# The Challenge of Sequential Modelling

Stanford University, 26 Sept. 2016

**Yoshua Bengio Montreal Institute for Learning Algorithms** PLUG: **Deep Learning**, MIT Press book in pre-sale, PLUG: **Deep Learning**, MIT Press book in pre-sale, Chapters will remain online Université de Montréal



**NLA** 

## Recurrent Neural Networks

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

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

$$s_{t-1} \xrightarrow{F_{\theta}} \underbrace{s_{t+1}}_{f_{\theta} \text{ shared over time}} \xrightarrow{F_{\theta}} \underbrace{s_{t+1}}_{f_{\theta} \text{ shared over time}} \xrightarrow{F_{\theta}} \underbrace{s_{t+1}}_{x_{t-1} x_{t}} \xrightarrow{F_{\theta}} \underbrace{s_{t+1}}_{x_{t+1}}$$

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

Generalizes naturally to new lengths not seen during training

## Generative RNNs

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



# Increasing the Expressive Power of RNNs with more Depth

ICLR 2014, How to construct deep recurrent neural networks  $\bullet$ 



Z<sub>t-1</sub>

h<sub>t-1</sub>

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

hŧ





+ skip connections for creating shorter paths

+ stacking

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

• The unfolded architecture needs not be a straight chain



# Multiplicative Interactions 1

(Wu et al, 2016, arXiv: 1606.06630, NIPS2016)

- Multiplicative Integration RNNs:
  - Replace  $\phi(\mathbf{W} oldsymbol{x} + \mathbf{U} oldsymbol{z} + \mathbf{b})$
  - By  $\phi(\mathbf{W} oldsymbol{x} \odot \mathbf{U} oldsymbol{z} + \mathbf{b})$
  - Or more general:

 $\phi(oldsymbol{lpha}\odot {f W}oldsymbol{x}\odot {f U}oldsymbol{z}+oldsymbol{eta}_1\odot {f U}oldsymbol{z}+oldsymbol{eta}_2\odot {f W}oldsymbol{x}+oldsymbol{b}ig)^5$  number of epochs

Experiments on speech recognition, but also on language modeling and learning word embeddings.

WSJ Corpus	CER	WER
DRNN+CTCbeamsearch [17]	10.0	14.1
Encoder-Decoder [18]	6.4	9.3
LSTM+CTCbeamsearch [19]	9.2	8.7
Eesen [20]	-	7.3
LSTM+CTC+WFST (ours)	6.5	8.7
MI-LSTM+CTC+WFST (ours)	6.0	8.2





Learning Long-Term Dependencies with Gradient Descent is Difficult



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

# Robustly storing 1 bit in the presence of bounded noise

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



## Storing Reliably → Vanishing gradients

- Reliably storing bits of information requires spectral radius<1</li>
- 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

## Why it hurts gradient-based learning

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

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

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

# 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].

$$\begin{split} L &= L(s_T(s_{T-1}(\ldots s_{t+1}(s_t,\ldots)))))\\ \frac{\partial L}{\partial s_t} &= \frac{\partial L}{\partial s_T} \frac{\partial s_T}{\partial s_{T-1}} \cdots \frac{\partial s_{t+1}}{\partial s_t} & \text{Storing bits}\\ \text{robustly requires}\\ \text{e-values<1} \end{split}$$

- Problems:
  - e-values of Jacobians > 1  $\rightarrow$  gradients explode
  - or e-values  $< 1 \rightarrow gradients shrink \& vanish$
  - or random  $\rightarrow$  variance grows exponentially



## Delays & Hierarchies to Reach Farther

• Delays and multiple time scales, *Elhihi & Bengio NIPS 1995*, *Koutnik et al ICML 2014*  $\bigcirc^{O}$ 

. . .

wow, i keep on bumping into you.

 $w_N^1$ 

- How to do this right?
- How to automatically and adaptively do it?

Hierarchical RNNs (words / sentences): Sordoni et al CIKM 2015, Serban et al AAAI 2016



yeah .

 $w_{N_{\alpha}}^2$ 

i hope your mango ' s ripe

## 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



# What RNN architectures do you use in your research?





Skip connections Graves et al. 2013



Stacked RNN



"Deep" RNN Pascanu et al. 2013

ΜΙΙ Δ



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



# **Depth in RNNs**

- Recurrent depth :

the length of the \*longest\* path over time, divide by time  $\Im_{i}(n) = l(n)$ 

$$d_r = \lim_{n \to +\infty} \frac{\mathfrak{D}_i(n)}{n} = \max_{\vartheta \in C(\mathcal{G}_c)} \frac{l(\vartheta)}{\sigma_s(\vartheta)}$$

- Feedforward depth :

consider the input-output nonlinearities while eliminating the effect of recurrent depth

$$d_f = \sup_{i,n \in \mathbb{Z}} \mathfrak{D}_i^*(n) - n \cdot d_r$$



## 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	sh	st	bu	td
PennTreebank	tanh RNN	1.54	1.59	1.54	1.49
	tanh RNN-SMALL	1.80	1.82	1.80	1.77
text8	tanh RNN-large	1.69	1.67	1.64	1.59
	LSTM-SMALL	1.65	1.66	1.65	1.63
	LSTM-LARGE	1.52	1.53	1.52	1.49

#### Impact of change in skip coefficient





						_					
RNN(tanh) s	s = 1	s = 5	s = 9	s = 13	s = 21		LSTM	s = 1	s = 3 s	$= 5 s = 10^{-10}$	7  s = 9
MNIST	34.9	46.9	74.9	85.4	87.8	_	MNIST	56.2	<b>87.2</b> 8	86.4 86.4	4 84.8
S	s = 1	s = 3	s = 5	s = 7	s = 9			s = 1	s = 3 s	= 4 s = 1	$5 \ s = 6$
<i>p</i> MNIST	49.8	79.1	84.3	<b>88.9</b>	88.0		pMNIST	28.5	25.0 6	60.8 62.2	2 <b>65.9</b>
					_						
Model		MNIS	ST p	MNIST		rah	itaatura	(1)	1 (2) 1	(2) k	$(A)$ $l_{a}$
<i>i</i> RNN[25]		97.0	)	$\approx 82.0$			itecture, s	(1),	$\frac{1}{2}$ (2), 1	$(3), \frac{1}{2}$	$(4), \kappa$
u R N N [24]		95 1		<b>Q1</b> <i>A</i>	1	MNI	ST k = 17	39.5	39.4	54.2	77 <b>.8</b>
		00.1		)1. <del>T</del>			k = 21	39.5	5 39.9	69.6	71.8
LSIM[24]		98.2		88.0		MN	IST $k = 5$	55 5	5 66 6	74.7	81.2
RNN(tanh)[2]	25]	$\approx$ 35.	0	$\approx$ 35.0	$p_{\perp}$	LVIIN.	K = J			74.7	01.2
stanh(s - 21)	11)	98.1		94 0			k = 9	55.5	) /1.1	78.6	86.9
$\sigma_{\text{unit}}(s - 21, 1)$	1 I J	70.1		<b>740</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)

# $\mathbf{W} = \mathbf{D}_3 \mathbf{R}_2 \mathcal{F}^{-1} \mathbf{D}_2 \mathbf{\Pi} \mathbf{R}_1 \mathcal{F} \mathbf{D}_1$

## 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





## 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



### Hierarchical Multiscale Recurrent Neural Networks (Chung, Ahn & Bengio arXiv:1609.01704)

- How to learn to update higher-level states at the right points in time, corresponding to slower time scales, but not necessarily on a fixed (clockwork) schedule, rather an adaptive, dynamic one?
- If each level updates once every 10 steps of lower level then 3<sup>rd</sup> level updates once every 1000 raw steps!
  - Advantage for propagating gradients for long-term dependencies
  - Computational advantage (but need HARD decisions)











### Hierarchical Multiscale Recurrent Neural Networks: Gating Mechanism

- Decide to copy the current level state vs updating it from lowerlevel state
- If updating it, restart the lower level RNN, conditioned on upper-  $z_{t-1}^{\ell}$  level state
- Need to take a binary decision in order to actually get a pure copy and avoid leakage of gradients
- Need that binary decision to be stochastic to have a chance to learn it



## How to backprop through stochastic binary decisions?

- REINFORCE: correlate the action with the reward, very high variance estimator
- Straight-through estimator (Hinton 2012, Bengio, Leonard & Courville arXiv 2013, Courbariaux et al NIPS 2015)

fprop

• Heuristic but worked very well in several settings



Anneal the slope during training

- MuProp (Gu et al ICLR 2016, arXiv 1511.05176) combines a soft decision with a hard one, unbiased baseline of REINFORCE uses the gradient of the loss wrt soft decision
- VIMCO (Mnih & Rezende ICLR 2016, arXiv:1602.06725)

### Hierarchical Multiscale Recurrent Neural Networks: Results (Chung, Ahn & Bengio arXiv:1609.01704)

- Automatically segmenting so as to better predict the next character
- SF-LSTM is not comparable (changes parameters on the fly)

Text8	
Model	BPC
<i>td</i> -LSTM (Zhang <i>et al.</i> , 2016)	1.63
HF-MRNN (Mikolov et al., 2012)	1.54
MI-RNN (Wu et al., 2016)	1.52
Skipping-RNN (Pachitariu and Sahani, 2013)	1.48
MI-LSTM (Wu et al., 2016)	1.44
Batch-normalized LSTM (Cooijmans et al., 2016)	1.36
HM-LSTM	1.30

Hutter Prize Wikipedia	
Model	BPC
SF-LSTM (Rocki, 2016b)*	1.37
Stacked LSTM (Graves, 2013)	1.67
MRNN (Sutskever et al., 2011)	1.60
GF-LSTM (Chung et al., 2015a)	1.58
Grid-LSTM (Kalchbrenner et al., 2015)	1.47
MI-LSTM (Wu et al., 2016)	1.44
Recurrent Highway Networks (Zilly et al., 2016)	1.42
Recurrent Memory Array Structures (Rocki, 2016a)	1.40
Layer-normalized LSTM (Ba et al., 2016) <sup>†</sup>	1.46
HM-LSTM	1.37

### Hierarchical Multiscale Recurrent Neural Networks: Results (Chung, Ahn & Bengio arXiv:1609.01704)

- Automatically segmenting so as to better predict the next character
- SF-LSTM is not comparable (changes parameters on the fly)

Text8	
Model	BPC
<i>td</i> -LSTM (Zhang <i>et al.</i> , 2016)	1.63
HF-MRNN (Mikolov et al., 2012)	1.54
MI-RNN (Wu et al., 2016)	1.52
Skipping-RNN (Pachitariu and Sahani, 2013)	1.48
MI-LSTM (Wu et al., 2016)	1.44
Batch-normalized LSTM (Cooijmans et al., 2016)	1.36
HM-LSTM	1.30

Hutter Prize Wikipedia	
Model	BPC
SF-LSTM (Rocki, 2016b)*	1.37
Stacked LSTM (Graves, 2013)	1.67
MRNN (Sutskever et al., 2011)	1.60
GF-LSTM (Chung et al., 2015a)	1.58
Grid-LSTM (Kalchbrenner et al., 2015)	1.47
MI-LSTM (Wu et al., 2016)	1.44
Recurrent Highway Networks (Zilly et al., 2016)	1.42
Recurrent Memory Array Structures (Rocki, 2016a)	1.40
Layer-normalized LSTM (Ba <i>et al.</i> , 2016) <sup>†</sup>	1.46
HM-LSTM	1.37









# Montreal Institute for Learning Algorithms

Universit

de Monti