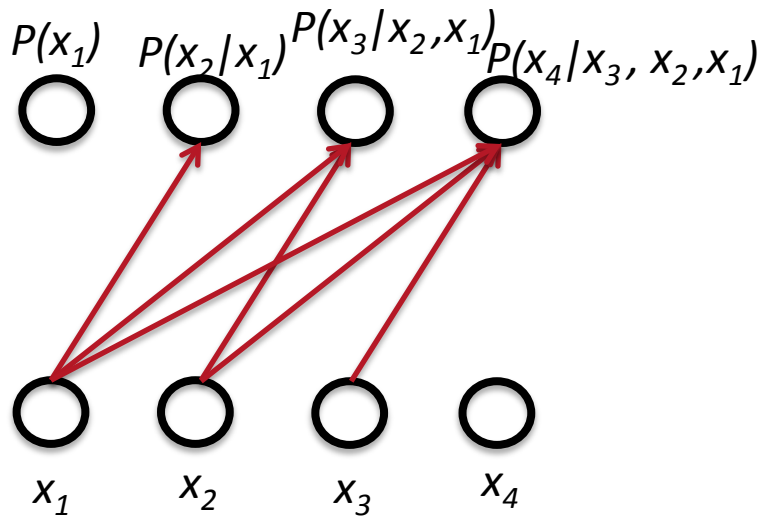


# Other Fully-Observed Neural Directed Graphical Models

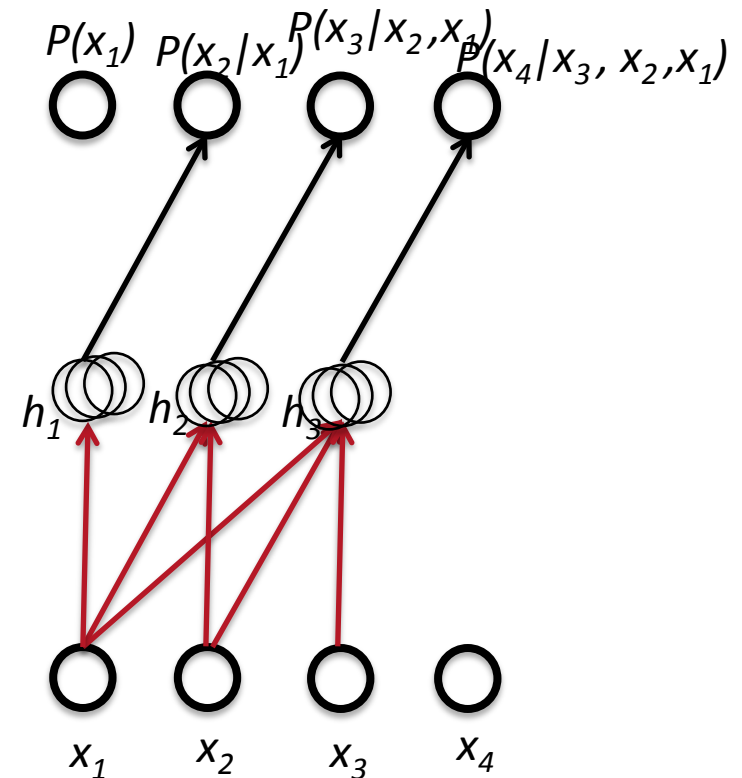
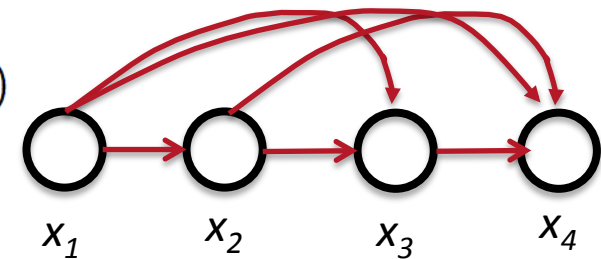
# Neural Auto-Regressive Models

$$P(\mathbf{x}) = P(x_1, \dots, x_T) = \prod_{t=1}^T P(x_t | x_{t-1}, x_{t-2}, \dots, x_1)$$

- Decomposes the joint of a fully observed directed model in terms of conditionals
- Logistic auto-regressive: (*Frey 1997*)



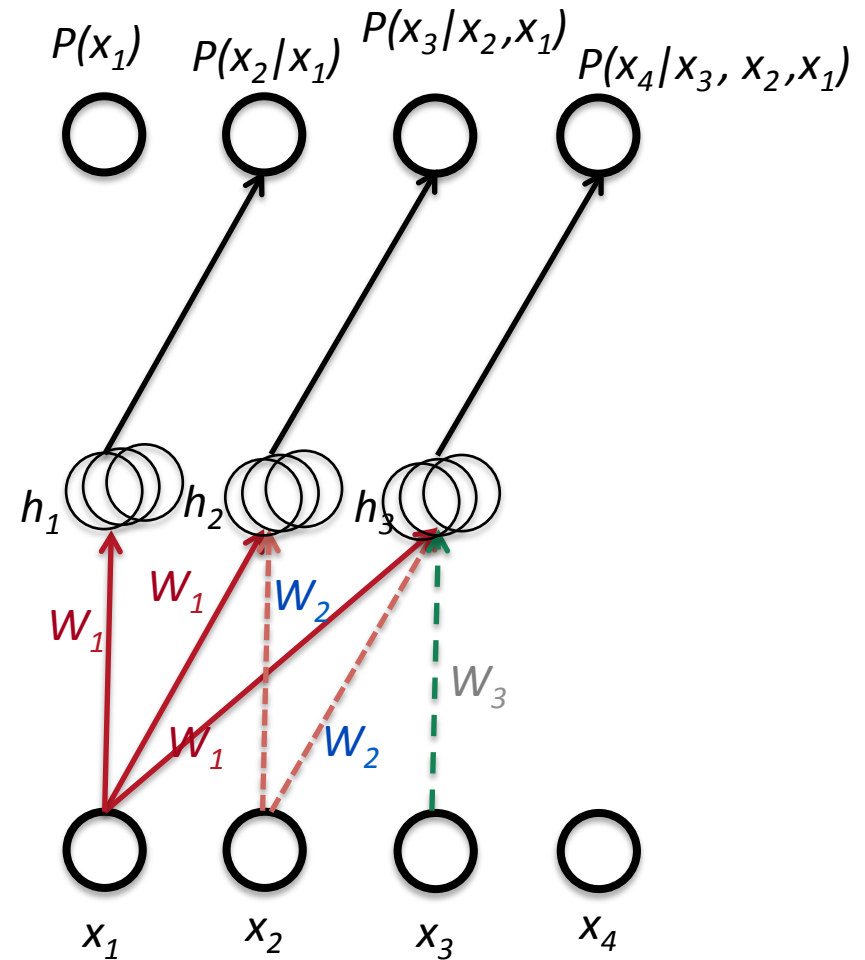
- First neural version: (*Bengio&Bengio NIPS'99*)



# NADE: Neural AutoRegressive Density Estimator

*(Larochelle & Murray AISTATS 2011)*

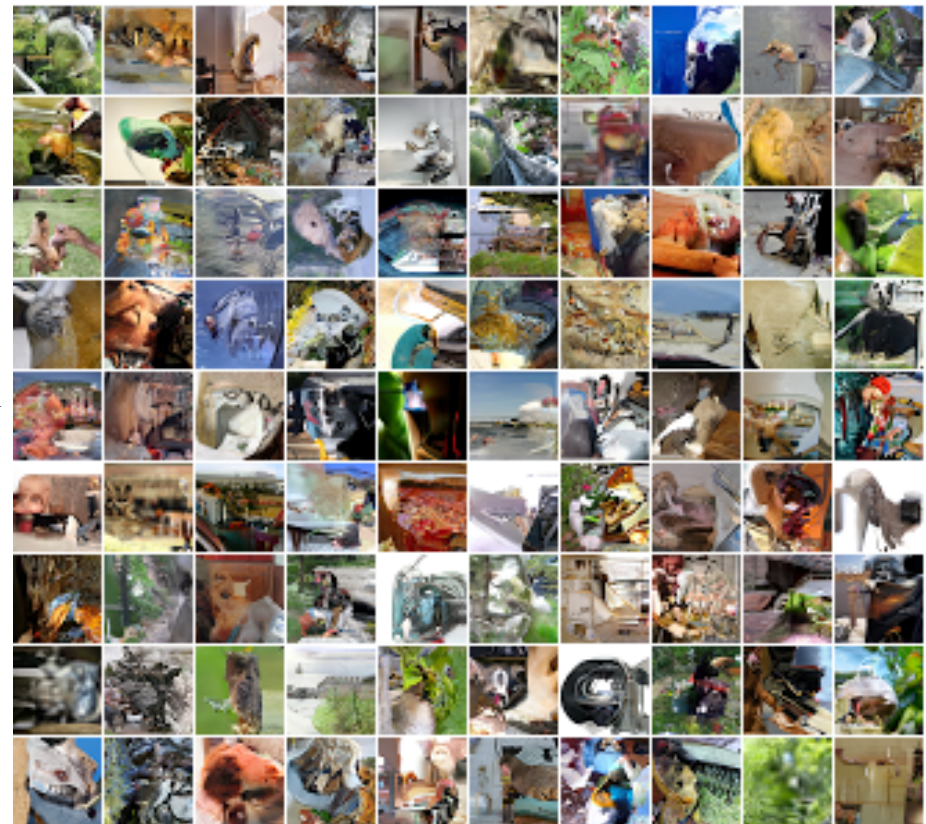
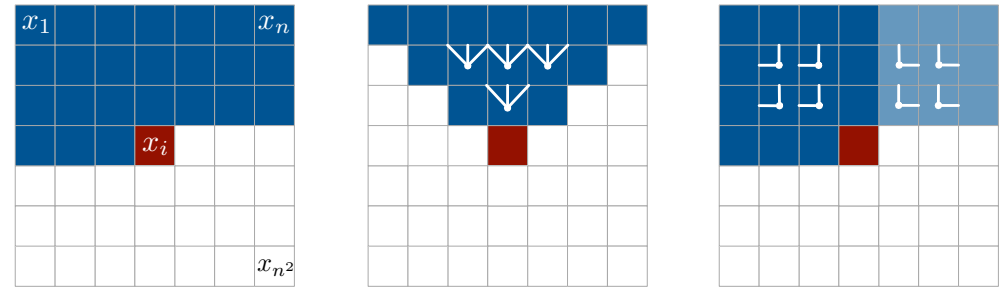
- Introduces smart sharing between some weights so that the different hidden groups use the same weights to the same input but look at more and more of the inputs.



# Pixel RNNs

(van den Oord et al ICML 2016, best paper)

- Similar to NADE and RNNs but for 2-D images
- Surprisingly sharp and realistic generation
- Gets texture right but not necessarily global structure



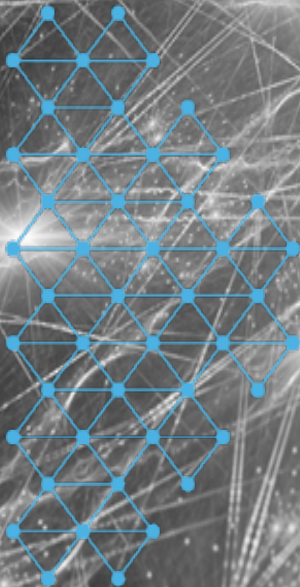
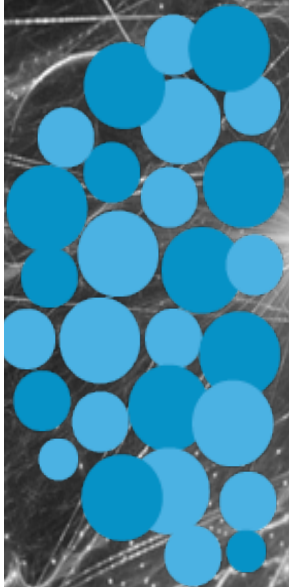
# Forward Computation of the Gradient

- BPTT does not seem biologically plausible and is memory-expensive
- RTRL (*Real-Time Recurrent Learning, Williams & Zipser 1989, Neural Comp.*)
  - Practically useful: online learning, no need to store all the past states and revisit history backwards (which is biologically weird)
  - Compute the gradients forward in time, rather than backwards
    - Think about multiplying many matrices left-to-right vs right-to-left
  - **BUT** exact computation is  $O(n_{\text{hidden}} \times n_{\text{weights}})$  instead of  $O(n_{\text{weights}})$ , to recursively compute  $dh(t)/dW \leftarrow$  all params
- Recently proposed, \*approximate\* the forward gradient using an efficient stochastic estimator (rank 1 estimator of  $dh/dW$  tensor) (*Training recurrent networks online without backtracking, Ollivier et al arXiv: 1507.07680*)





# Montreal Institute for Learning Algorithms



MILA

Université   
de Montréal