## From Attention to Memory and towards Longer-Term Dependencies

### **Yoshua Bengio**

December 12, 2015

NIPS'2015 Reasoning, Attention & Memory Workshop PLUG: **Deep Learning**, NILL Press book in preparation, draft chapters online for feedback 



### Encoder-Decoder Framework

- Intermediate representation of meaning
  - = 'universal representation'
- Encoder: from word sequence to sentence representation
- Decoder: from representation to word sequence distribution



### Attention Mechanism for Deep Learning

- 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



### Content-Based & Location-Based Attention Mechanisms

- *(Graves 2013)*: location-based, handwriting generation
- (Bahdanau et al 2014): content-based, machine translation
- (Weston et al 2014, Graves et al 2014, Chorowski et al 2014 & NIPS 2015): content-based + location-based, speech recognition



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

(Bahdanau et al 2014, Jean et al 2014, Gulcehre et al 2015, Jean et al 2015)

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

	NMT(A)	Google	P-SMT
NMT	32.68	30.6*	
+Cand	33.28	_	27 02•
+UNK	33.99	32.7°	57.05
+Ens	36.71	<b>36.9</b> °	

(a) English → French (WMT-14)

#### (b) English→German (WMT-15)

(c) English→Czech (WMT-15)

Model	Note	Model	Note
24.8	Neural MT	18.3	Neural MT
24.0	U.Edinburgh, Syntactic SMT	18.2	JHU, SMT+LM+OSM+Sparse
23.6	LIMSI/KIT	17.6	CU, Phrase SMT
22.8	U.Edinburgh, Phrase SMT	17.4	U.Edinburgh, Phrase SMT
22.7	KIT, Phrase SMT	16.1	U.Edinburgh, Syntactic SMT

### IWSLT 2015 - Luong & Manning (2015) TED talk MT, English-German



BLEU (CASED)





### Image-to-Text: Caption Generation with Attention



(Xu et al., 2015), (Yao et al., 2015)

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <



## Speaking about what one sees



is(0.22)



with(0.28),



the(0.21)



on(0.25)







background(0.11)



a(0.21)



mountain(0.44)



.(0.13)





road(0.26)





### The Good



A woman is throwing a <u>frisbee</u> in a park.



A <u>dog</u> is standing on a hardwood floor.



A <u>stop</u> sign is on a road with a mountain in the background.



A little <u>girl</u> sitting on a bed with a teddy bear.



A group of <u>people</u> sitting on a boat in the water.



A giraffe standing in a forest with <u>trees</u> in the background.

### And the Bad



A large white <u>bird</u> standing in a forest.



A woman holding a <u>clock</u> in her hand.



A man wearing a hat and a hat on a <u>skateboard</u>.



A person is standing on a beach with a <u>surfboard.</u>

A woman is sitting at a table with a large <u>pizza.</u>

A man is talking on his cell phone while another man watches.

## Attention on Memory Elements

Recurrent networks cannot remember things for very long

▶The cortex only remember things for 20 seconds

We need a "hippocampus" (a separate memory module)

- LSTM [Hochreiter 1997], registers
- Memory networks [Weston et 2014] (FAIR), associative memory

NTM [Graves et al. 2014], "tape".



### Ongoing Project: Knowledge Extraction

- Learn to fill the memory network from natural language descriptions of facts
- Force the neural net to understand language
- Extract knowledge from documents into a usable form



# Long-Term Dependencies

 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].

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

Gradient

clipping

- Problems:
  - sing. values of Jacobians > 1  $\rightarrow$  gradients explode
  - or sing. values  $< 1 \rightarrow gradients shrink \& vanish$  (Hochreiter 1991)
  - or random  $\rightarrow$  variance grows exponentially

## Gated Recurrent Units & LSTM

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



## Delays & Hierarchies to Reach Farther

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

 $Q_{t+1}$ 



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

- A mental state stored in an external memory can stay for arbitrarily long durations, until evoked for read or write
- Forgetting = vanishing gradient.
- Memory = larger state, avoiding the need for forgetting/vanishing



### Paths where gradient flows unhampered → leigenvalue = 1

- Self-loop with weight = 1
- Copy operation in memory network

For each scalar that is so preserved, gradient can flow back unhampered

- And there must exist a direction in which the Jacobian preserves vector magnitude
- Corresponds to an eigenvalue whose magnitude is 1

## Unitary Evolution RNNs

### Martin Arjowski, Amar Shah & Yoshua Bengio arXiv 1511.06464 Submitted to ICLR 2016

Large state: cannot afford  $O(n^2)$  recurrent computation

Orthogonal matrices C Unitary matrices

Non-trivial to parametrize efficiently

## What if W had all eigenvalues=1

$$z_{t+1} = \mathbf{W}_t h_t + \mathbf{V}_t x_{t+1}$$
$$h_{t+1} = \sigma(z_{t+1})$$

 $\frac{\partial C}{\partial h_t} = \frac{\partial C}{\partial h_T} \frac{\partial h_T}{\partial h_t} = \frac{\partial C}{\partial h_T} \prod_{k=t}^{T-1} \frac{\partial h_{k+1}}{\partial h_k} = \frac{\partial C}{\partial h_T} \prod_{k=t}^{T-1} \mathbf{D}_{k+1} \mathbf{W}_k^T,$  $\left\| \frac{\partial C}{\partial h_t} \right\| = \left\| \frac{\partial C}{\partial h_T} \prod_{k=t}^{T-1} \mathbf{D}_{k+1} \mathbf{W}_k^T \right\| \le \left\| \frac{\partial C}{\partial h_T} \right\| \prod_{k=t}^{T-1} \| \mathbf{D}_{k+1} \mathbf{W}_k^T \| = \left\| \frac{\partial C}{\partial h_T} \right\| \prod_{k=t}^{T-1} \| \mathbf{D}_{k+1} \|$  $\left\| \frac{\partial C}{\partial h_t} \right\| \le \left\| \frac{\partial C}{\partial h_T} \right\| \prod_{k=t}^{T-1} \| \mathbf{D}_{k+1} \| = \left\| \frac{\partial C}{\partial h_T} \right\|$ 

Guaranteed no explosion, no need for gradient clipping

## Building Blocks for Unitary W

 $U = V D V^*$  with V fixed would have a bad ratio of Computation  $O(n^2)$  to parameters O(n), and D complex even for orthogonal matrices

- **D**, a diagonal matrix with  $\mathbf{D}_{j,j} = e^{iw_j}$ , with parameters  $w_j \in \mathbb{R}$ ,
- $\mathbf{R} = \mathbf{I} 2 \frac{vv^*}{\|v\|^2}$ , a reflection matrix in the complex vector  $v \in \mathbb{C}^n$ ,
- $\Pi$ , a fixed random index permutation matrix, and
- $\mathcal{F}$  and  $\mathcal{F}^{-1}$ , the Fourier and inverse Fourier transforms.

O(n) computation & params

O(n log n) computation

## Building Blocks for Unitary W

- **D**, a diagonal matrix with  $\mathbf{D}_{j,j} = e^{iw_j}$ , with parameters  $w_j \in \mathbb{R}$ ,
- $\mathbf{R} = \mathbf{I} 2 \frac{vv^*}{\|v\|^2}$ , a reflection matrix in the complex vector  $v \in \mathbb{C}^n$ ,
- $\Pi$ , a fixed random index permutation matrix, and
- $\mathcal{F}$  and  $\mathcal{F}^{-1}$ , the Fourier and inverse Fourier transforms.

A Recipe Inspired by FastFood (Le et al ICML 2013) and ACDC (Moczulski et al arXiv 2015)  $W = D_3 R_2 \mathcal{F}^{-1} D_2 \Pi R_1 \mathcal{F} D_1$ 



We do not want to destroy the information present in the phase

$$\sigma_{\text{modReLU}}(z) = \begin{cases} (|z|+b)\frac{z}{|z|} & \text{if } |z|+b \ge 0\\ 0 & \text{if } |z|+b < 0 \end{cases}$$

$$\sigma_{\text{modReLU}}(z) = \sigma_{\text{ReLU}}(|z|+b)\frac{z}{|z|}$$



### Adding Problem







What is going on? URNN forgets at a Lower rate and LSTM state gets stuck at some point



### MILA: Montreal Institute for Learning Algorithms

