

UCU ML Summer Workshops 2017

Natural Language Processing: Part 2



DataRobot

Yuriy Guts
ML Engineer

Remaining Agenda

Day 4: Short Text Similarity & Paraphrase Detection

- Paraphrase identification problem
- Traditional and DL approaches to semantic similarity
- Architectural principles for state-of-the-art neural NLP
- DL feature extraction for NLP

Day 5: Sequence-to-Sequence Modeling

- Attention mechanisms in DL for NLP
- Sequence-to-sequence neural models

Paraphrase Identification



Where can I get very professional and reliable envelope printing service in Sydney?

Where can I get very affordable branded envelope printing service in Sydney?



Why are doctors always late?

Why doctors always make you wait for 15-20 minutes before they see you?

Recap: BoW Representations

Documents



Vector-space
representation

However, complexity
We will see how small
Given a function-based
Using entropy of traffic
We study the complexity
of influencing elections
through bribery. How
computationally complex
is it for an external actor
to determine whether by
a certain amount of
bribing voters a specified
candidate can be made
the election's winner? We
study this problem for
election systems as varied
as scoring ...

	D1	D2	D3	D4	D5
complexity	2		3	2	3
algorithm	3			4	4
entropy	1			2	
traffic		2	3		
network		1	4		

Term-document matrix

TFIDF

For a term i in document j :

$$w_{i,j} = tf_{i,j} \times \log \left(\frac{N}{df_i} \right)$$

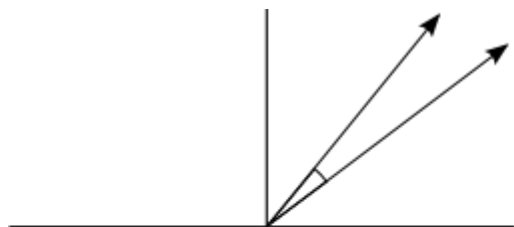
$tf_{i,j}$ = number of occurrences of i in j

df_i = number of documents containing i

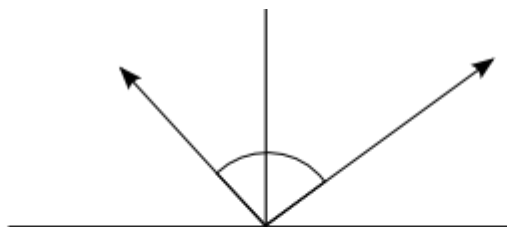
N = total number of documents

Recap: Cosine Similarity

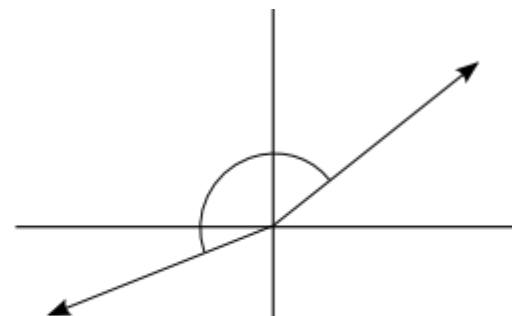
$$\text{similarity} = \cos(\theta) = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\|_2 \|\mathbf{B}\|_2} = \frac{\sum_{i=1}^n A_i B_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}}$$



Similar scores
Score Vectors in same direction
Angle between them is near 0 deg.
Cosine of angle is near 1 i.e. 100%

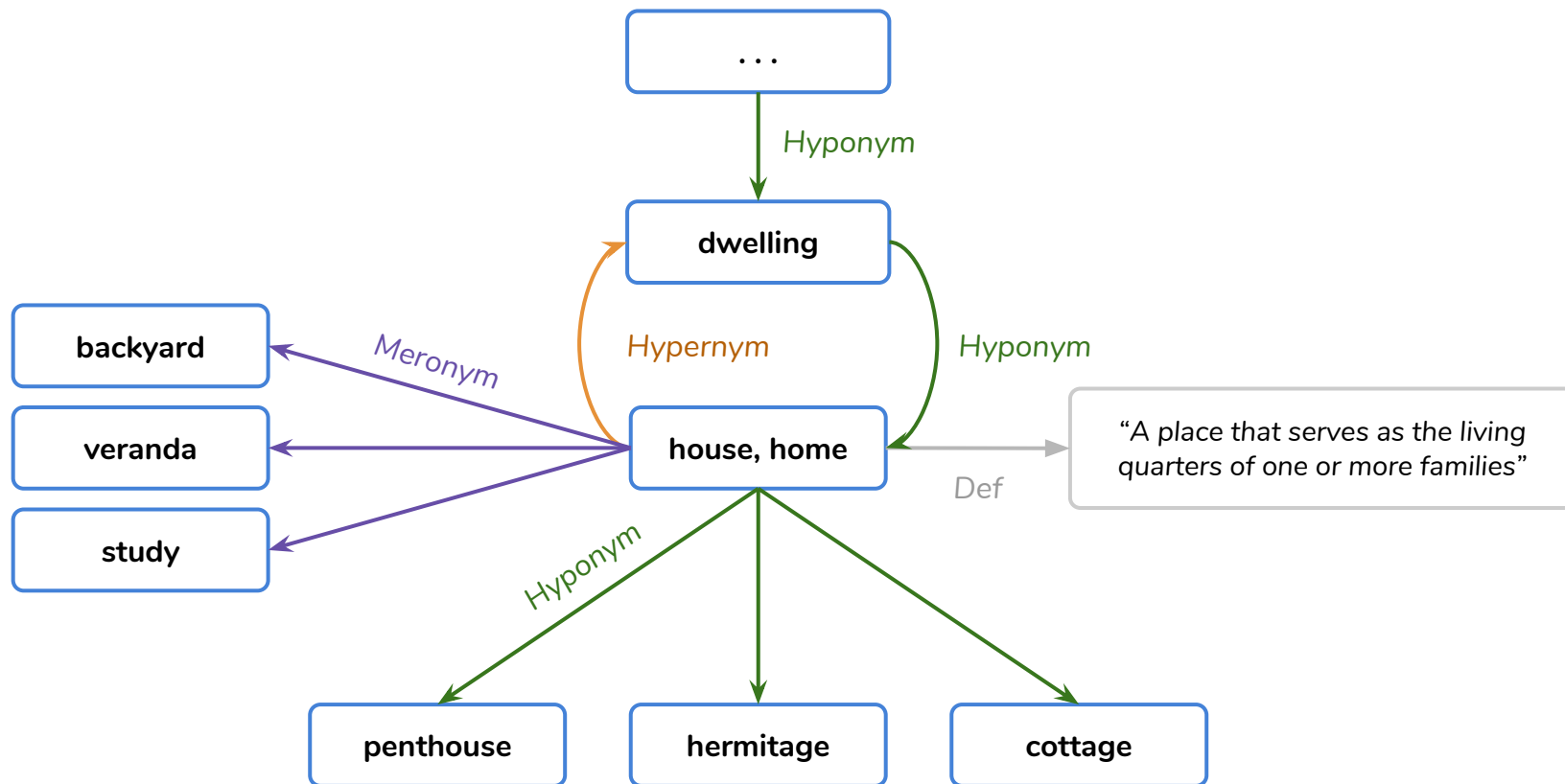


Unrelated scores
Score Vectors are nearly orthogonal
Angle between them is near 90 deg.
Cosine of angle is near 0 i.e. 0%

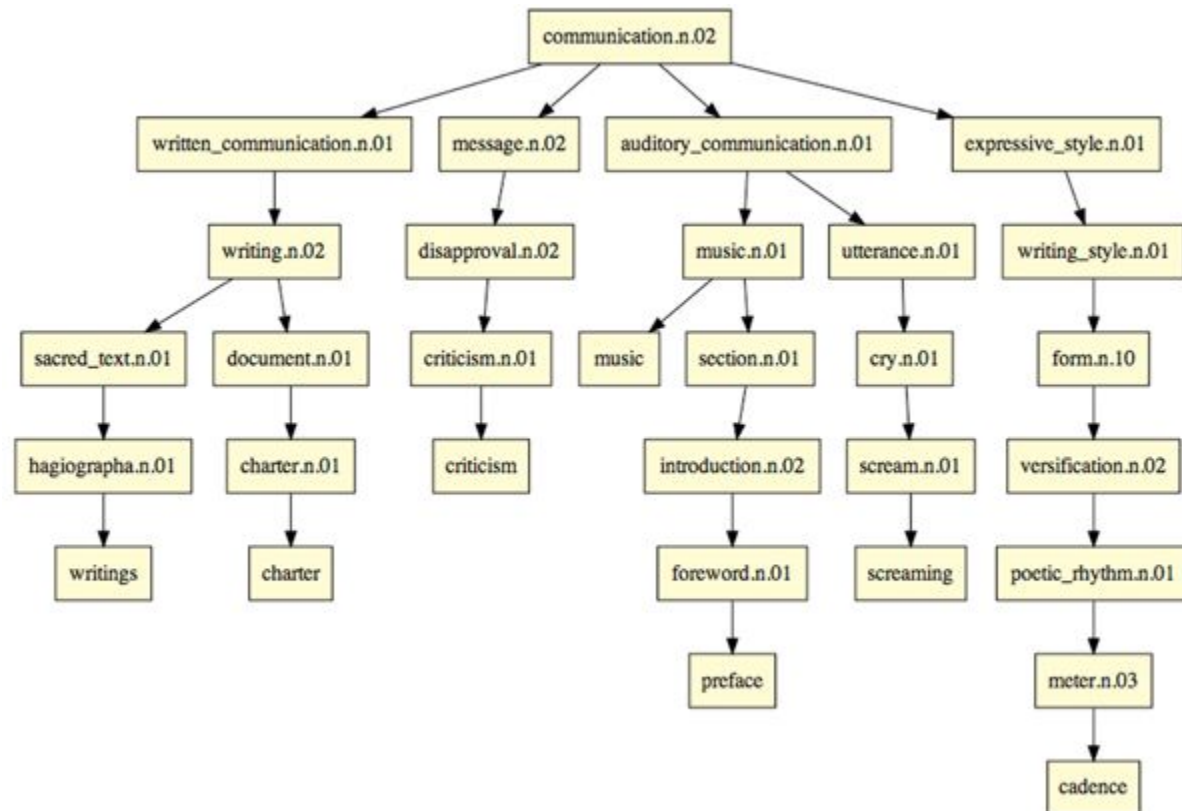


Opposite scores
Score Vectors in opposite direction
Angle between them is near 180 deg.
Cosine of angle is near -1 i.e. -100%

Lexical Databases and Ontologies



WordNet



Distributed Hypothesis of Language

“The complete meaning of a word is always contextual, and no study of meaning apart from context can be taken seriously”

John R. Firth. The technique of semantics.

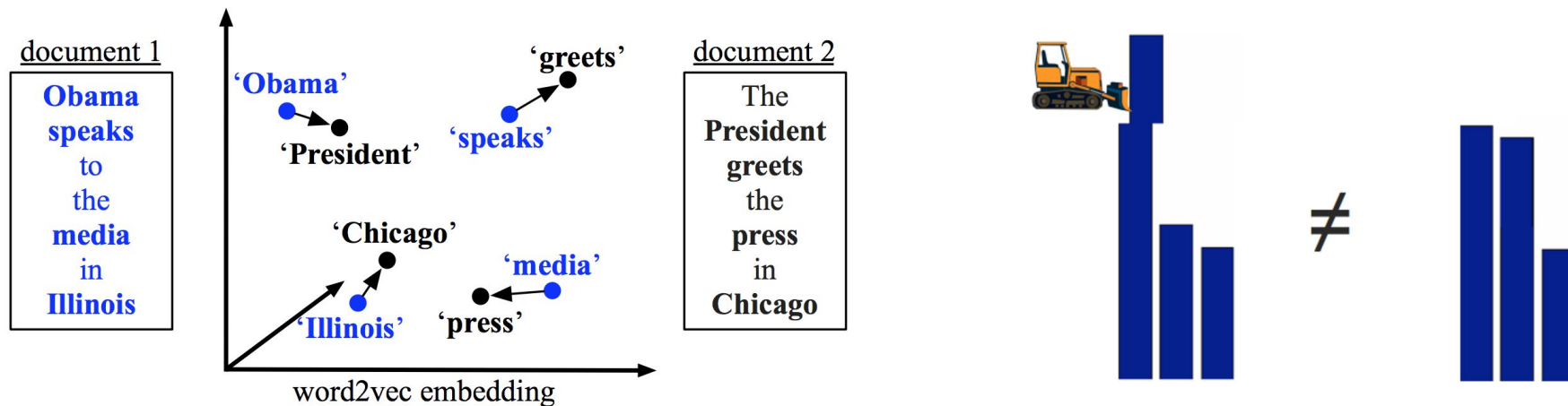
Transactions of the Philological Society, 1935.

Distributional Hypothesis of Language [Harris, 1954]

“Words that occur in the same contexts tend to have similar meanings”

The functional interplay of philosophy and **science** should, as a minimum, guarantee...
...and among works of dystopian **science** fiction...
 The rapid advance in **science** today suggests...
...calculus, which are more popular in **science** -oriented schools.
 But because **science** is based on mathematics...
...the value of opinions formed in **science** as well as in the religions...
 ...if **science** can discover the laws of human nature...
 ...is an art, not an exact **science** .
...factors shaping the future of our civilization: **science** and religion.
...certainty which every new discovery in **science** either replaces or reshapes.
...if the new technology of computer **science** is to grow significantly
 He got a **science** scholarship to Yale.
...frightened by the powers of destruction **science** has given...
...but there is also specialization in **science** and technology...

Word Mover's Distance (WMD)

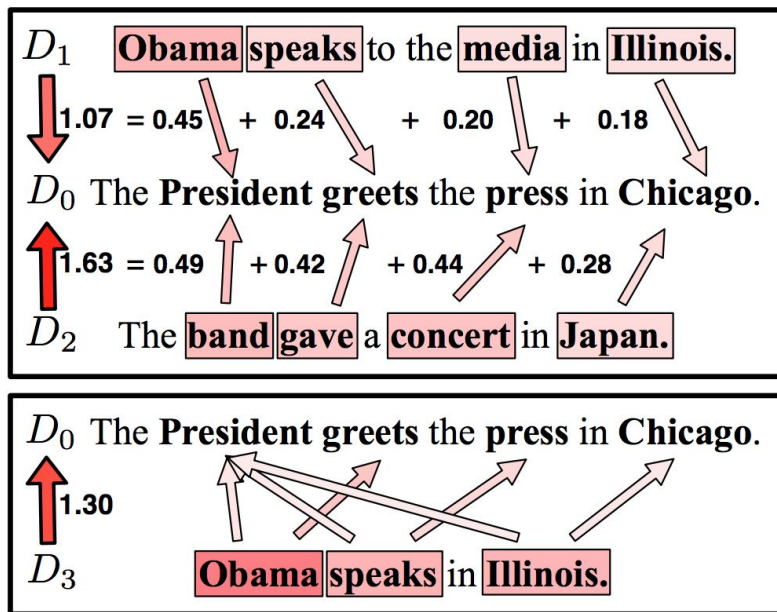


The minimum amount of “work” needed to transform document **1** to document **2**. Inspired by Earth Mover’s Distance, a well-studied transportation problem.

M. Kusner et al. “From Word Embeddings to Document Distances”, 2015.

<http://proceedings.mlr.press/v37/kusnerb15.pdf>

WMD: Linear Optimization Problem



$$d_i = \frac{c_i}{\sum_{j=1}^n c_j} \quad \text{nBOW frequency of the } i\text{-th word in the document}$$

$$\mathbf{T}_{ij} \geq 0 \quad \text{"How much" of word } i \text{ travels to word } j$$

$$\min_{\mathbf{T} \geq 0} \sum_{i,j=1}^n \mathbf{T}_{ij} c(i, j)$$

$$\text{subject to: } \sum_{j=1}^n \mathbf{T}_{ij} = d_i \quad \forall i \in \{1, \dots, n\}$$

$$\sum_{i=1}^n \mathbf{T}_{ij} = d'_j \quad \forall j \in \{1, \dots, n\}.$$

M. Kusner et al. "From Word Embeddings to Document Distances", 2015.

<http://proceedings.mlr.press/v37/kusnerb15.pdf>

Architectural Principles

for State-of-the-Art Neural NLP

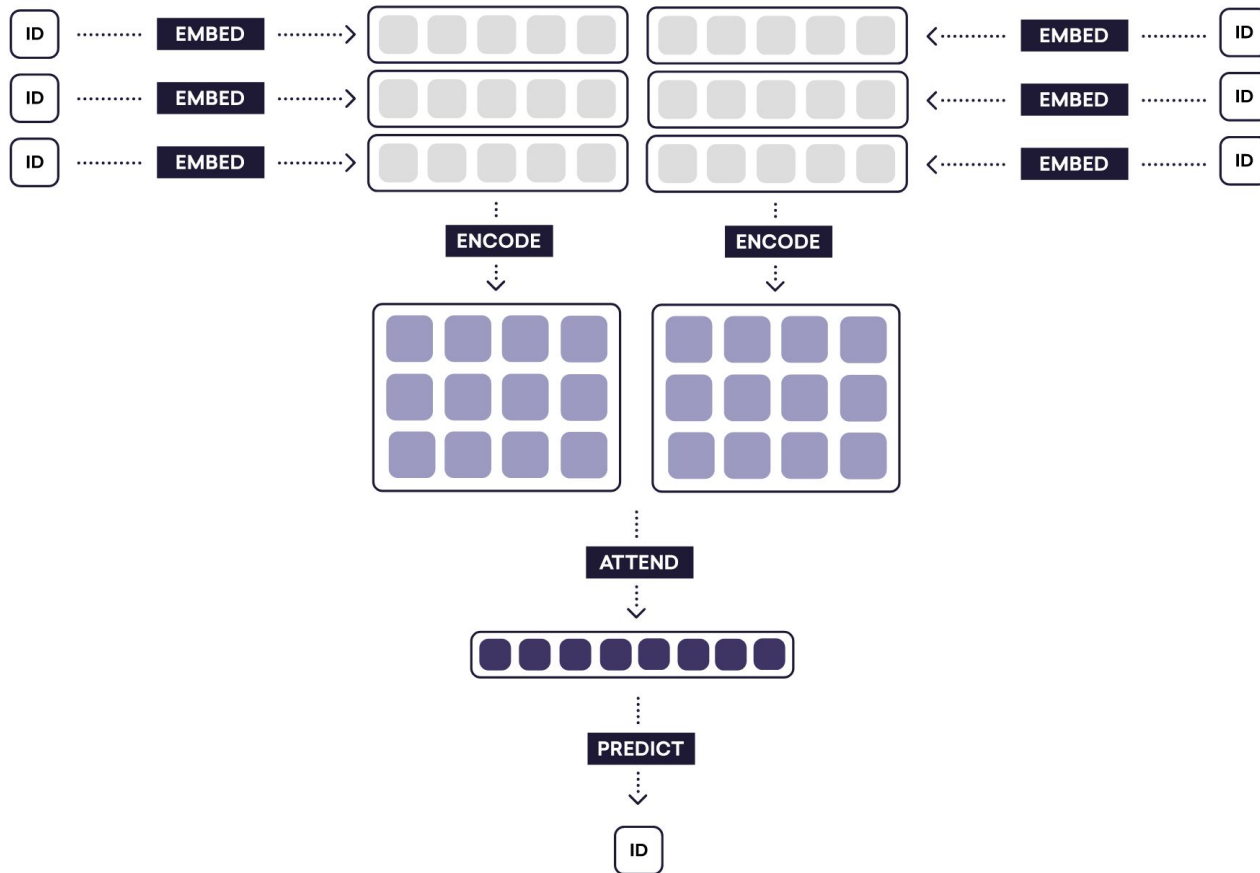
Embed

Encode

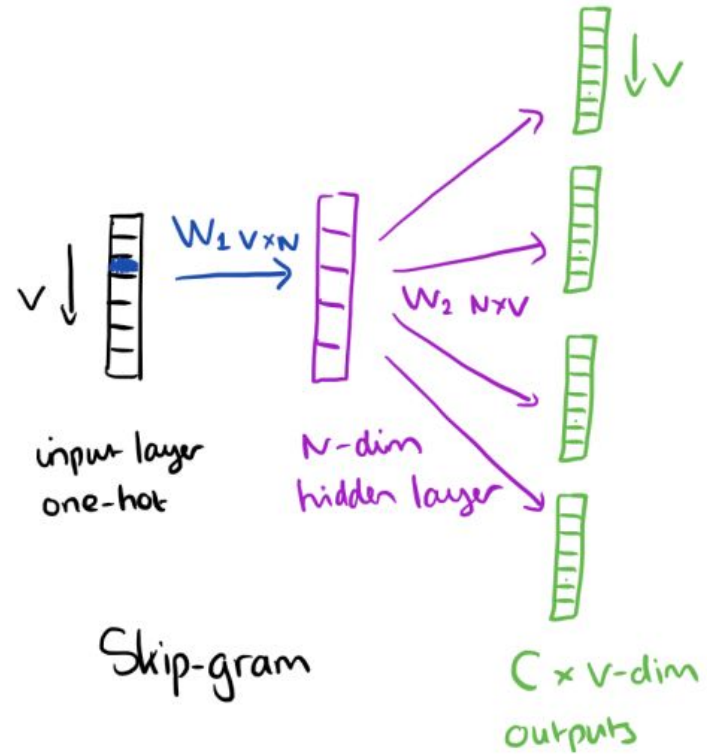
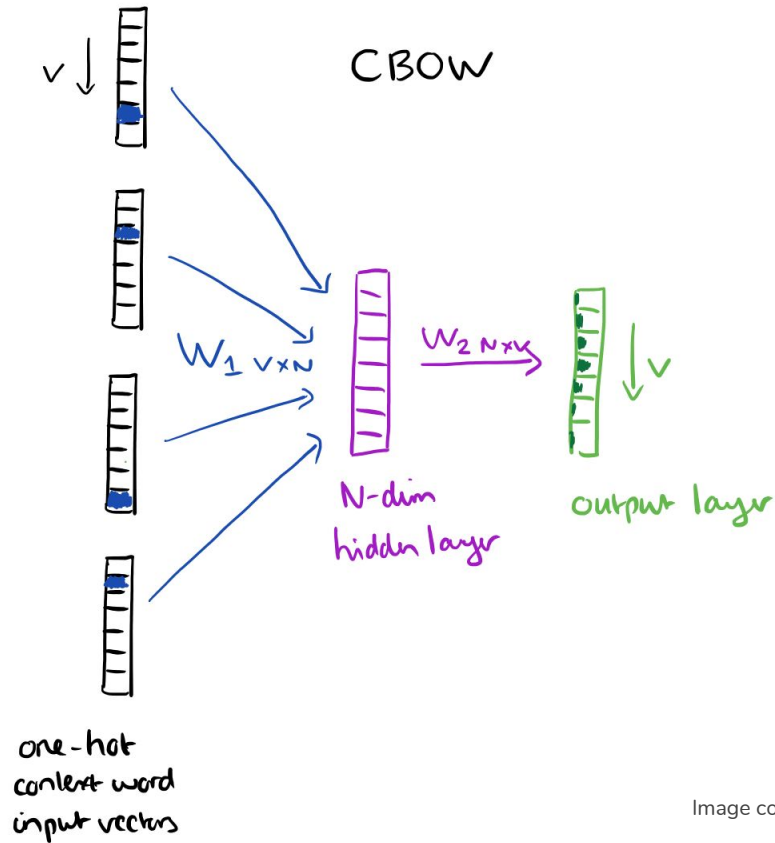
Attend

Predict

State-of-the-Art NLP Pipeline



Embed



Encode: Convolutions over Word Embeddings

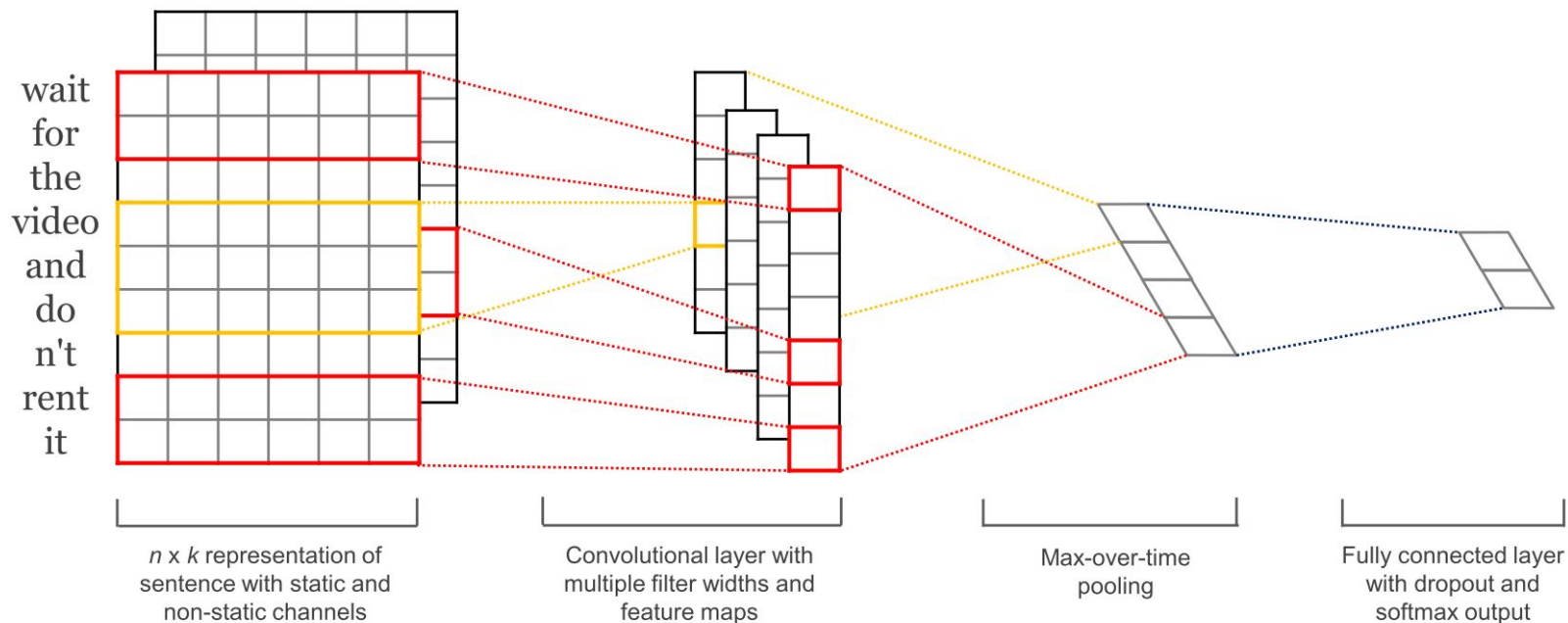
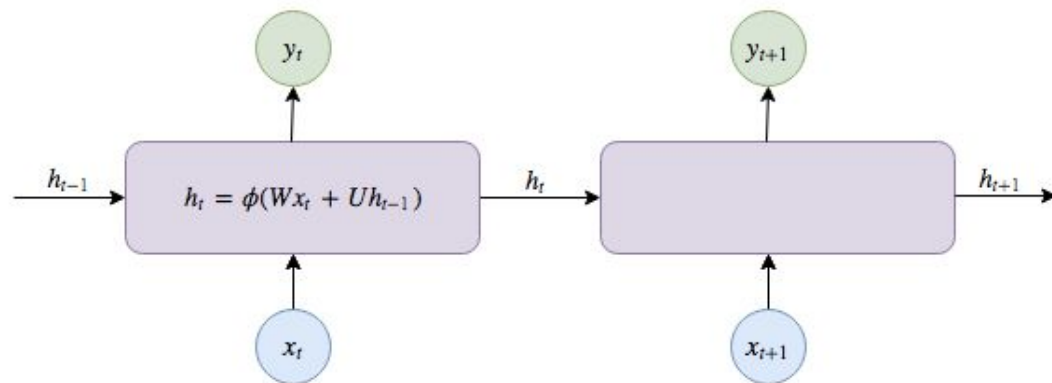
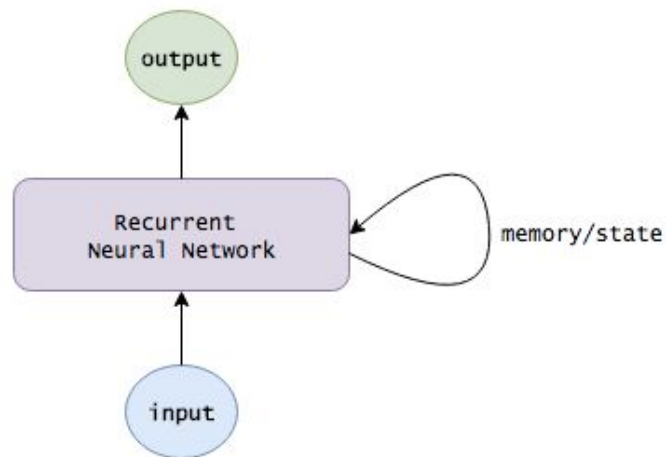


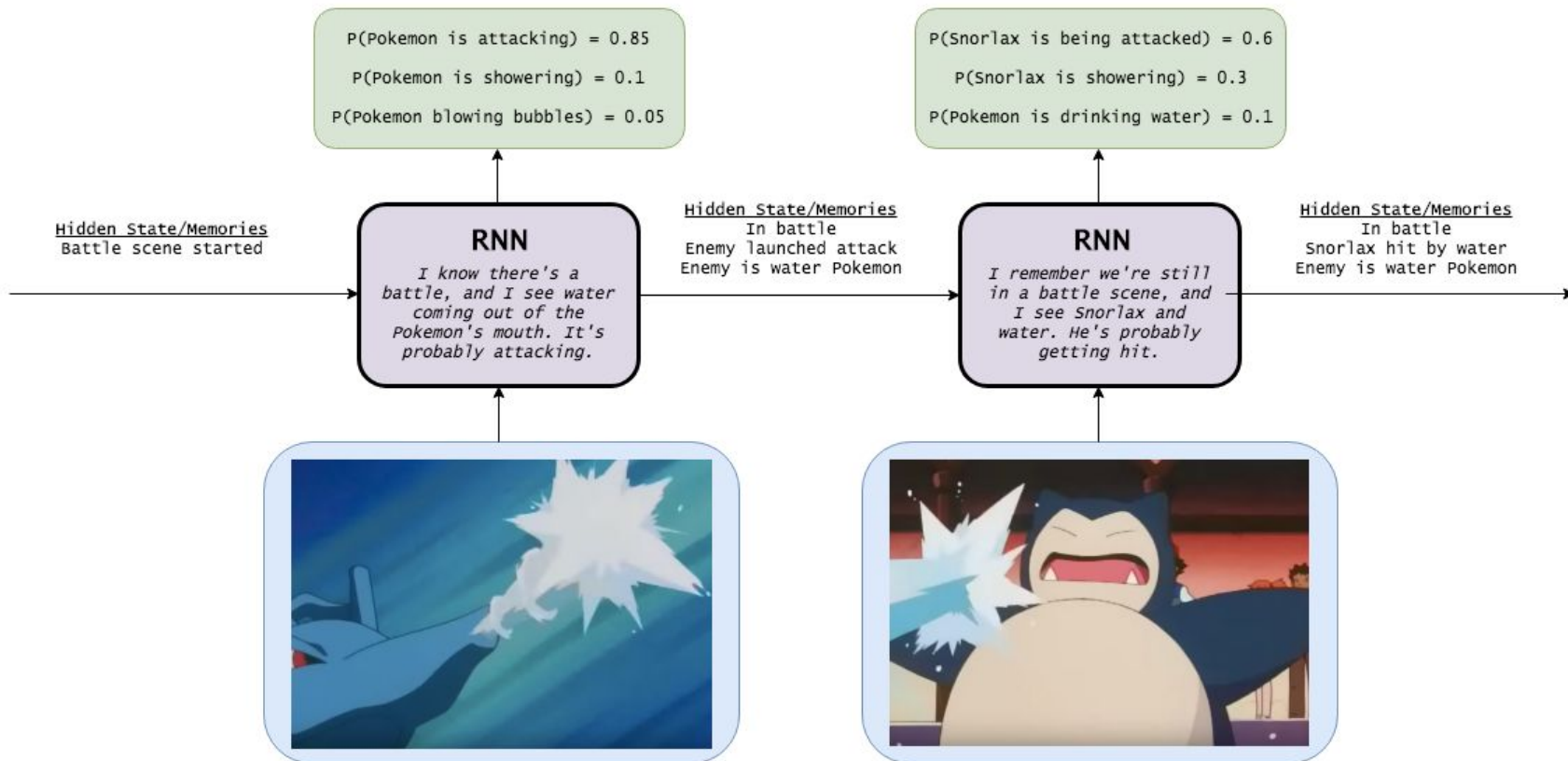
Image copyright © WildML

Y. Kim. "Convolutional Neural Networks for Sentence Classification" [2014]

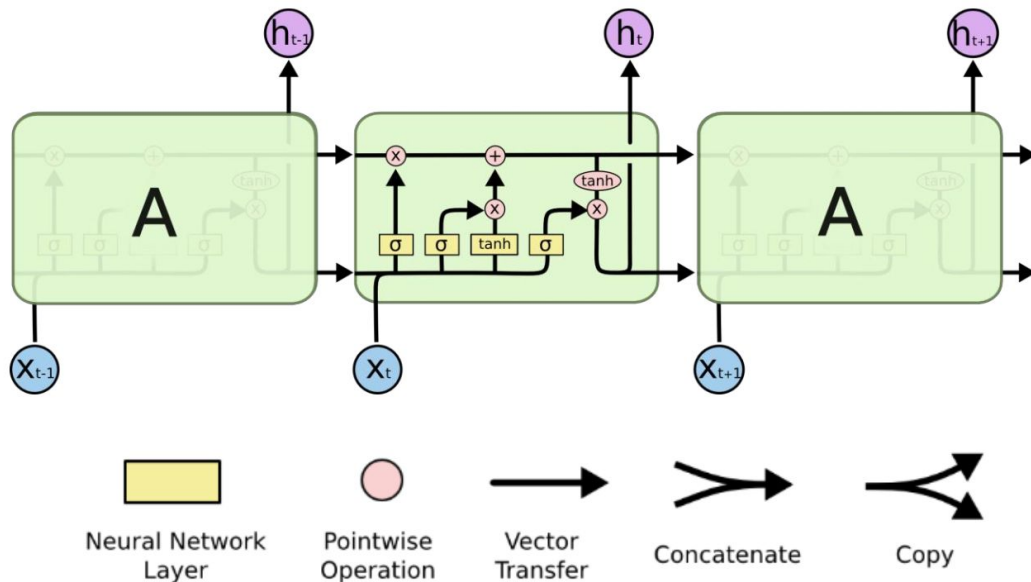
Encode: RNNs over Word Embeddings



Encode: RNNs over Word Embeddings



Encode: Gated Recurrent Architectures



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

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

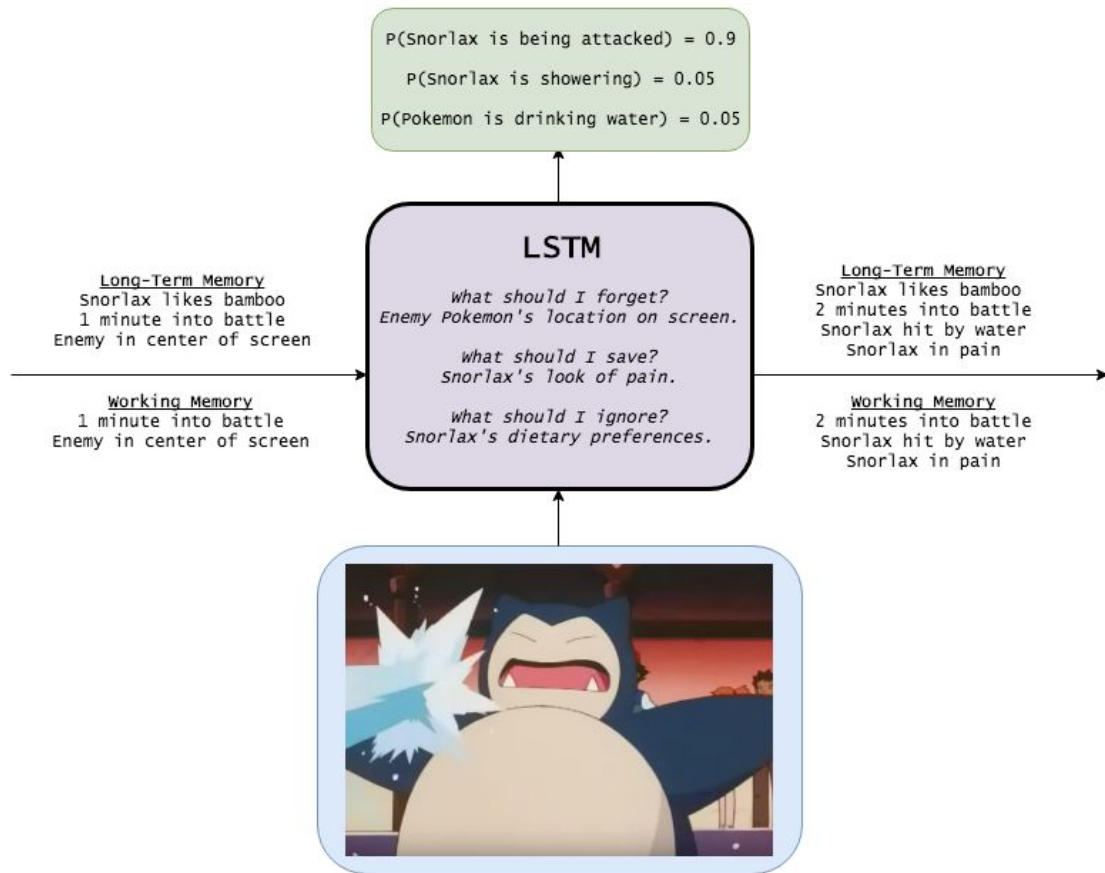
$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

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

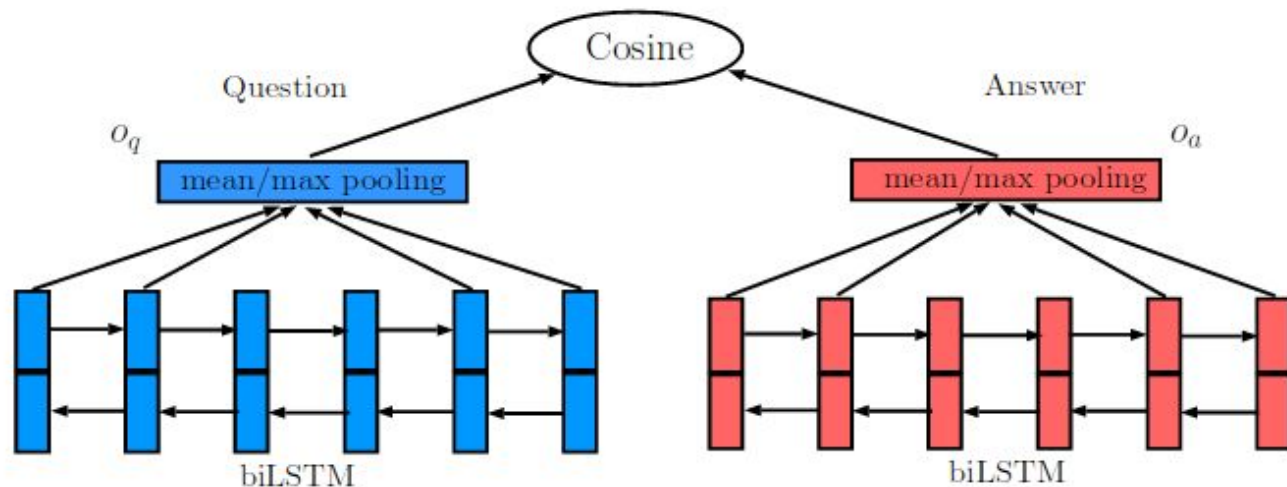
$$o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o)$$

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

Encode: Gated Recurrent Architectures



Two Input Documents? Siamese Network



Run the same encoding step for every input.

Share encoder weights, don't learn W_{E1} and W_{E2} separately.

Distance Learning, Contrastive Loss

$$L(W, Y, \vec{X}_1, \vec{X}_2) = (1 - Y) \frac{1}{2} (D_W)^2 + (Y) \frac{1}{2} \{ \max(0, m - D_W) \}^2$$

Penalize similar pairs by a monotonically **increasing** function of their learned distance.

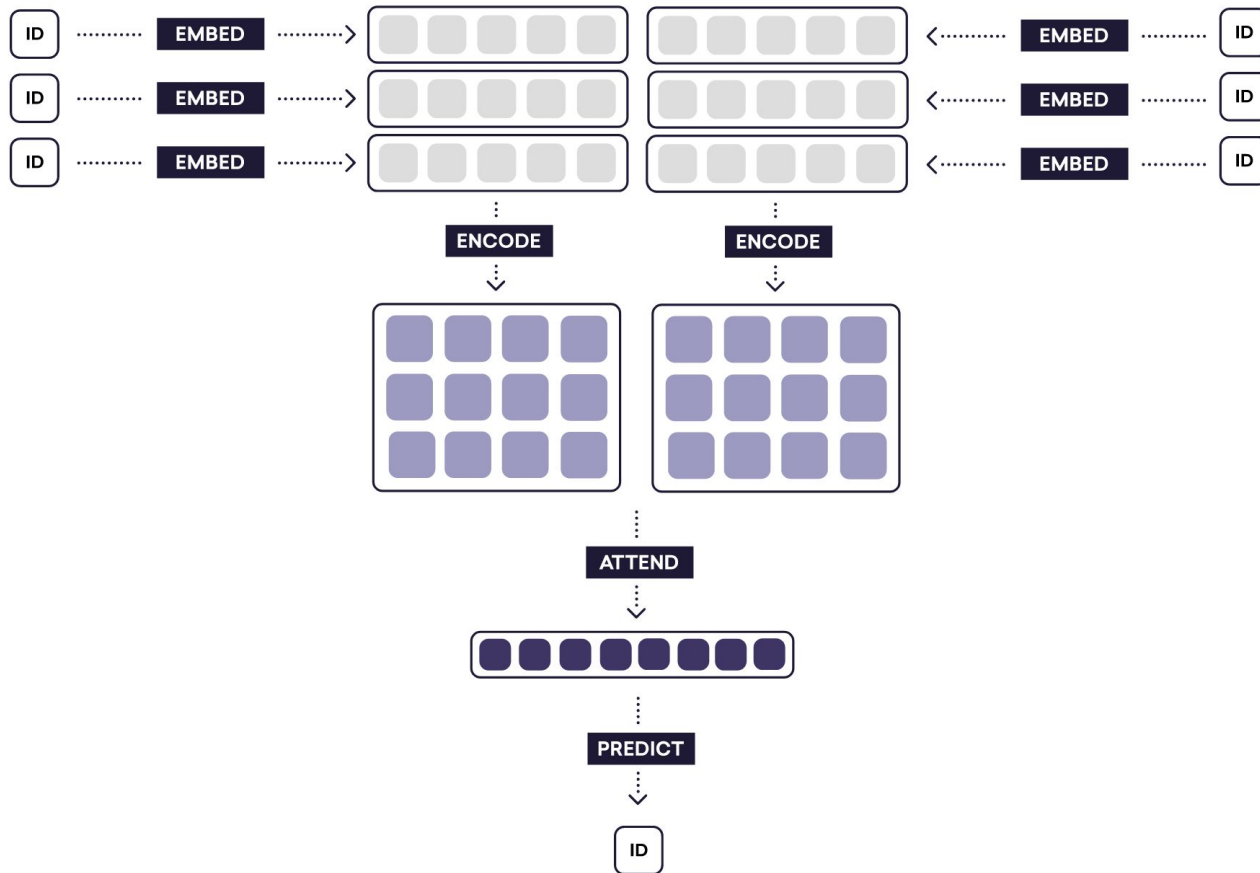
Penalize different pairs by a monotonically **decreasing** function of their learned distance.

Hadsell, Chopra, LeCun “Dimensionality Reduction by Learning an Invariant Mapping”, 2006.

<http://yann.lecun.com/exdb/publis/pdf/hadsell-chopra-lecun-06.pdf>

Sequence-to-Sequence Modeling and Attention Mechanisms in Neural NLP

State-of-the-Art NLP Pipeline



Relative Complexity of NLP Problems



Easy

- Chunking
- Part-of-Speech Tagging
- Named Entity Recognition
- Spam Detection
- Spellchecking



Medium

- Syntactic Parsing
- Paraphrase Identification
- Sentiment Analysis
- Topic Modeling
- Information Retrieval



Hard

- Machine Translation
- Text Generation
- Automatic Summarization
- Question Answering
- Conversational Interfaces

Relative Complexity of NLP Problems



Easy

- Chunking
- Part-of-Speech Tagging
- Named Entity Recognition
- Spam Detection
- Spellchecking



Medium

- Syntactic Parsing
- Paraphrase Identification
- Sentiment Analysis
- Topic Modeling
- Information Retrieval

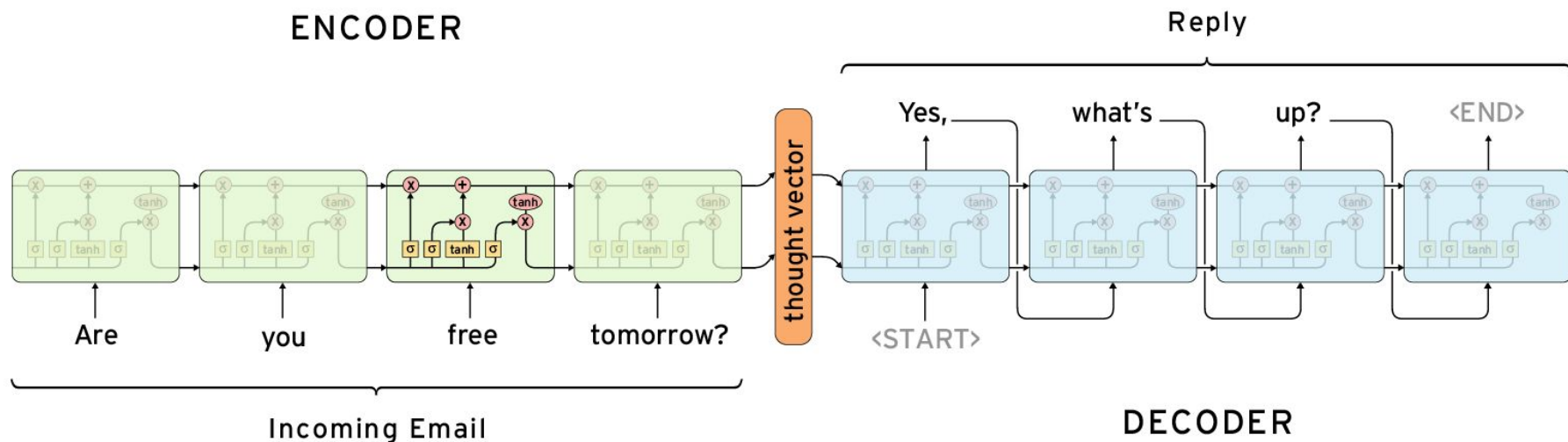


Hard

- Machine Translation
- Text Generation
- Automatic Summarization
- Question Answering
- Conversational Interfaces

Seq2Seq can be applied here

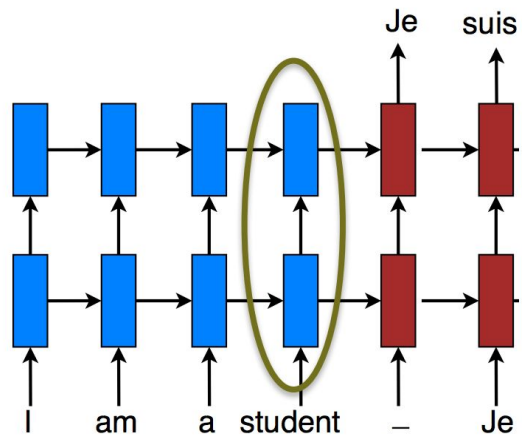
Seq2Seq Modeling at a Glance



- Two different models: an **encoder** and a **decoder** (typically RNNs).
- Encoder represents an input document as a fixed size “thought vector”.
- Decoder unwinds the thought vector into an output sequence.

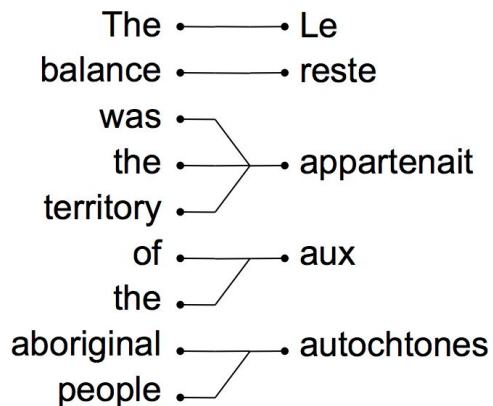
Vanilla Seq2Seq: Challenges

- Hard for encoder to compress the source documents into a fixed size vector.
- Performance deteriorates rapidly as the size of the document increases (e.g. good local grammar, but lots of word repetition).



Traditional MT: “Word Alignment”

Phrase-based SMT aligned words in a preprocessing-step, usually using EM



	Le	reste	appartenait	aux	autochtones
The					
balance					
was					
the					
territory					
of					
the					
aboriginal					
people					

Jointly Learning to Align and Translate

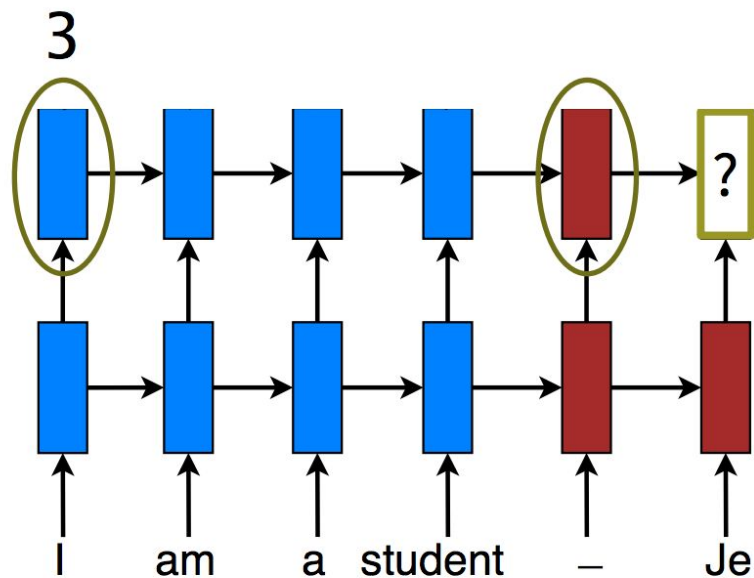
It would be a lot more beneficial to let the decoder **peek at the original input** at **each step** of the decoding phase.

D. Bahdanau, K. Cho et al. “Neural Machine Translation by Jointly Learning to Align and Translate”, 2014.

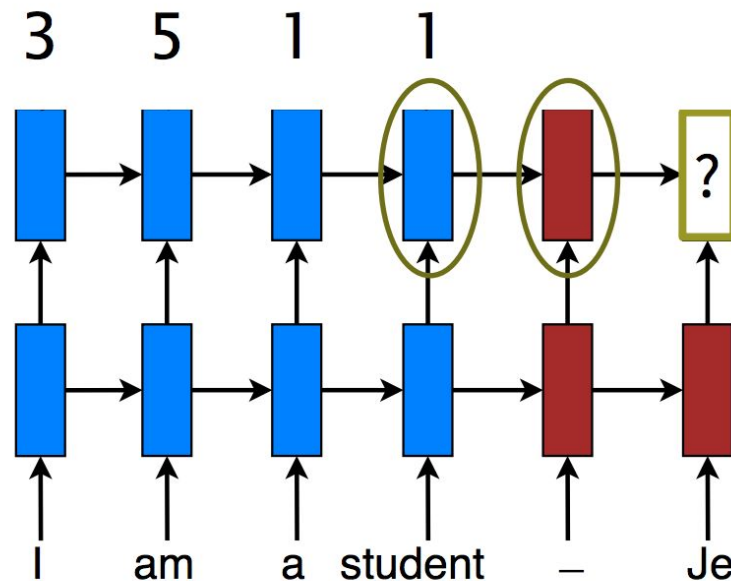
<https://arxiv.org/pdf/1409.0473.pdf>

Neural Attention: Scoring

$$\text{score}(\mathbf{h}_{t-1}, \bar{\mathbf{h}}_s)$$



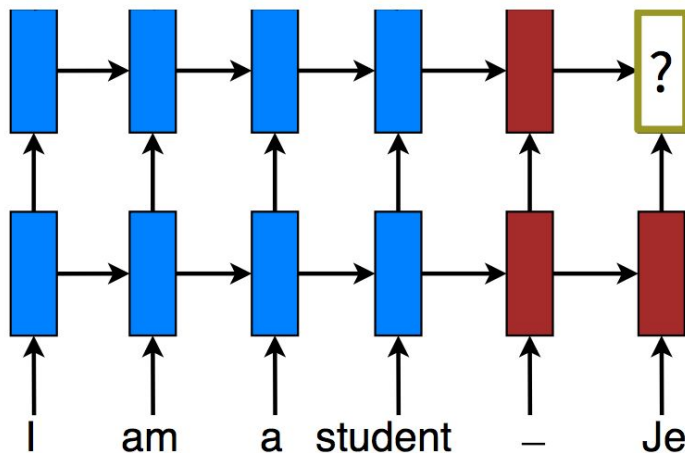
$$\text{score}(\mathbf{h}_{t-1}, \bar{\mathbf{h}}_s)$$



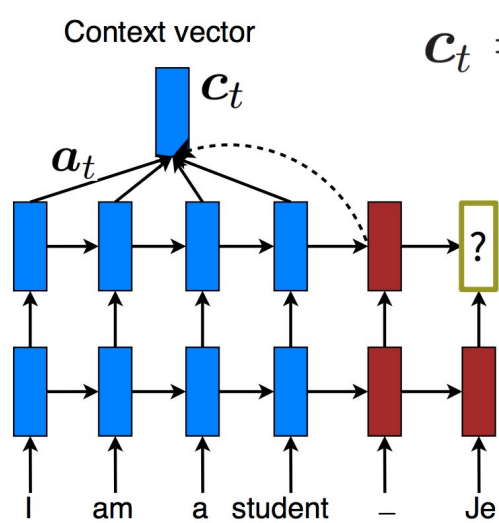
Neural Attention: Scoring

$$\mathbf{a}_t(s) = \frac{e^{\text{score}(s)}}{\sum_{s'} e^{\text{score}(s')}}$$

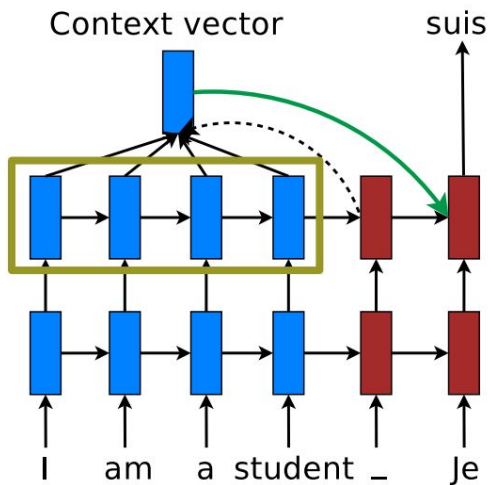
\mathbf{a}_t 0.3 0.5 0.1 0.1



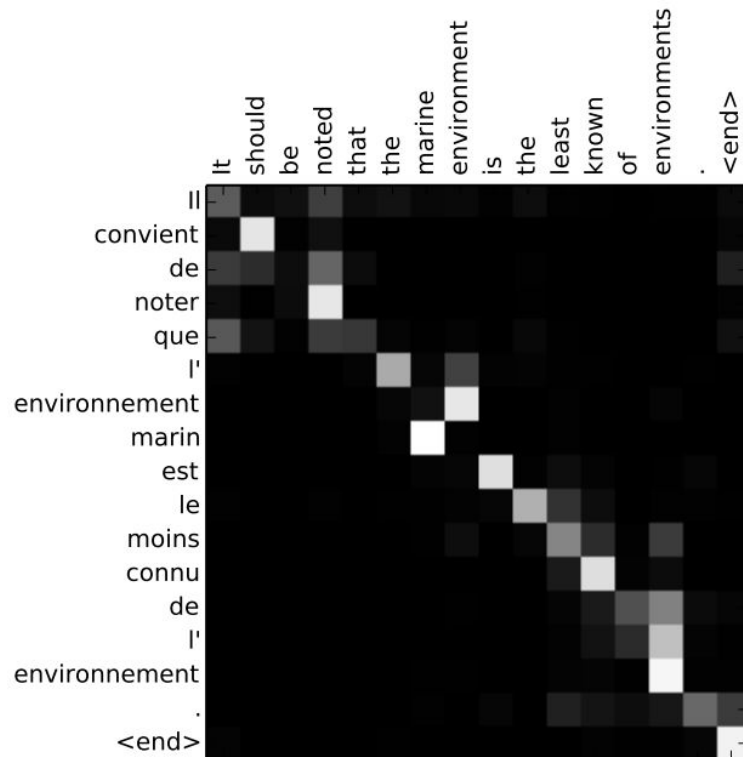
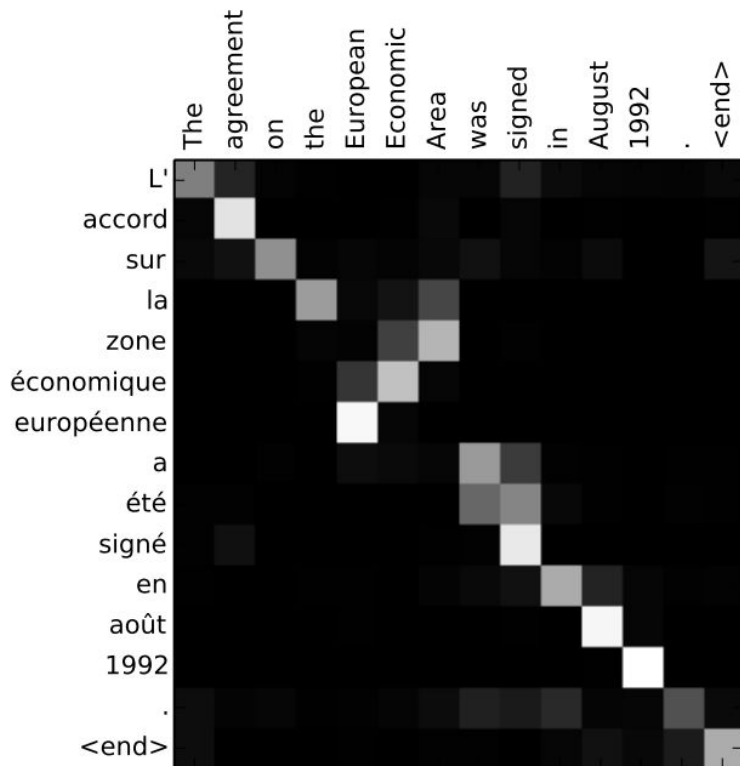
Attention Mechanism: Context



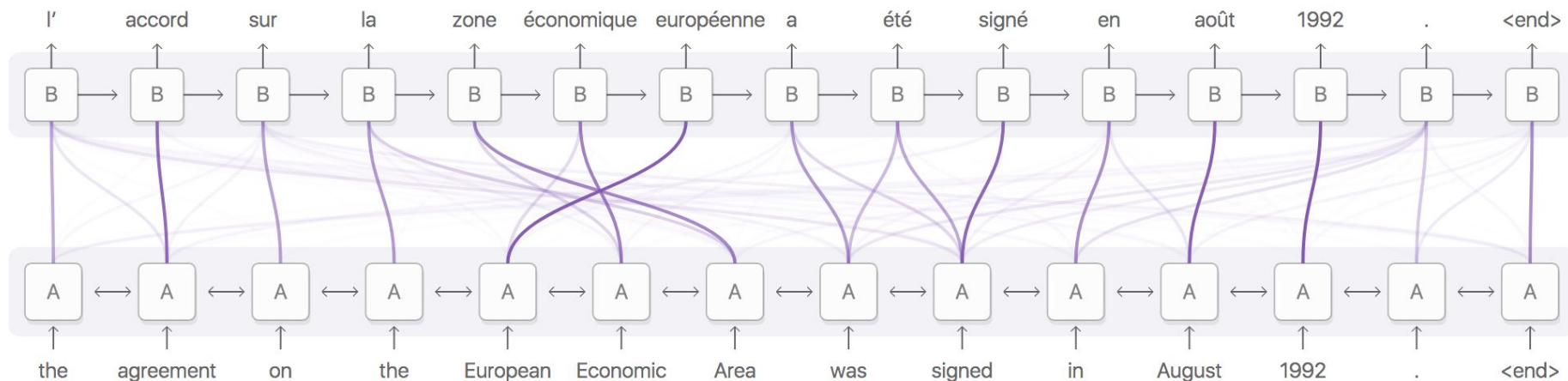
$$c_t = \sum_s a_t(s) \bar{h}_s$$



Interpretability of Attention Scores

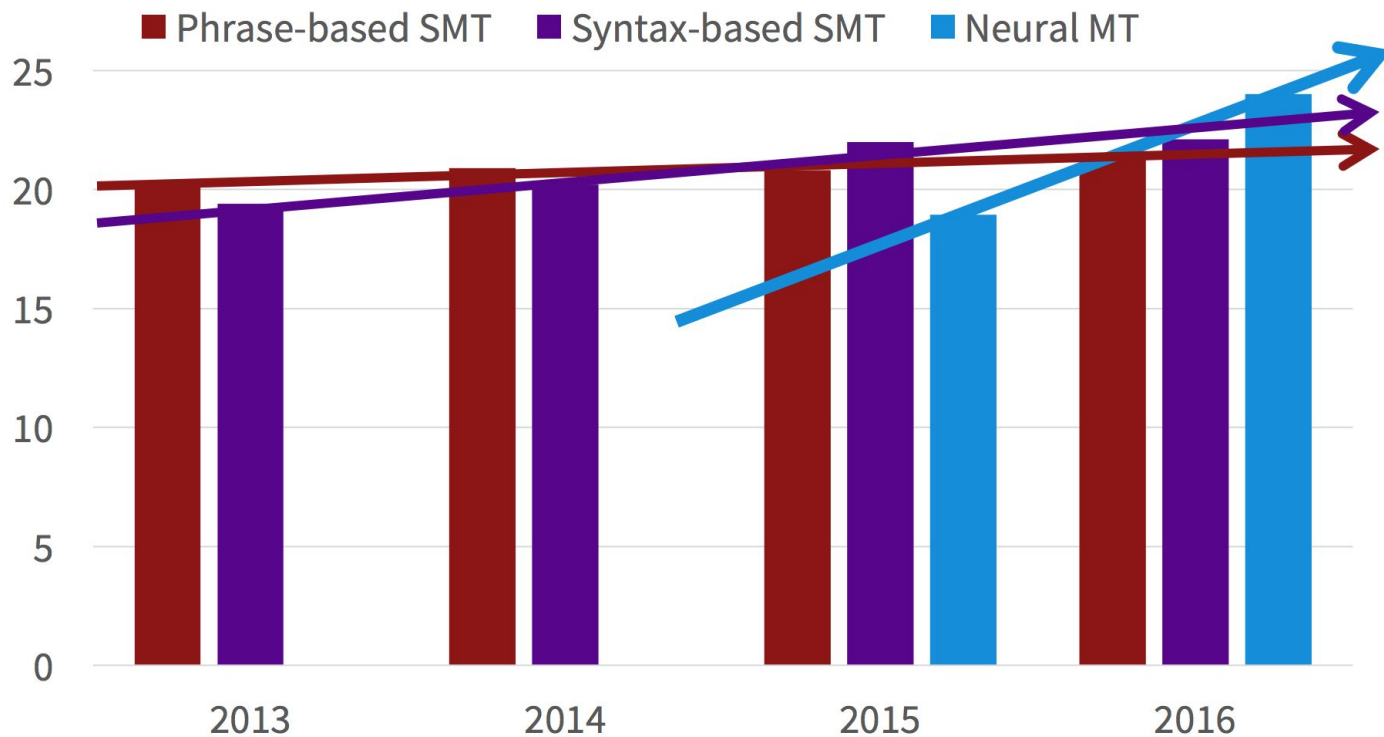


Interpretability of Attention Scores



MT Performance

[Edinburgh En-De WMT newstest2013 Cased BLEU; NMT 2015 from U. Montréal]



Attention-Only NMT

Recent paper: SotA NMT using only stacked attention units and positional embeddings (no CNN, RNN etc.).
Very computationally efficient.

Vaswani, Shazeer et al. (Google Brain)
“Attention is All You Need”, Jun 2017.

<https://arxiv.org/pdf/1706.03762.pdf>

