

# Deep Learning

## Lecture 6

Thomas Hofmann, ETH Zurich

31 October 2019

# Section 11

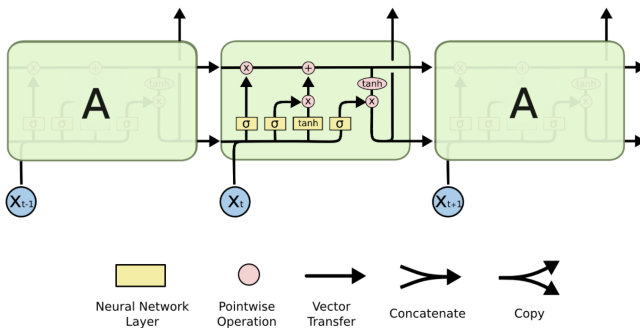
## Differentiable Memory

# Long-Term Dependencies

- 1 Sometimes: important to model long-term dependencies  $\implies$  network needs to **memorize** features from the distant past
- 2 Recurrent network: hidden state needs to preserve memory
- 3 Conflicts with short-term fluctuations and vanishing gradients
- 4 Conclusion: difficult to learn long-term dependencies with standard recurrent network (see DL Section 10.7)
- 5 Popular remedy: **gated units**

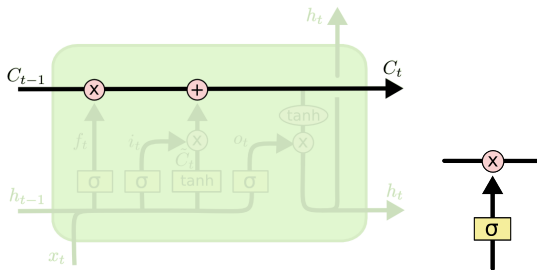
# LSTM: Overall Architecture

**Long-Short-Term-Memory:** unit for memory management



from Christopher Olah: <http://colah.github.io/posts/2015-08-Understanding-LSTMs>

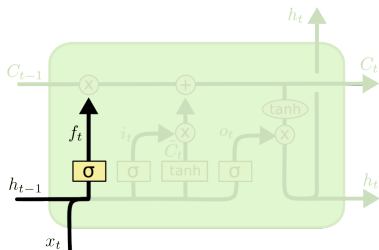
# LSTM: Flow of Information



- ① Information propagates along the chain like on a conveyor belt.
- ② Information can flow unchanged and is only selectively changed (vector addition) by  $\sigma$ -gates.

from Christopher Olah: <http://colah.github.io/posts/2015-08-Understanding-LSTMs>

# LSTM: Forget Gate

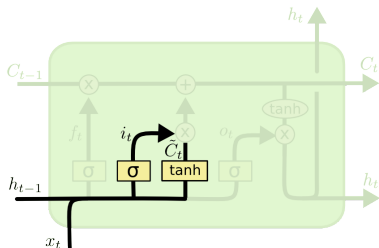


$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

## 1 Keeping or forgetting of stored content?

from Christopher Olah: <http://colah.github.io/posts/2015-08-Understanding-LSTMs>

# LSTM: Input $\rightarrow$ Memory Value



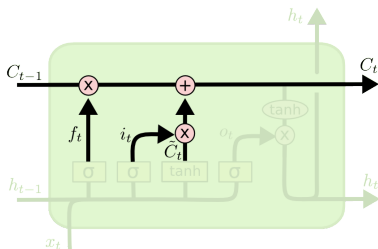
$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

- 1 Preparing new input information to be added to the memory.

from Christopher Olah: <http://colah.github.io/posts/2015-08-Understanding-LSTMs>

# LSTM: Updating Memory



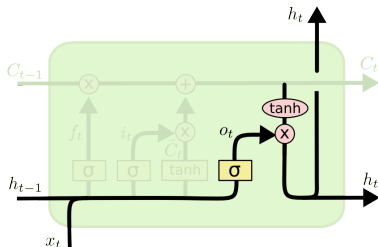
$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

- 1 Combining stored and new information.

from Christopher Olah: <http://colah.github.io/posts/2015-08-Understanding-LSTMs>



# LSTM: Output Gate



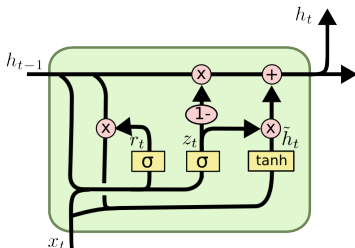
$$o_t = \sigma(W_o [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * \tanh(C_t)$$

- 1 Computing output selectively.

from Christopher Olah: <http://colah.github.io/posts/2015-08-Understanding-LSTMs>

# Gated Memory Units



$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t])$$

$$r_t = \sigma(W_r \cdot [h_{t-1}, x_t])$$

$$\tilde{h}_t = \tanh(W \cdot [r_t * h_{t-1}, x_t])$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

- ① Memory state = output. Modifications to logic (Cho et al, 2014).
- ② Convex combination of old and new information.

from Christopher Olah: <http://colah.github.io/posts/2015-08-Understanding-LSTMs>

# Gated Memory Units

- ① GRUs and LSTMs can learn active memory strategies: what to memorize, overwrite and recall when.
- ② Successful use cases
  - handwriting recognition
  - speech recognition (also: Google)
  - machine translation
  - image captioning
- ③ Notoriously difficult to understand what units learn...  
Resource-hungry. Slow in learning.

# Language Modeling

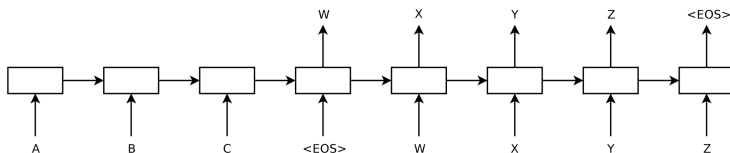
MODEL	TEST PERPLEXITY	NUMBER OF PARAMS [BILLIONS]
SIGMOID-RNN-2048 (JI ET AL., 2015A)	68.3	4.1
INTERPOLATED KN 5-GRAM, 1.1B N-GRAMS (CHELBA ET AL., 2013)	67.6	1.76
SPARSE NON-NEGATIVE MATRIX LM (SHAZEER ET AL., 2015)	52.9	33
RNN-1024 + MAXENT 9-GRAM FEATURES (CHELBA ET AL., 2013)	51.3	20
LSTM-512-512	54.1	0.82
LSTM-1024-512	48.2	0.82
LSTM-2048-512	43.7	0.83
LSTM-8192-2048 (NO DROPOUT)	37.9	3.3
LSTM-8192-2048 (50% DROPOUT)	32.2	3.3
2-LAYER LSTM-8192-1024 (BIG LSTM)	30.6	1.8
BIG LSTM+CNN INPUTS	<b>30.0</b>	<b>1.04</b>
BIG LSTM+CNN INPUTS + CNN SOFTMAX	39.8	<b>0.29</b>
BIG LSTM+CNN INPUTS + CNN SOFTMAX + 128-DIM CORRECTION	35.8	<b>0.39</b>
BIG LSTM+CNN INPUTS + CHAR LSTM PREDICTIONS	47.9	<b>0.23</b>

from Jozefowicz et al, 2016

- 1 evaluation on corpus w/ 1B words
- 2 number of parameters can be in the 100Ms or even Bs!
- 3 ensembles can reduce perplexity to  $\approx 23$  (best result 06/2016)

# Sequence to Sequence Learning

- 1 Important use of memory units: **sequence to sequence learning**. Seminal paper: **Sutskever, Vinyals & Le, 2014**
- 2 **Encoder-decoder architecture**



Encode sequence (e.g. sentence) into vector

Decode sequence (e.g. translate) from vector (w/  
autoregressive output feedback)

# RNN Encoder/Decoder

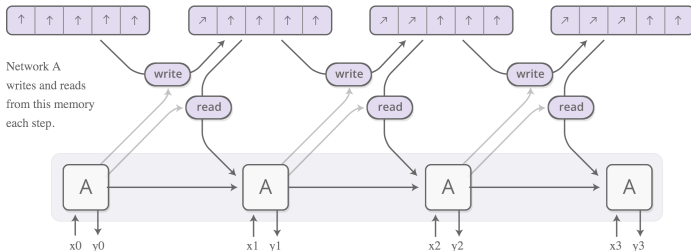
How to make this work? (Sutskever, Vinyals & Le, 2014)

- 1 Deep LSTMs (multiple layers, e.g. 4)
- 2 Different RNNs for encoding and decoding
- 3 Teacher forcing (maximum likelihood) during training
- 4 Beam search for decoding at test time
- 5 Reverse order of source sequence
- 6 Ensemble-ing
- 7  $\implies$  State-of-the art results on WMT benchmarks at the time

Today: use of attention-based models (see next section)

# Neural Turing Machine: Architecture

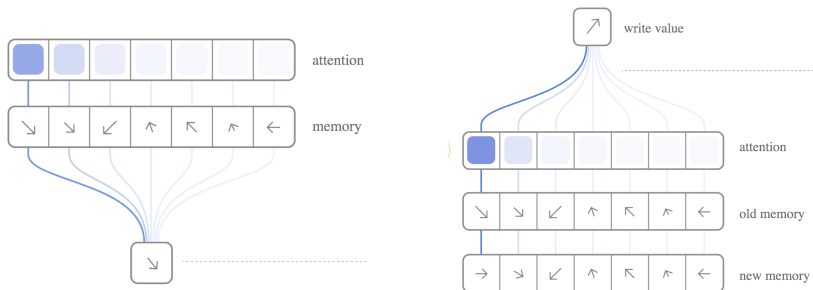
Memory is an array of vectors.



- 1 RNN controls an external memory bank
- 2 Reminiscent of Turing machine, but: each cell  $M_i \in \mathbb{R}^d$

from Olah & Carter, 2016: <http://distill.pub/2016/augmented-rnns>

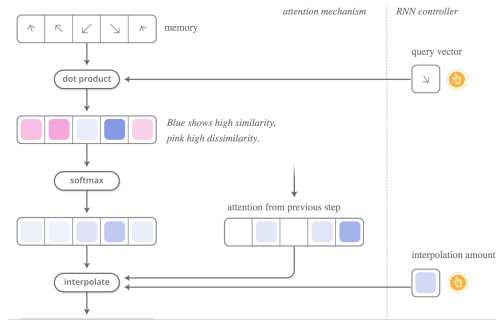
# Neural Turing Machine: Differentiable Memory



- 1 Compute attention distribution  $(\alpha_i)_i$ ,  $\alpha_i \geq 0$  s.t.  $\sum_i \alpha_i = 1$ .
- 2 Read out expected memory content:  $r \leftarrow \sum_i \alpha_i M_i$ .
- 3 Write uses weights  $(\beta_i)_i$ ,  $\beta_i \in [0; 1]$ ,  $M_i \leftarrow (1 - \beta_i)M_i + \beta_i w$ .

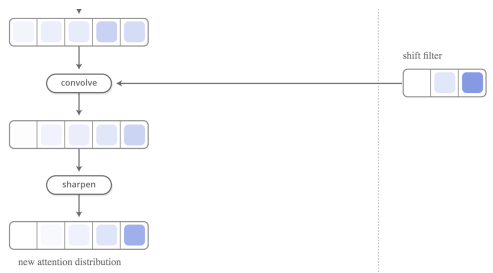


# Neural Turing Machine: Memory Controller



- 1 Associative memory access with query (key)  $q \in \mathbb{R}^d$ .
- 2 Cells are scored via softmax.

# Neural Turing Machine: Memory Controller



- ① Ability to shift relative to content-selected locations (convolution).
- ② Additional accentuation to sharpen attention distribution.

from Olah & Carter, 2016: <http://distill.pub/2016/augmented-rnns>

# Differentiable Memory: Discussion

- ① NTM architectures can learn loops and simple programs.  
However, no real-world applications.
- ② Other architectures:
  - Neural random access machines ([Kurach et al, 2015](#))
  - Differentiable data structures like stacks and queues ([Grefenstette et al., 2015](#); [Joulin & Mikolov, 2015](#))
- ③ Direction of future research

# Section 12

## Attention

# Attention Mechanisms

- 1 Simple way to overcome some challenges of RNN-based memorization: **attention mechanism**

Selectively attend to **inputs** or **feature representations** computed from inputs.

- 2 **RNNs**: learn to encode information relevant for the future.

vs.

**Attention**: select what is relevant from the past in hindsight!

Both ideas can be combined

# Gating Function

## Definition (Softmax Gating Function)

A softmax gating function  $f_\phi$  takes as input a query vector  $\xi \in \mathbb{R}^n$  as well as a set of values  $\mathbf{x}^t \in \mathbb{R}^m$  ( $t = 1, \dots, s$ ) and is defined as

$$f_\phi(\xi, (\mathbf{x}^1, \dots, \mathbf{x}^s)) = \frac{1}{\sum_j \exp[\phi(\xi, \mathbf{x}^j)]} \begin{pmatrix} \exp[\phi(\xi, \mathbf{x}^1)] \\ \vdots \\ \exp[\phi(\xi, \mathbf{x}^s)] \end{pmatrix}$$

for some similarity or compatibility function  $\phi : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}$ .

- 1  $\phi$  can often be learned in a black-box manner via a MLP
- 2 simplest choice for  $n = m$ :  $\phi(\xi, \mathbf{x}) = \xi^\top \mathbf{x}$  (inner product)
- 3 every restriction  $f_\phi(\xi, \cdot)$  maps to the interior of a simplex

# Self-Gated Attention

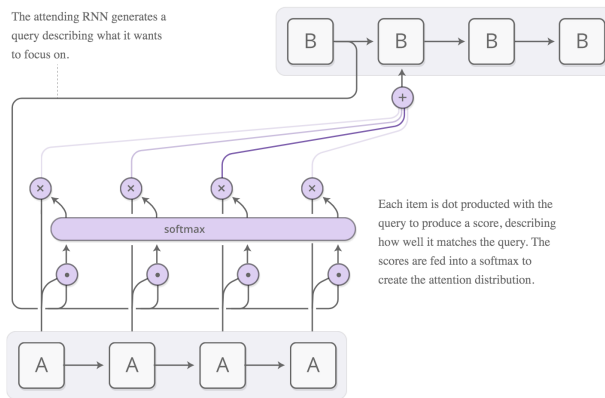
## Definition (Self-Gated Attention)

Given a query  $\xi \in \mathbb{R}^m$  and a set of values  $\mathbf{x}_i \in \mathbb{R}^n$  ( $i = 1, \dots, k$ ). The self-gated attention is defined as

$$\underbrace{F(\xi, (\mathbf{x}_1, \dots, \mathbf{x}_k))}_{\in \mathbb{R}^k} = \underbrace{[\mathbf{x}_1 \quad \mathbf{x}_2 \quad \dots \quad \mathbf{x}_k]}_{\in \mathbb{R}^{k \times n}} \cdot \underbrace{f_\phi(\xi, (\mathbf{x}_1, \dots, \mathbf{x}_k))}_{\in \mathbb{R}^n}$$

where  $f_\phi$  is a gating function.

# Seq2seq with Attention: Schematic



from Olah & Carter, 2016: <http://distill.pub/2016/augmented-rnns>



# Seq2seq with Attention

- 1 Attend to the hidden state of the encoding RNN, i.e. values  $(\mathbf{h}_e^1, \dots, \mathbf{h}_e^s)$ .
- 2 Decoding RNN produces query at each time, i.e.  $(\xi^1, \dots, \xi^{s'})$ .
- 3 Self-gated attention produces “read-out”  $\mathbf{z}^t$  from encoder sequence
- 4 Used as input to the decoding RNN:  $(\mathbf{h}_d^t, \mathbf{z}^t) \mapsto \mathbf{h}_d^{t+1}$

# Seq2seq with Attention: MT Example

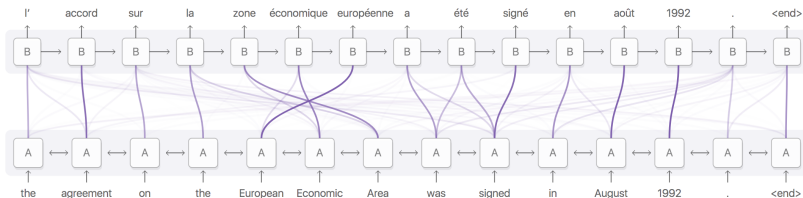
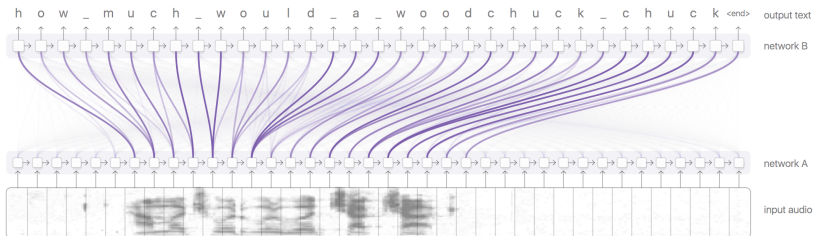


Diagram derived from Fig. 3 of Bahdanau, et al. 2014

from Olah & Carter, 2016: <http://distill.pub/2016/augmented-rnns>

- 1 Interpretable attention model (akin to alignments)  
Bahdanau, Cho & Bengio, 2014
- 2 Bi-directional GRU encoder, left-to-right GRU decoder

# Seq2seq with Attention: Speech Recognition



from Olah & Carter, 2016: <http://distill.pub/2016/augmented-rnns>

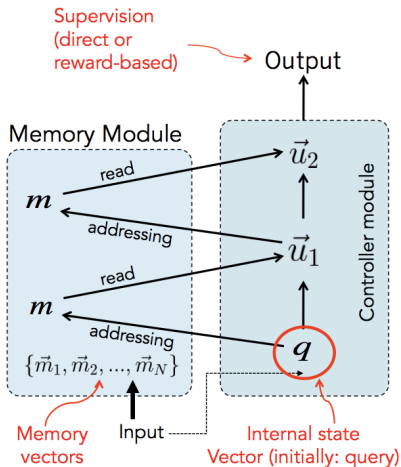
- 1 Listen, Attend and Spell Modell (Chan et al, 2015)
- 2 Bi-directional, pyramidal LSTM encoder

# Memory Networks

Memory networks (Weston et al, 2014; Kumar et al, 2015)

- ① operate over large data corpus (e.g. text documents)
- ② given question, learn to infer answers (QA retrieval)
- ③ simplest form: memory is not altered
- ④ **recursive associative recall**: given query  $q$ , find best matching memory cell  $i$ ; use  $M_i$  and  $x$  as new key; repeat

# Memory Networks



from Weston, ICML Tutorial 2016: <http://www.thespermwhale.com/jaseweston/icml2016/>

# Key-Value Attention

## Definition (Key-Value Attention)

Given a query  $\xi \in \mathbb{R}^n$ , key-value pairs  $(\mathbf{x}^t, \mathbf{z}^t) \in \mathbb{R}^n \times \mathbb{R}^m$ ,  $t = 1 \dots, s$  and a gating function  $f$ . The  $(n, m)$ -dimensional key-value attention map is defined as

$$F(\xi, ((\mathbf{x}^1, \mathbf{z}^1) \dots, (\mathbf{x}^s, \mathbf{z}^s))) = [\mathbf{z}^1 \quad \mathbf{z}^2 \quad \dots \quad \mathbf{z}^s] \cdot f(\xi, (\mathbf{x}^1, \dots, \mathbf{x}^s))$$

- 1 attention weights are computed based on keys
- 2 produced value is linear (or convex) combination of values
- 3 keys determine where to look, values determine what features get extracted

# Dot-Product Attention

## Definition (Scaled Dot-Product Attention)

The attention map induced by

$$f(\boldsymbol{\xi}, \mathbf{x}) = \frac{\boldsymbol{\xi}^\top \mathbf{x}}{\sqrt{n}}$$

is called scaled dot-product attention.

- 1 simple dot-product similarity between query and key, not necessarily convex (soft-max)
- 2 motivation for normalization: assume  $\boldsymbol{\xi}, \mathbf{x}$  are random  $n$ -vectors with zero mean and unit variances, then

$$\mathbf{E} \left[ \boldsymbol{\xi}^\top \mathbf{x} \right] = 0 \quad \text{and} \quad \mathbf{E} \left[ (\boldsymbol{\xi}^\top \mathbf{x})^2 \right] = n$$

# Multi-Headed Attention

## Definition (Multi-Headed Attention)

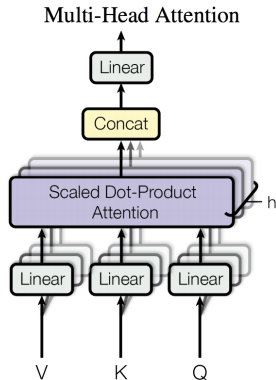
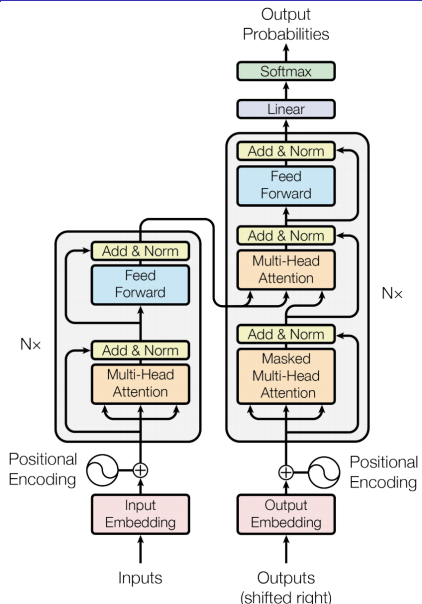
Let  $F_j$ ,  $1 \leq j \leq r$  be  $(n, m)$ -dimensional key-value attention map. An  $r$  multi-headed  $(N, M)$ -dimensional attention map  $G$  is defined as follows:

$$G(\boldsymbol{\xi}, (\mathbf{x}^t, \mathbf{z}^t)_{t=1}^s) = \mathbf{W} \begin{bmatrix} F_1(\mathbf{W}_1^q \boldsymbol{\xi}, (\mathbf{W}_1^x \mathbf{x}^t, \mathbf{W}_1^z \mathbf{z}^t)_{t=1}^s) \\ \vdots \\ F_r(\mathbf{W}_r^q \boldsymbol{\xi}, (\mathbf{W}_r^x \mathbf{x}^t, \mathbf{W}_r^z \mathbf{z}^t)_{t=1}^s) \end{bmatrix}$$

- 1 matrices  $\mathbf{W}_i^q, \mathbf{W}_i^x \in \mathbb{R}^{n \times N}$  and  $\mathbf{W}_i^z \in \mathbb{R}^{m \times M}$  are linear dimension-reduction matrices (typically:  $n < N$  and  $m < M$ ).
- 2  $\mathbf{W} \in \mathbb{R}^{M \times r \cdot m}$  adjusts the dimension (typically: reduction)
- 3 example: design choice in (Vaswani et al, 2017):  $r = 8$ ,  $n = m = 64$ ,  $N = M = 512$ .



# Transformer Architecture: Overview



(Vaswani et al, 2017)

# Transformer Architecture: Other Design Choices

- 1 Fully-connected feedforward networks (specifically: ReLU with layer width  $512 \mapsto 2048 \mapsto 512$ , cf.  $1 \times 1$  convolution)
- 2 Positional encoding: learned or fixed (sine-functions of different frequency)
- 3 Layer normalization (Ba, Krios, Hinton, 2017), cf. later section on activity re-normalization
- 4 Skip connections with adds (cf. residual layers)