

Natural Language Processing

Sequence to Sequence Learning

Instructor: Yangqiu Song

Sequence to Sequence

- Speech recognition



<http://nlp.stanford.edu/courses/lsa352/>

Sequence to Sequence

- Question answering



Sequence to Sequence

- Machine translation

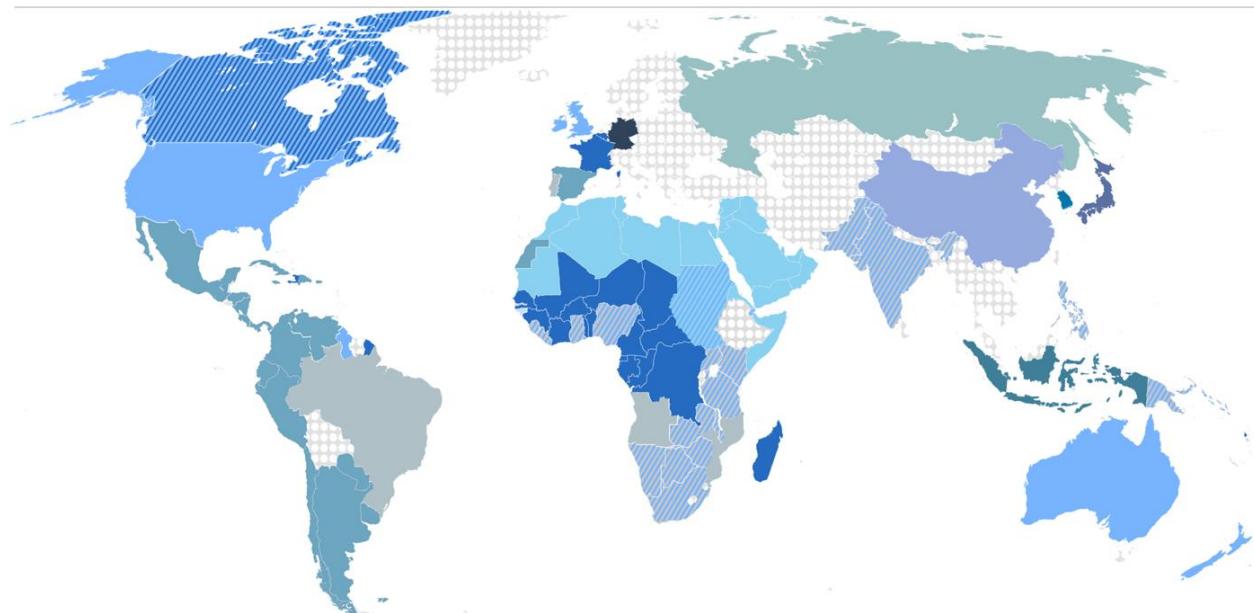
مرحبا بكم في درس التعلم العميق



Welcome to the deep learning class

7 billion people, 7000 languages

Top Languages on the Internet



Number of Internet users by Language - mln people
The bars' heights correspond with the figure



Internet Penetration by Language

English - 43%
Chinese - 37%
Spanish - 39%
Japanese - 78%
Portuguese - 32%
German - 79%
Arabic - 18%
French - 17%
Russian - 42%
Korean - 55%
Indonesian - 16%

World population by Language (mln)

English - 1302
Chinese - 1372
Spanish - 423
Japanese - 126
Portuguese - 253
German - 94
Arabic - 347
French - 347
Russian - 139
Korean - 71
Indonesian - 245

Source: Internet World Stats

Machine Translation

- Machine Translation (MT) is the task of translating a sentence x from one language (the source language) to a sentence y in another language (the target language).
- 1950s: Early Machine Translation
 - Mostly Russian → English (motivated by the Cold War!)
 - Systems were mostly rule-based, using a bilingual dictionary to map Russian words to their English counterparts



Source: <https://youtu.be/K-HfpsHPmvw>

1990s-2010s: Statistical Machine Translation

- Core idea: Learn a probabilistic model from data
- Suppose we're translating Language X → English.
- We want to find best English sentence y , given Language X sentence x

$$\operatorname{argmax}_y P(y|x)$$

- Use Bayes Rule to break this down into two components to be learnt separately:

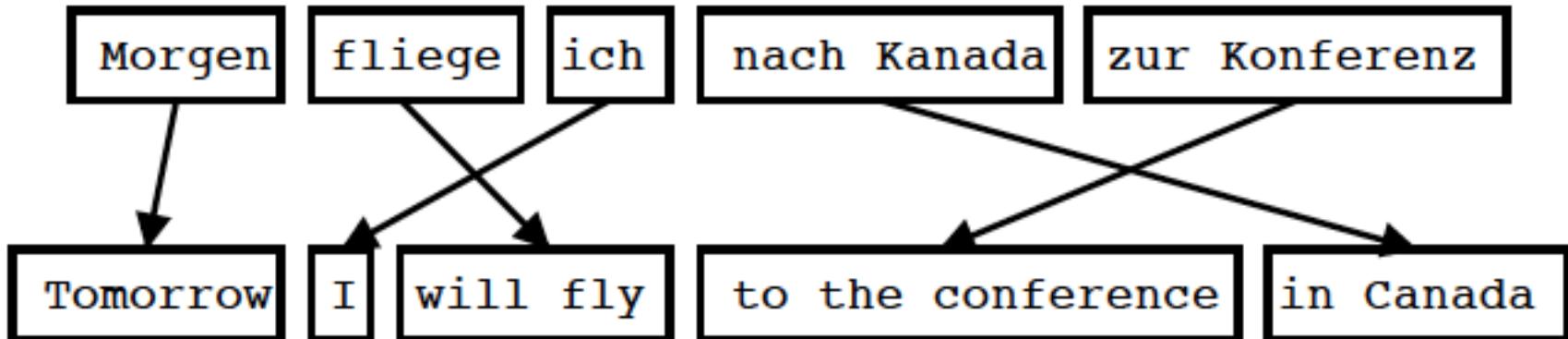
$$= \operatorname{argmax}_y P(x|y)P(y)$$

 **Translation Model**
Models how words and phrases
should be translated.
Learnt from parallel data.

 **Language Model**
Models how to write good English.
Learnt from monolingual data.

Statistical Machine Translation

- Translation model
- Input is Segmented in Phrases
- Each Phrase is Translated into English
- Phrases are Reordered



Statistical Machine Translation

- Language Model

Goal of the Language Model: Detect good English

For Example: Trigram Model

```
Mary did not slap the green witch
```

```
Mary => p(Mary)
```

```
Mary did => p(did|Mary)
```

```
Mary did not => p(not|Mary did)
```

```
did not slap => p(slap|did not)
```

```
not slap the => p(the|not slap)
```

```
slap the green => p(green|slap the)
```

```
the green witch => p(witch|the green)
```

1990s-2010s: Statistical Machine Translation

$$\operatorname{argmax}_y P(x|y)P(y)$$

Question: Translation Model Language Model
How to compute
this argmax?

- We could enumerate every possible y and calculate the probability? → Too expensive!
- **Answer:** Use a heuristic search algorithm to gradually build up the translation, discarding hypotheses that are too low probability

Statistical Machine Translation

- Decoding

Goal of the decoding algorithm: Put models to work, perform the actual translation

Maria	no	dio	una	bofetada	a	la	bruja	verde
<u>Mary</u>	<u>not</u>	<u>give</u>	<u>a</u>	<u>slap</u>	<u>to</u>	<u>the</u>	<u>witch</u>	<u>green</u>
	<u>did not</u>			<u>a slap</u>		<u>by</u>		<u>green witch</u>
	<u>no</u>		<u>slap</u>			<u>to the</u>		
	<u>did not give</u>					<u>to</u>		
				<u>slap</u>			<u>the witch</u>	

```
e:  
f: -----  
p: 1
```

Statistical Machine Translation

- Decoding

Goal of the decoding algorithm: Put models to work, perform the actual translation

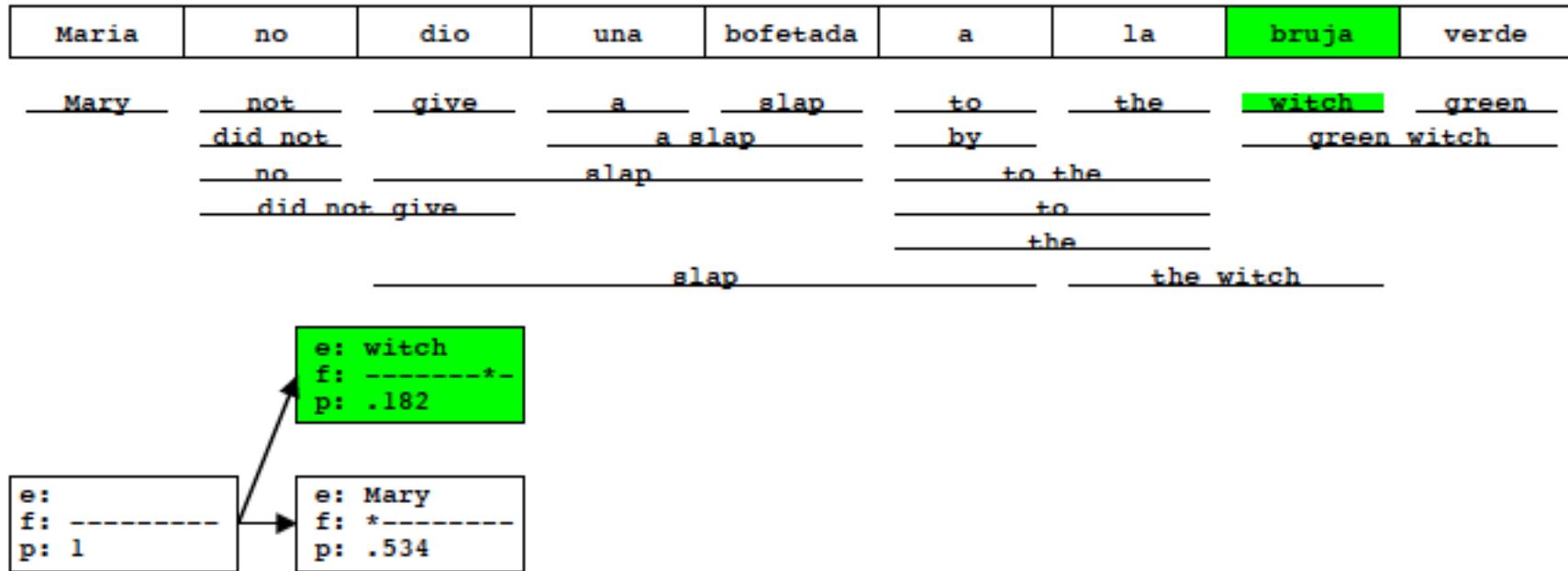
Maria	no	dio	una	bofetada	a	la	bruja	verde
Mary	not	give	a	slap	to	the	witch	green
	did not		a slap		by		green	witch
	no		slap		to the			
	did not give				to			
				slap		the		
						the	witch	



Statistical Machine Translation

- Decoding

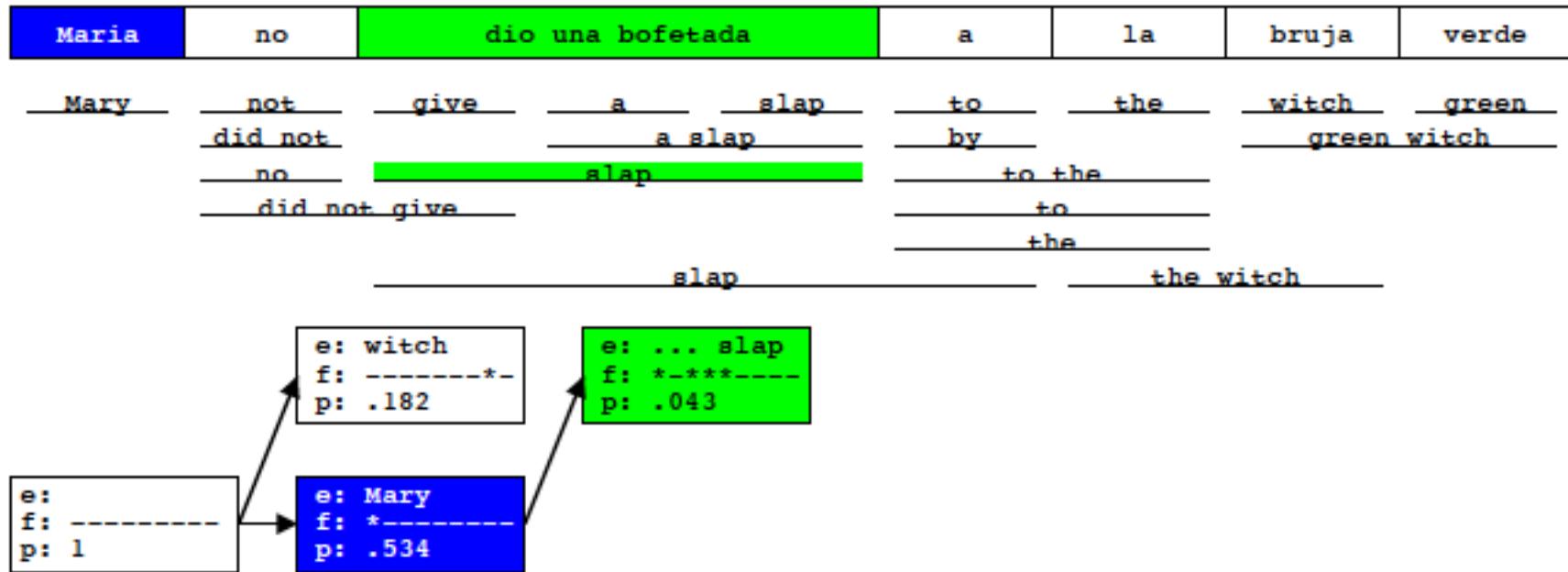
Goal of the decoding algorithm: Put models to work, perform the actual translation



Statistical Machine Translation

- Decoding

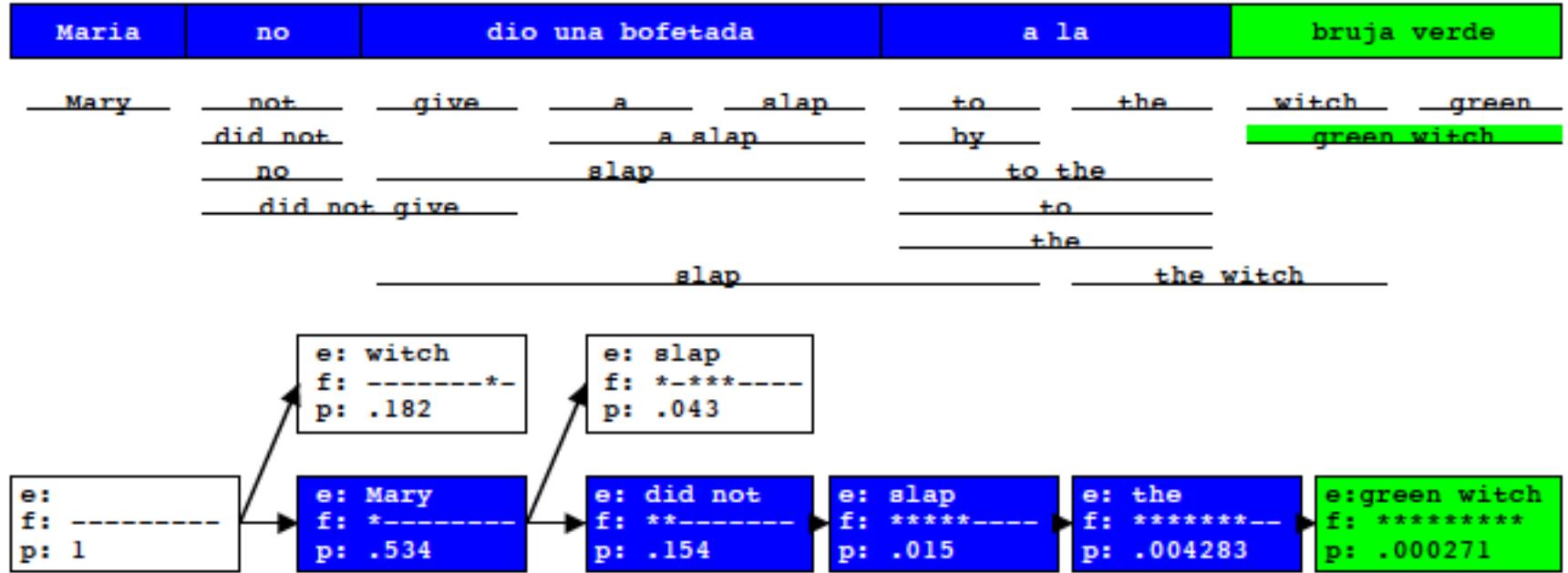
Goal of the decoding algorithm: Put models to work, perform the actual translation



Statistical Machine Translation

- Decoding

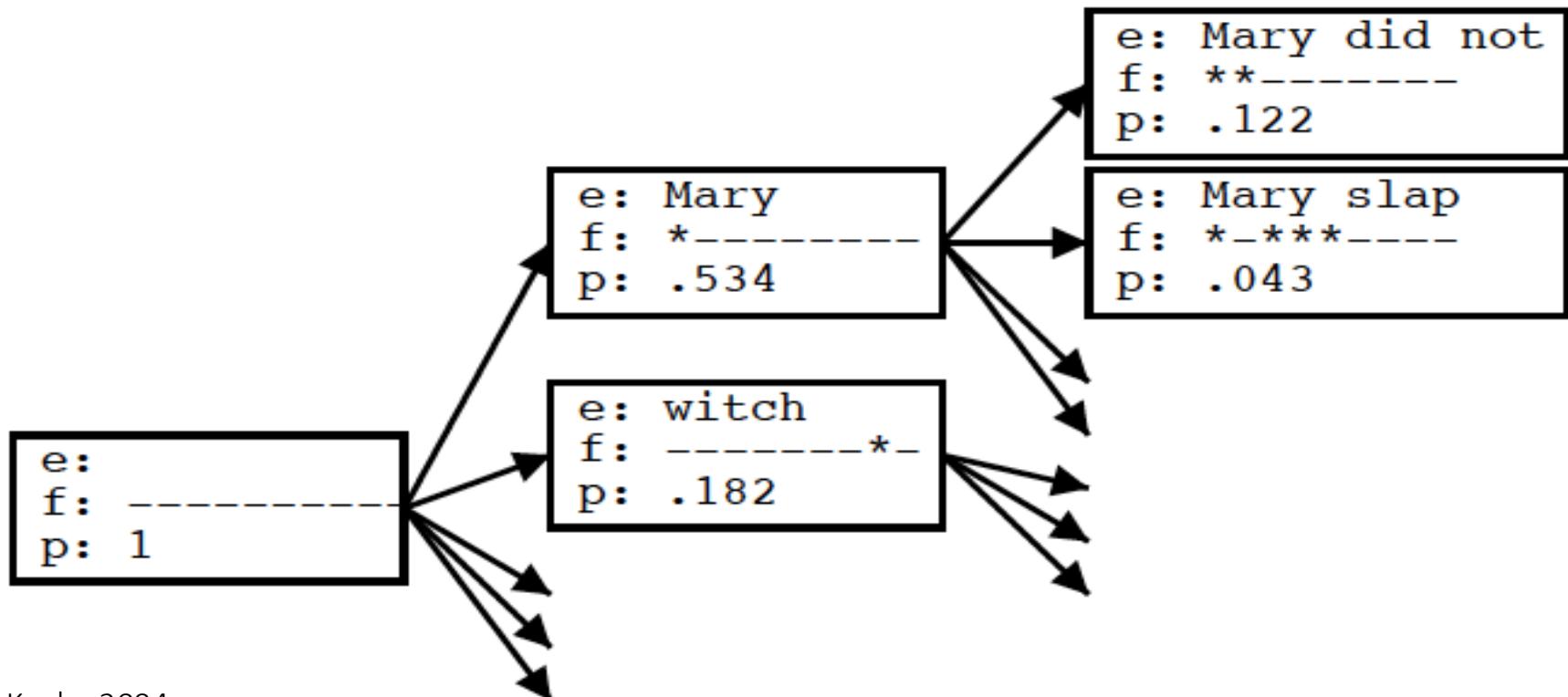
Goal of the decoding algorithm: Put models to work, perform the actual translation



Statistical Machine Translation

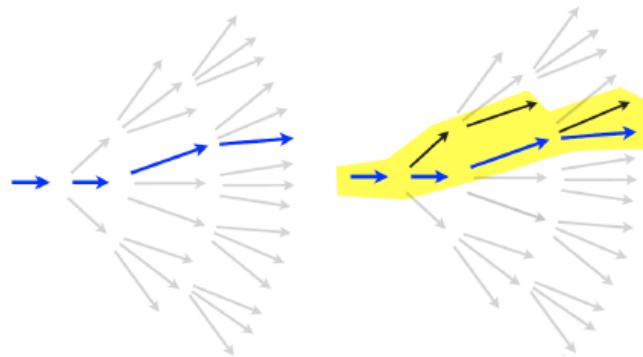
- Decoding

Goal of the decoding algorithm: Put models to work, perform the actual translation



Decoding

- Global solution: Dynamic programming
- Approximate solutions: Beam inference
 - At each position keep the top k complete sequences
 - Extend each sequence in each local way
 - The extensions compete for the k slots at the next position



- Advantages
 - Fast; beam sizes of 3-5 are almost as good as exact inference in many cases
 - Easy to implement (no dynamic programming required)
- Disadvantage
 - Inexact: the globally best sequence can fall off the beam

1990s-2010s: Statistical Machine Translation

- SMT is a huge research field
- The best systems are **extremely complex**
 - Hundreds of important details we haven't mentioned here
 - Systems have many **separately-designed subcomponents**
 - Lots of **feature engineering**
 - Need to design features to capture particular language phenomena
 - Require compiling and maintaining **extra resources**
 - Like tables of equivalent phrases
 - Lots of **human effort** to maintain
 - Repeated effort for each language pair!

Neural Machine Translation

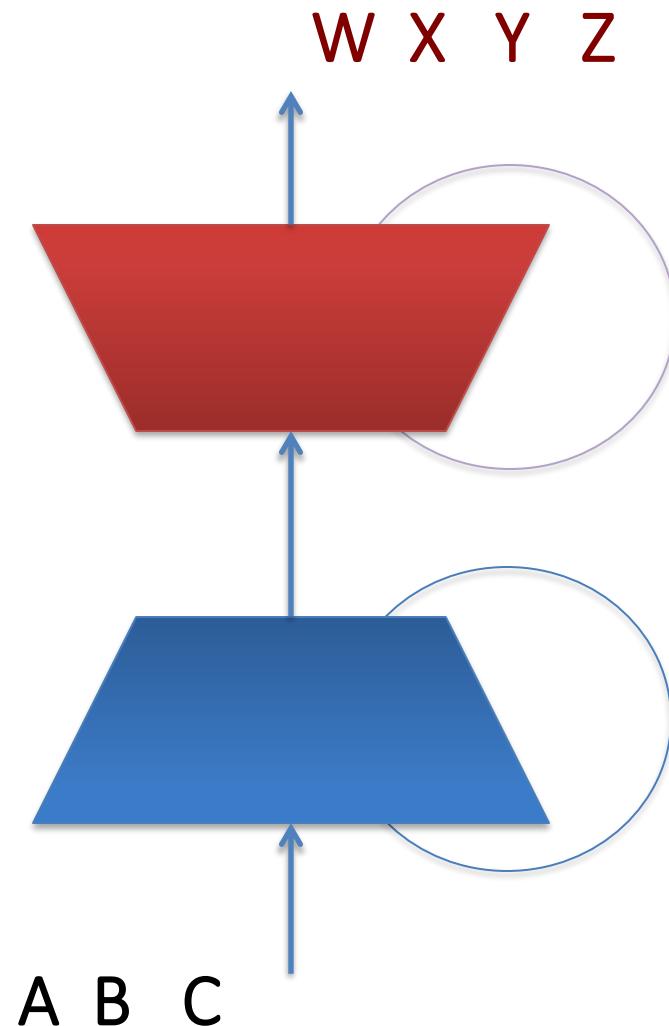
- Neural Machine Translation (NMT) is a way to do Machine Translation with a *single neural network*
- The neural network architecture is called *sequence-to-sequence* (aka seq2seq) and it involves *two RNNs*.

Sequence to Sequence Learning with Neural Networks

- Ilya Sutskever, Oriol Vinyals, Quoc V. Le

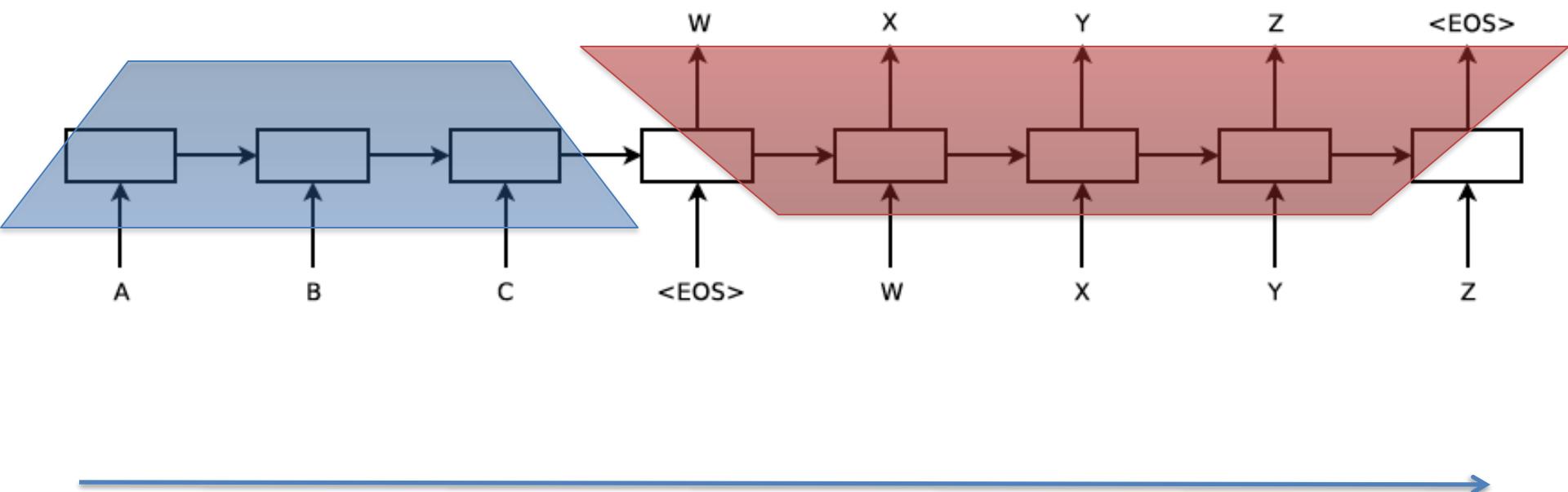
Neural Machine Translation

- Model



Neural Machine Translation

- Model



Neural Machine Translation (NMT)

- The sequence-to-sequence model is an example of a **Conditional Language Model**.
 - **Language Model** because the decoder is predicting the next word of the target sentence y
 - **Conditional** because its predictions are *also* conditioned on the source sentence x
- NMT directly calculates $P(y|x)$

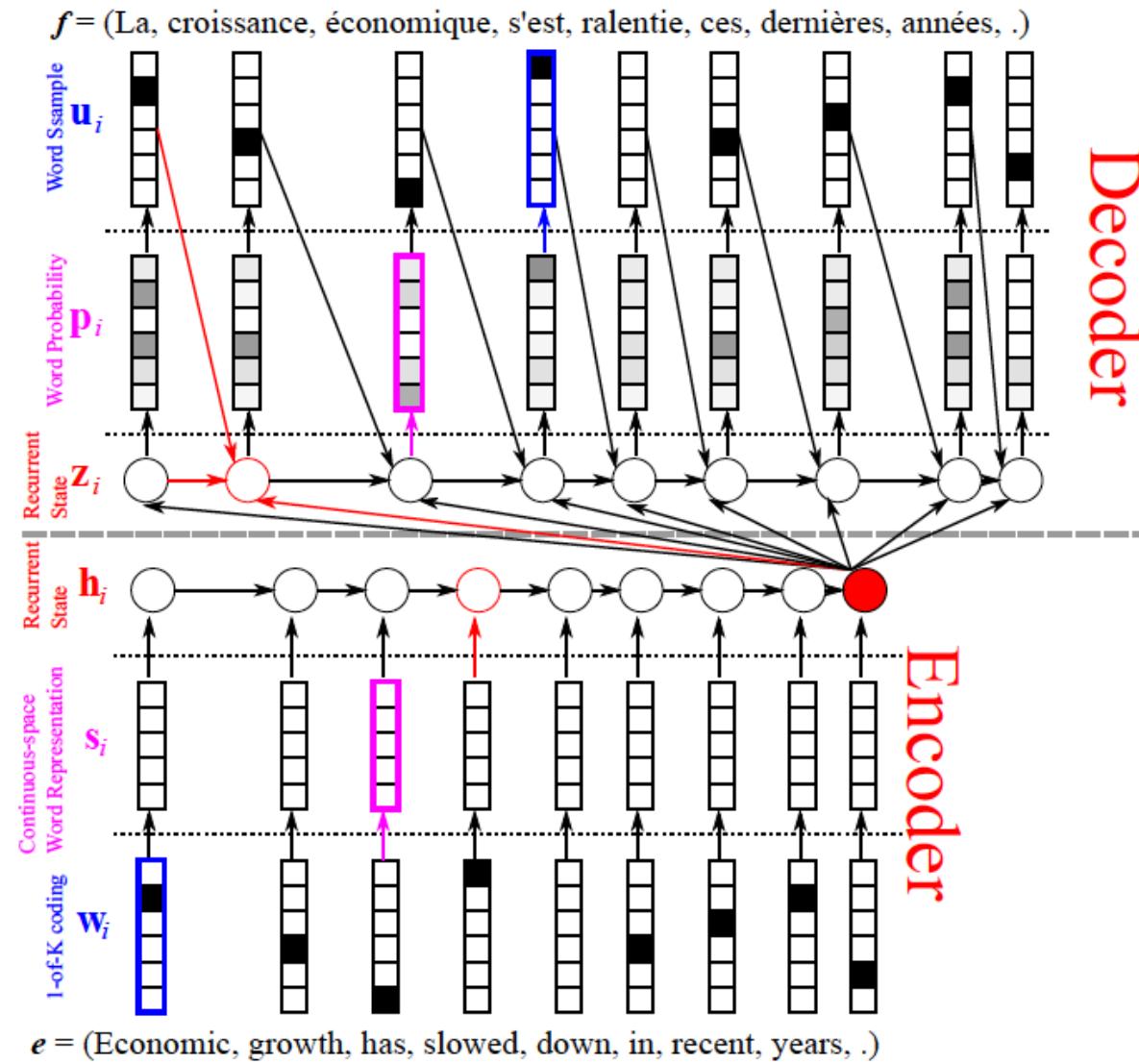
$$P(y|x) = P(y_1|x) P(y_2|y_1, x) P(y_3|y_1, y_2, x) \dots, \underline{P(y_T|y_1, \dots, y_{T-1}, x)}$$

- **Question:** How to train a NMT system?
- **Answer:** Get a big parallel corpus...

Probability of next target word, given target words so far and source sentence x

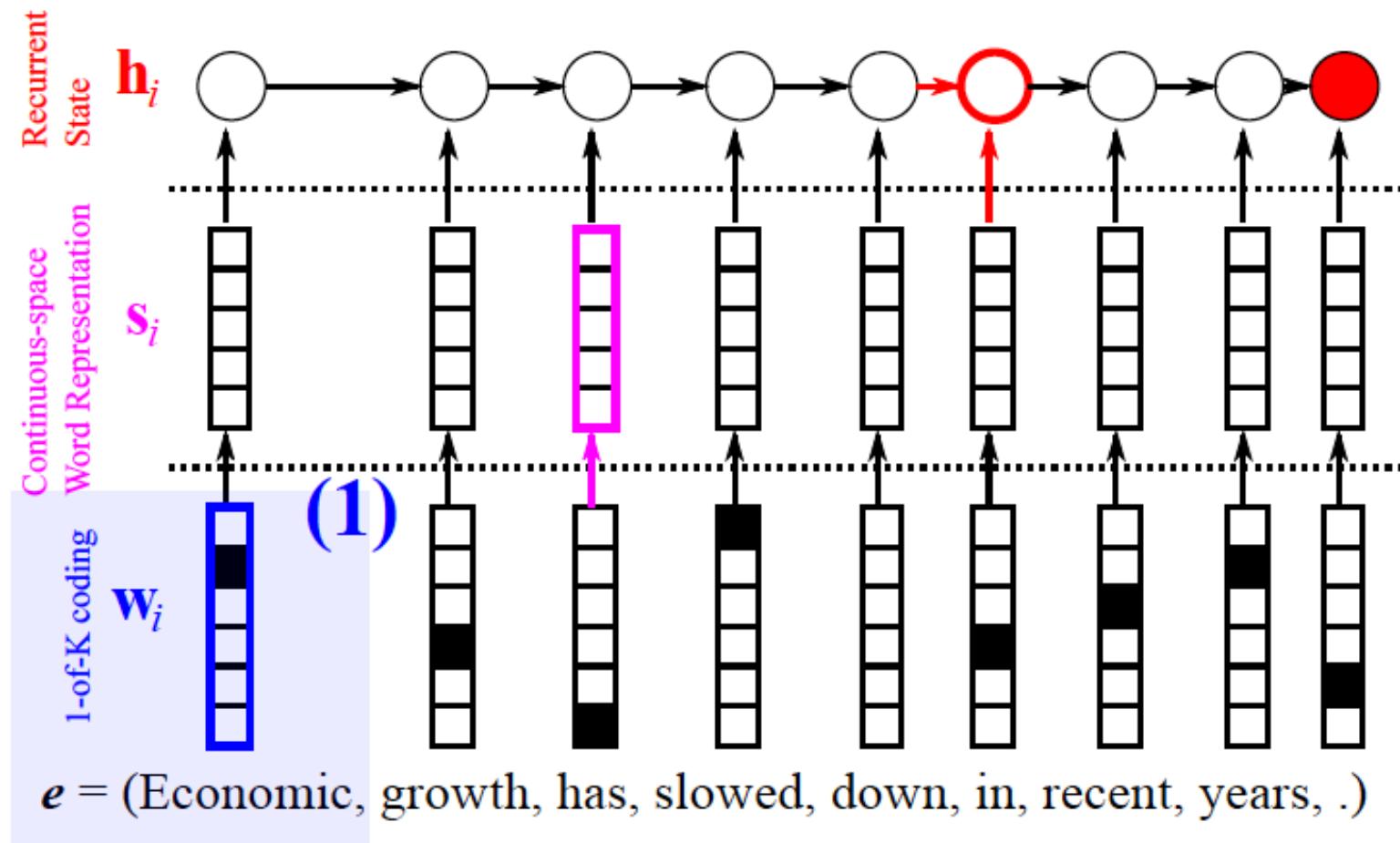
Neural Machine Translation

- Model



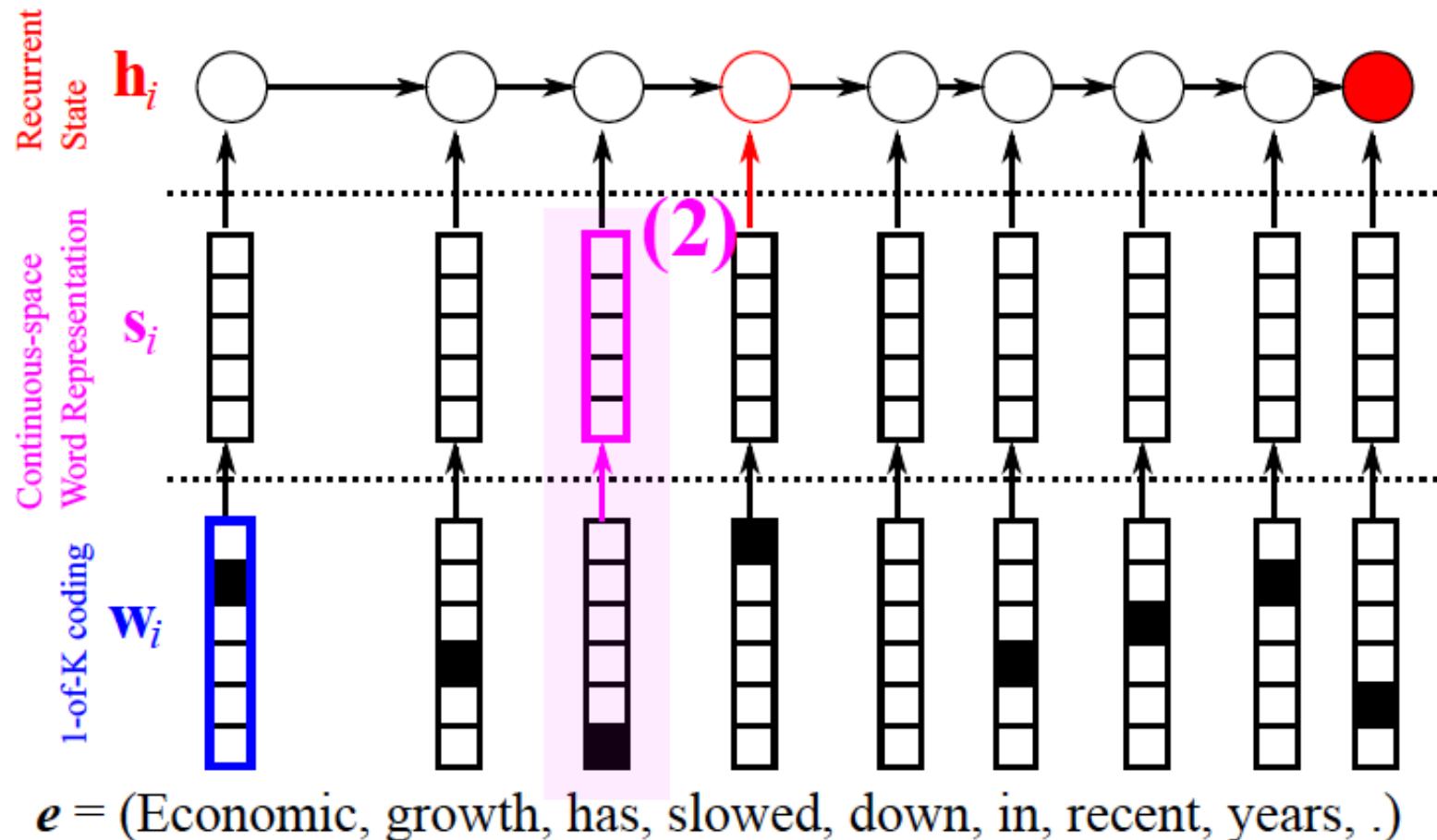
Neural Machine Translation

- Model- *encoder*



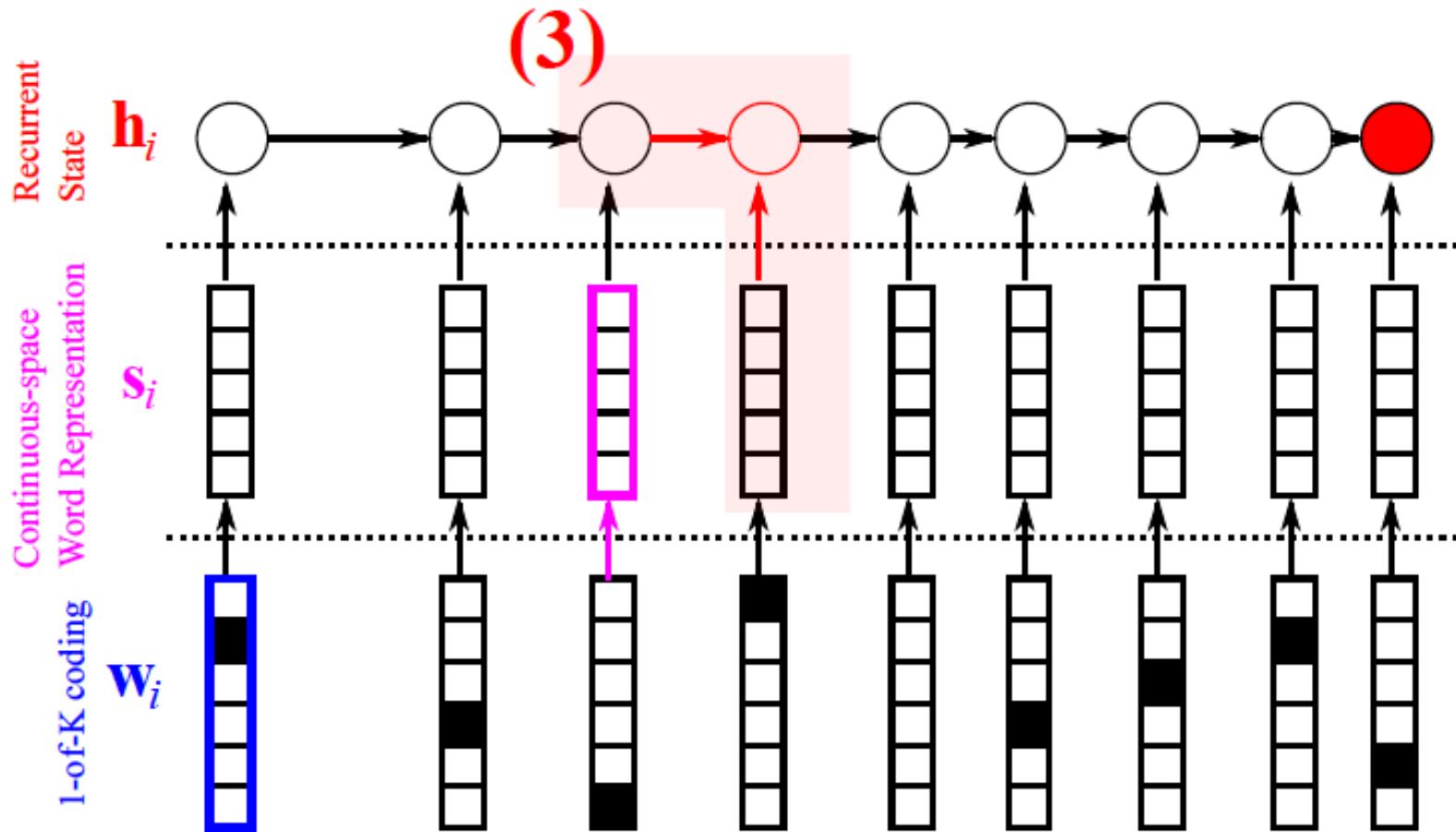
Neural Machine Translation

- Model- *encoder*



Neural Machine Translation

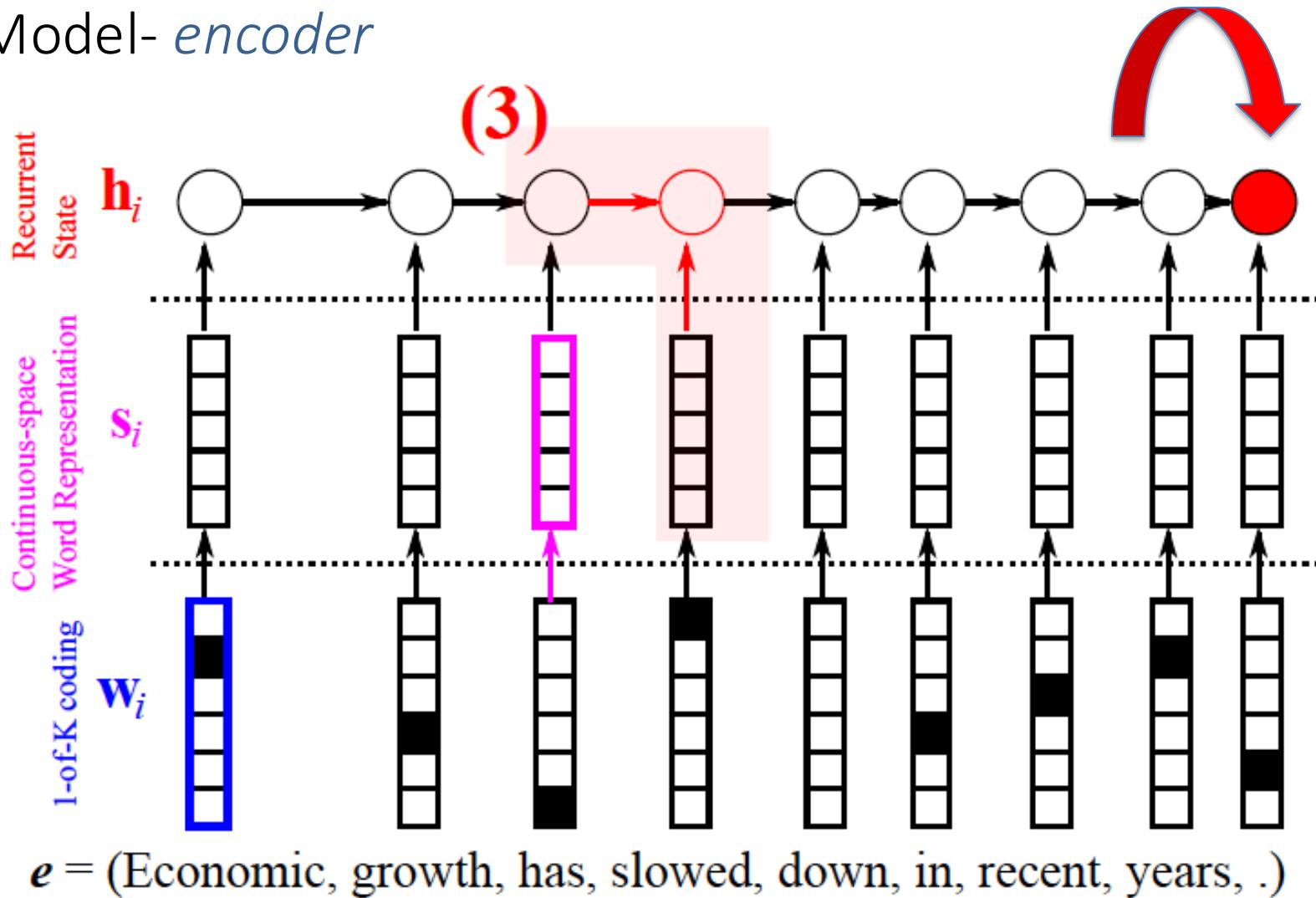
- Model- *encoder*



$e = (\text{Economic, growth, has, slowed, down, in, recent, years, .})$

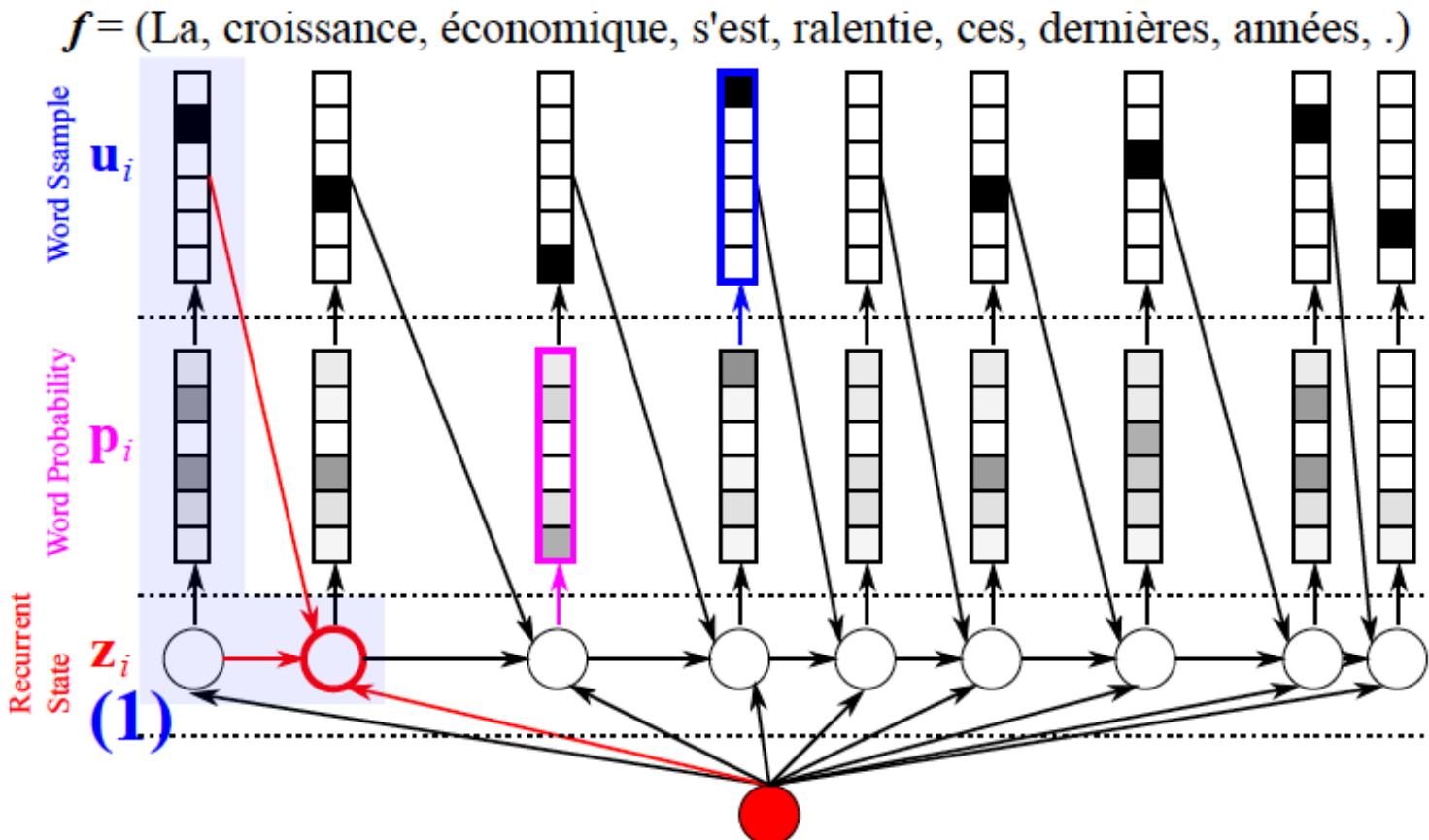
Neural Machine Translation

- Model- *encoder*



Neural Machine Translation

- Model- *decoder*

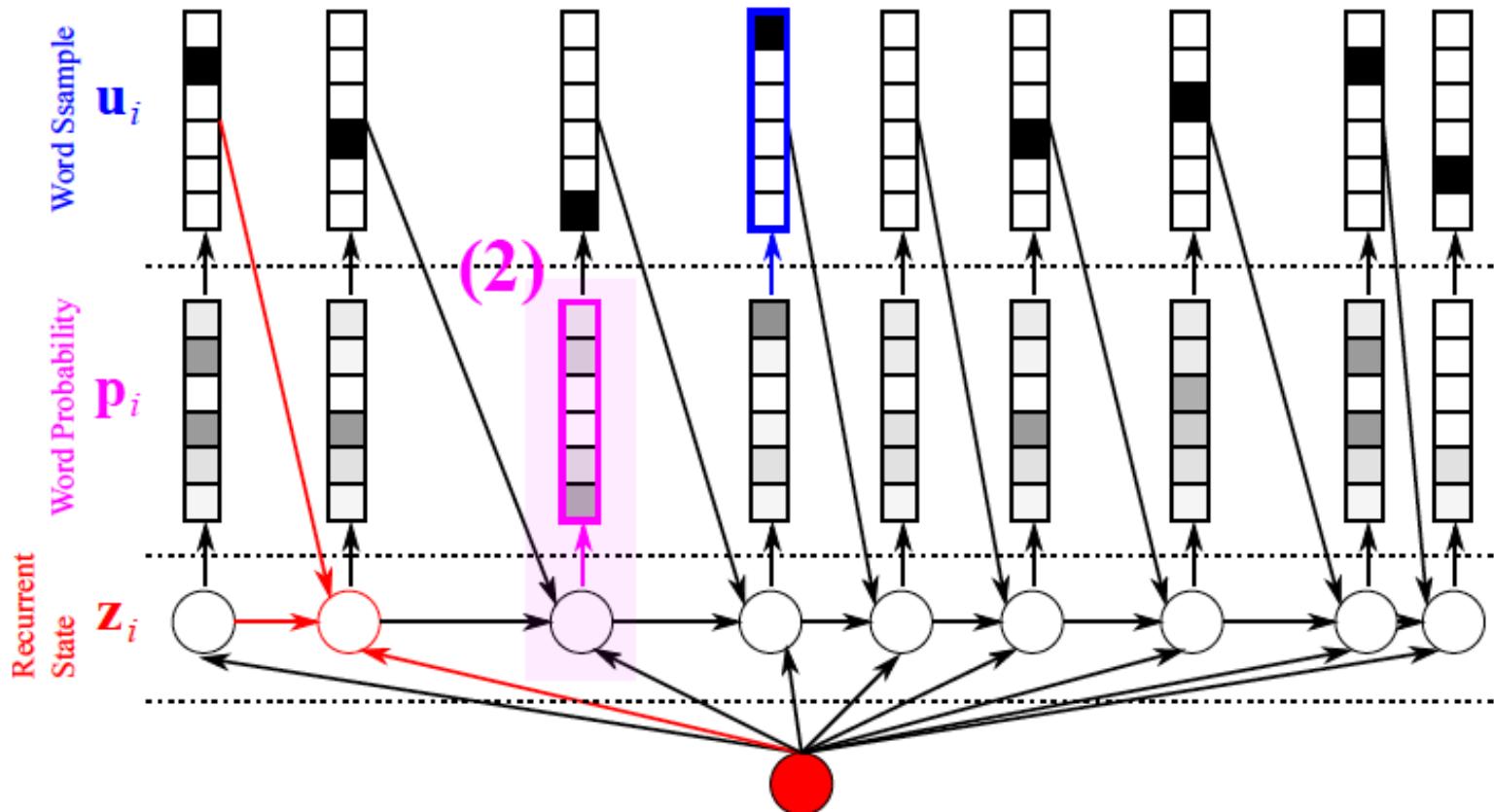


$e = (\text{Economic, growth, has, slowed, down, in, recent, years, .})$

Neural Machine Translation

- Model- *decoder*

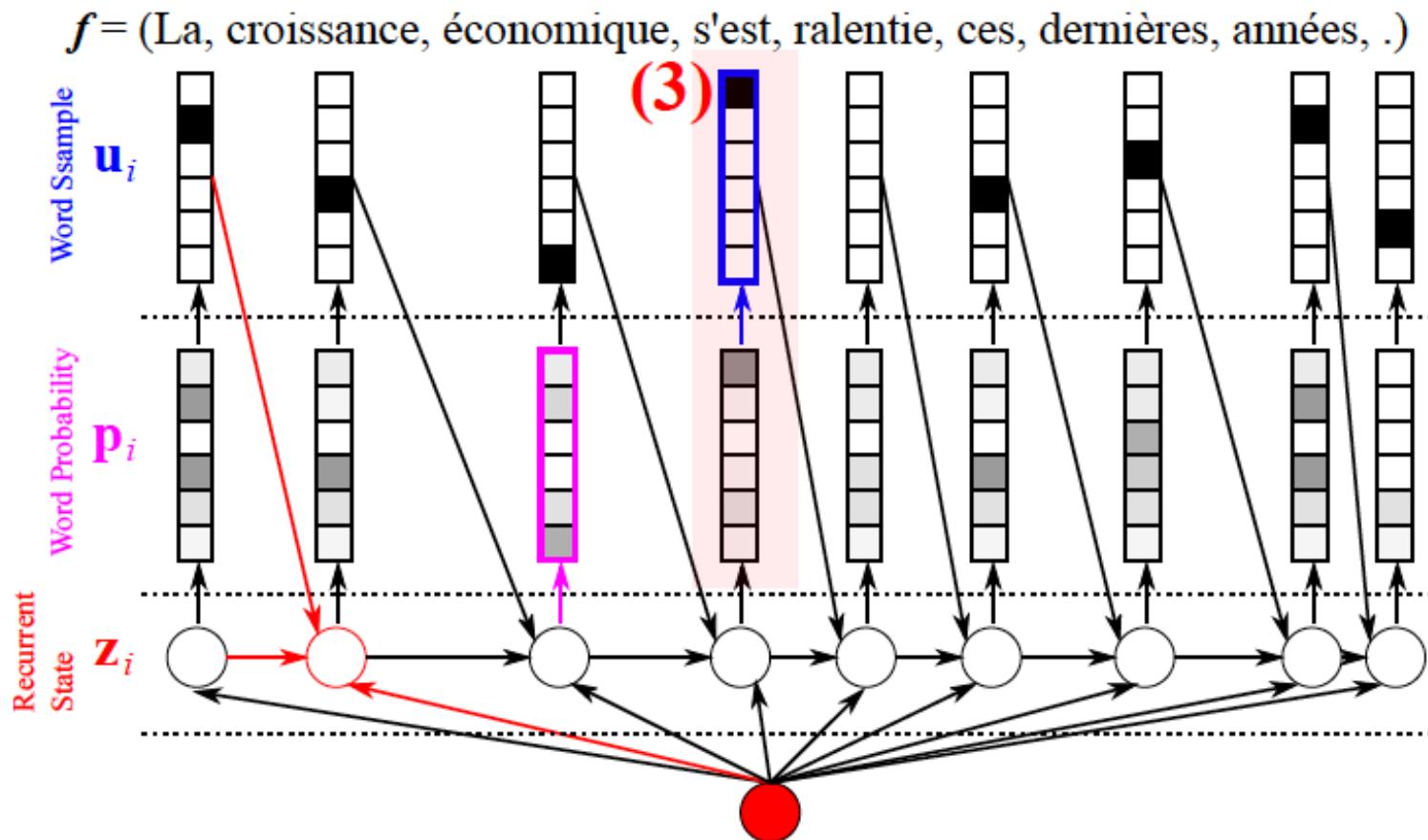
$f = (\text{La, croissance, économique, s'est, ralenti, ces, dernières, années, .})$



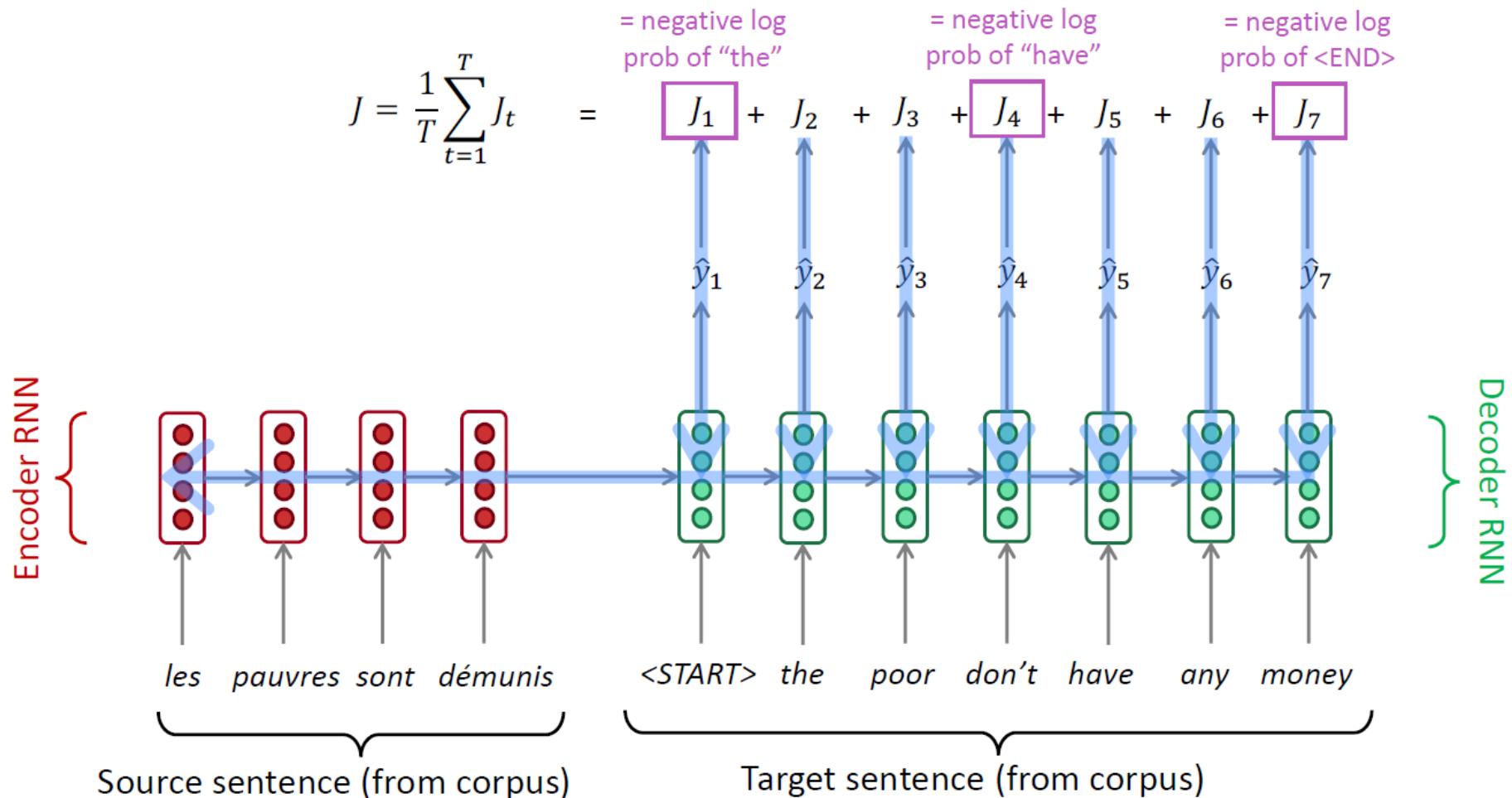
$e = (\text{Economic, growth, has, slowed, down, in, recent, years, .})$

Neural Machine Translation

- Model- *decoder*



Training a Neural Machine Translation system



Seq2seq is optimized as a **single system**.
Backpropagation operates “*end to end*”.

Advantages of NMT

Compared to SMT, NMT has many advantages:

- Better performance
 - More fluent
 - Better use of context
 - Better use of phrase similarities
- A single neural network to be optimized end-to-end
 - No subcomponents to be individually optimized
- Requires much less human engineering effort
 - No feature engineering
 - Same method for all language pairs

Disadvantages of NMT?

Compared to SMT:

- NMT is less interpretable
 - Hard to debug
- NMT is difficult to control
 - For example, can't easily specify rules or guidelines for translation
 - Safety concerns!

Evaluation (Machine Translation)

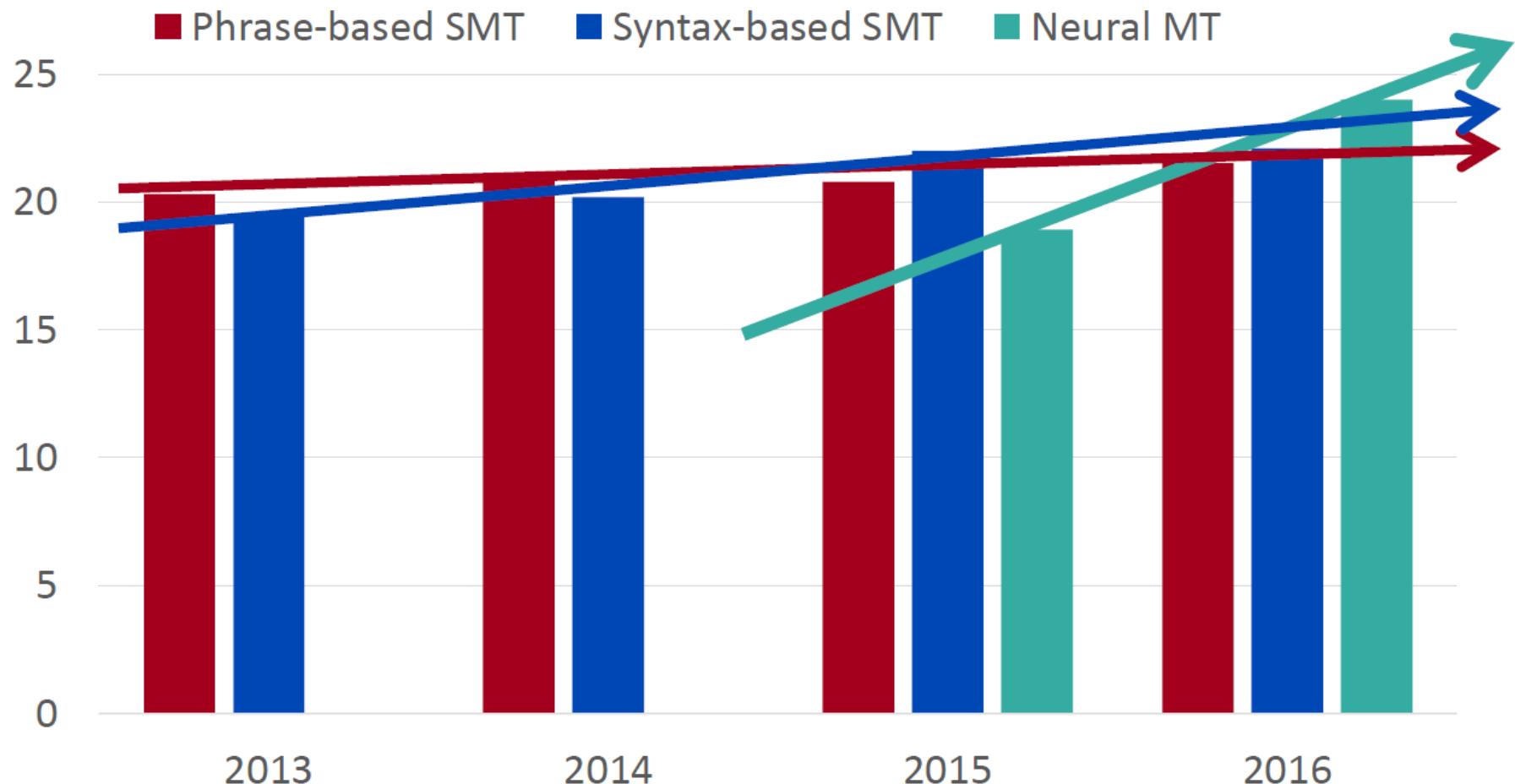
- BLEU (bilingual evaluation understudy) (Papineni et al. (2002))
 - BLEU compares the machine-written translation to **one or several** human-written translation(s), and computes a similarity score based on:
 - *n*-gram precision (usually up to 3 or 4-grams)
 - Penalty for too-short system translations
 - BLEUs output is always a number between 0 and 1
 - 1 means identical to the reference translations

BLEU is useful but imperfect

- BLEU was one of the first metrics to claim a high correlation with human judgements of quality, and remains one of the most popular automated and inexpensive metrics
- There are many valid ways to translate a sentence
- So a good translation can get a poor BLEU score because it has low n -gram overlap with the human translation L

MT progress over time

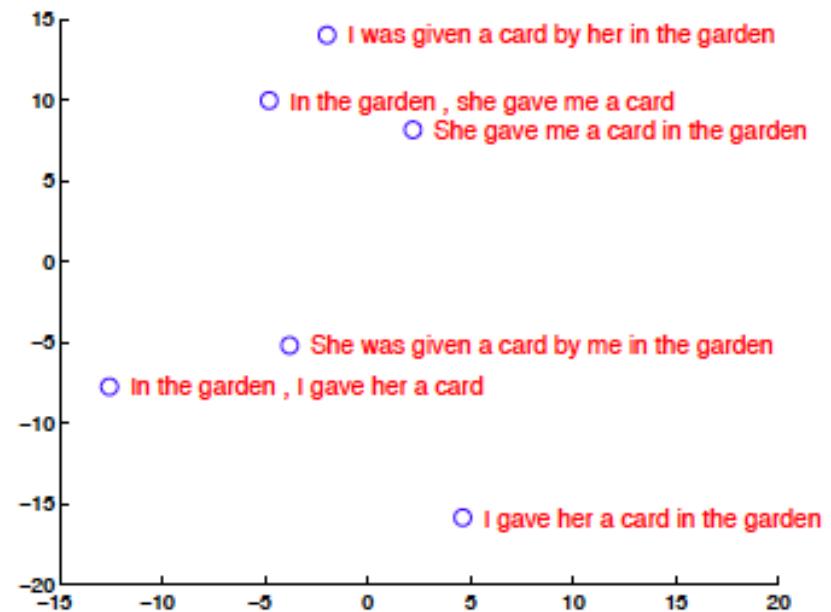
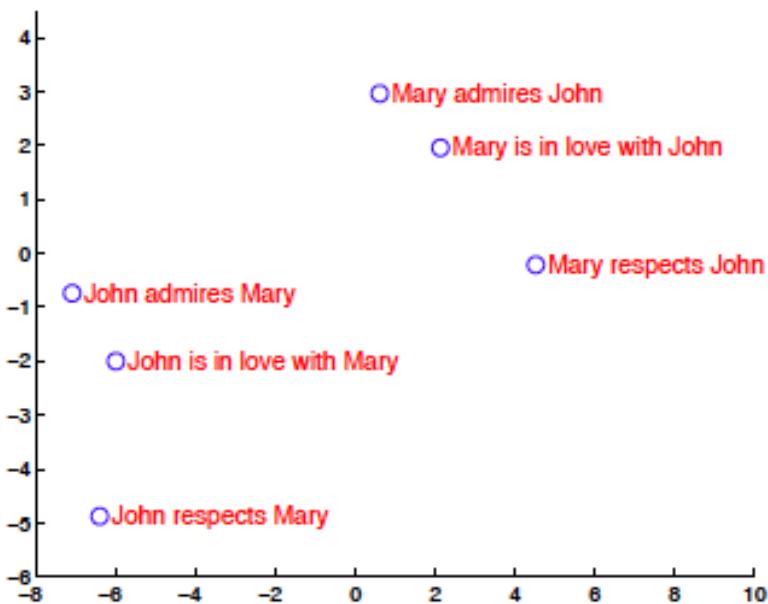
[Edinburgh En-De WMT newstest2013 Cased BLEU; NMT 2015 from U. Montréal]



Source: http://www.meta-net.eu/events/meta-forum-2016/slides/09_sennrich.pdf

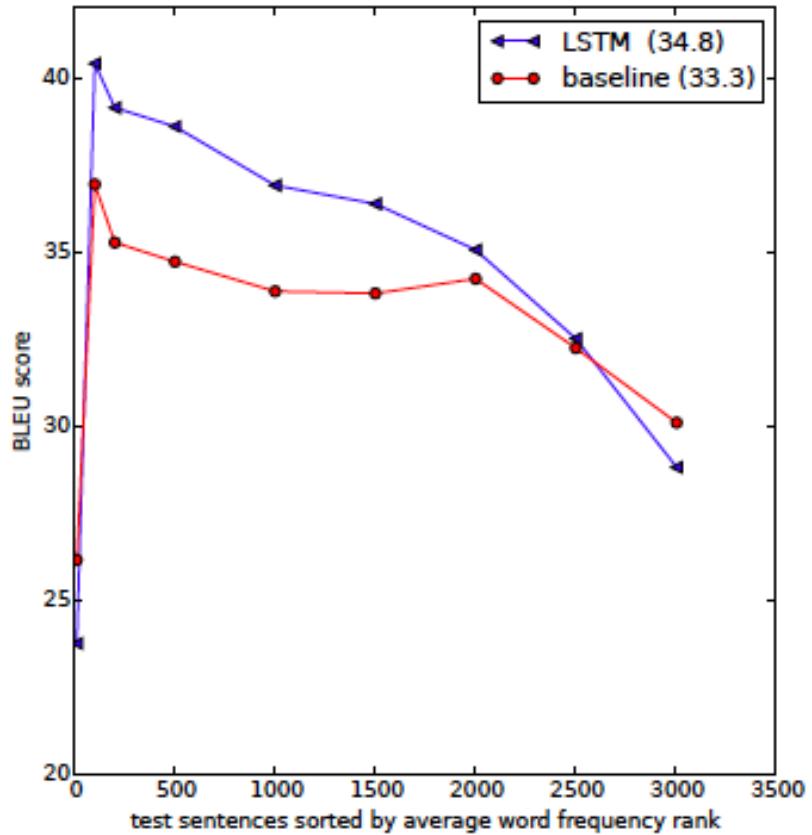
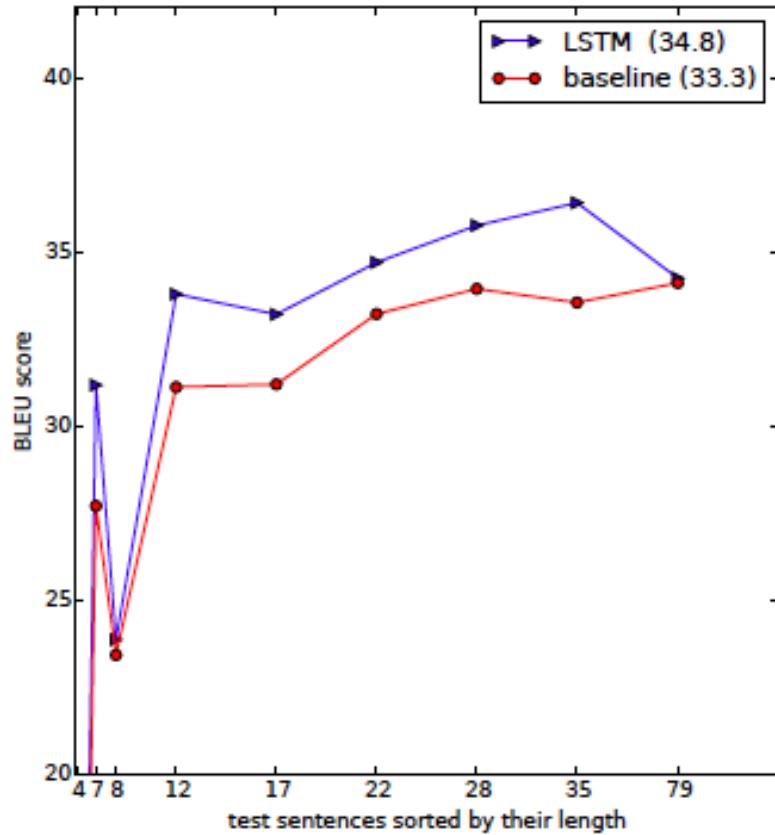
NMT Analysis

- Model Analysis



NMT Analysis

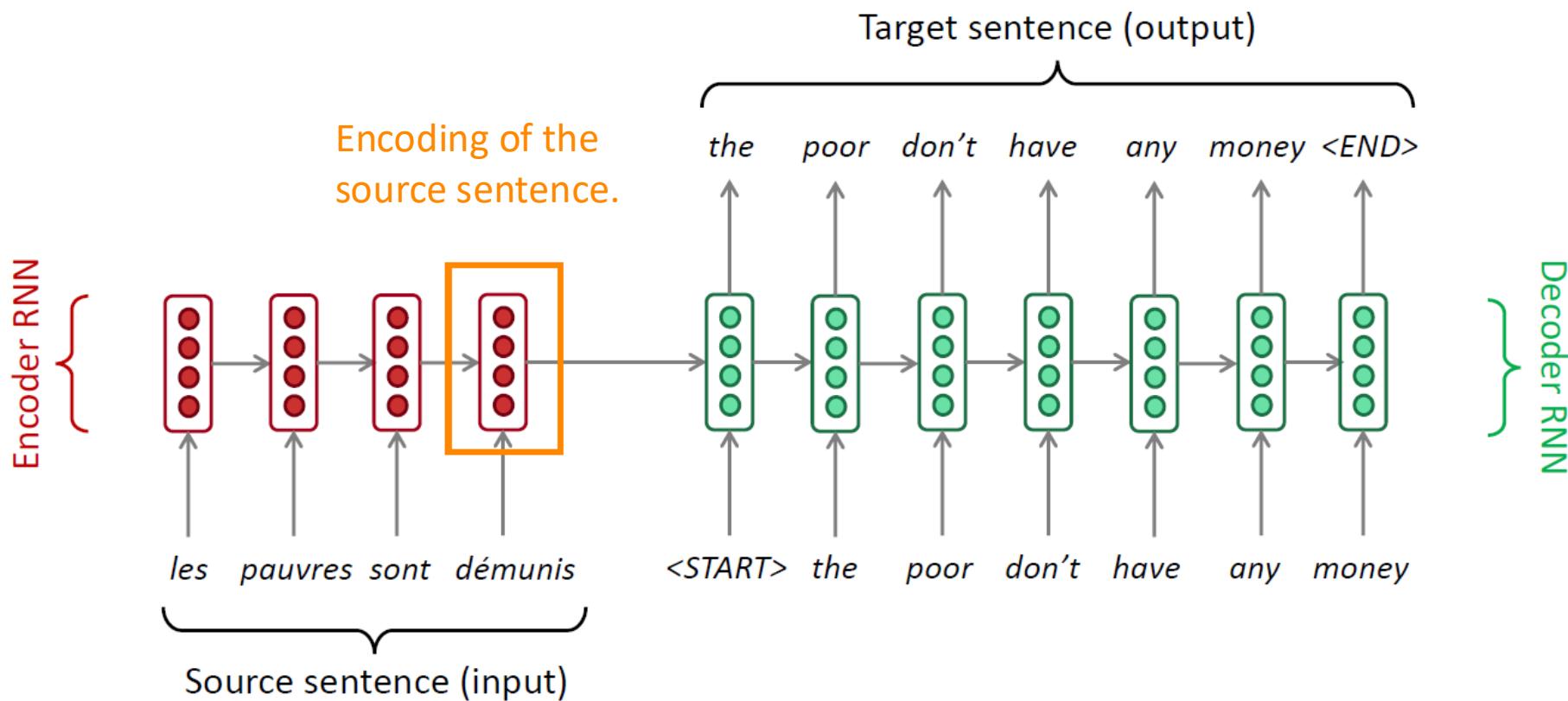
- Long sentences



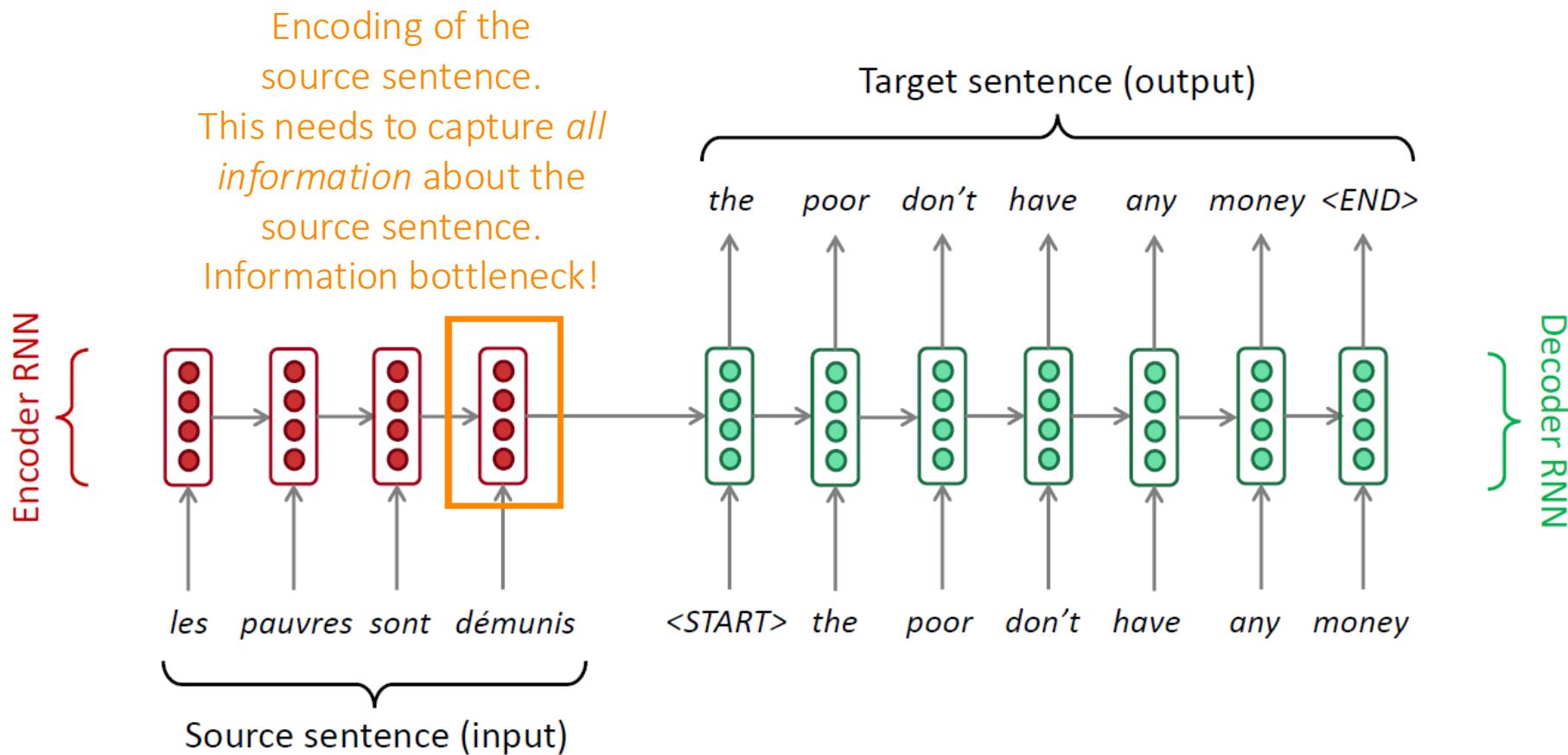
NMT: the biggest success story of NLP Deep Learning

- Neural Machine Translation went from a fringe research activity in **2014** to the leading standard method in **2016**
 - **2014**: First seq2seq paper published
 - **2016**: Google Translate switches from SMT to NMT
- This is amazing!
 - **SMT** systems, built by hundreds of engineers over many years, outperformed by NMT systems trained by a handful of engineers in a few months

Sequence-to-sequence: the bottleneck problem



Sequence-to-sequence: the bottleneck problem



Attention

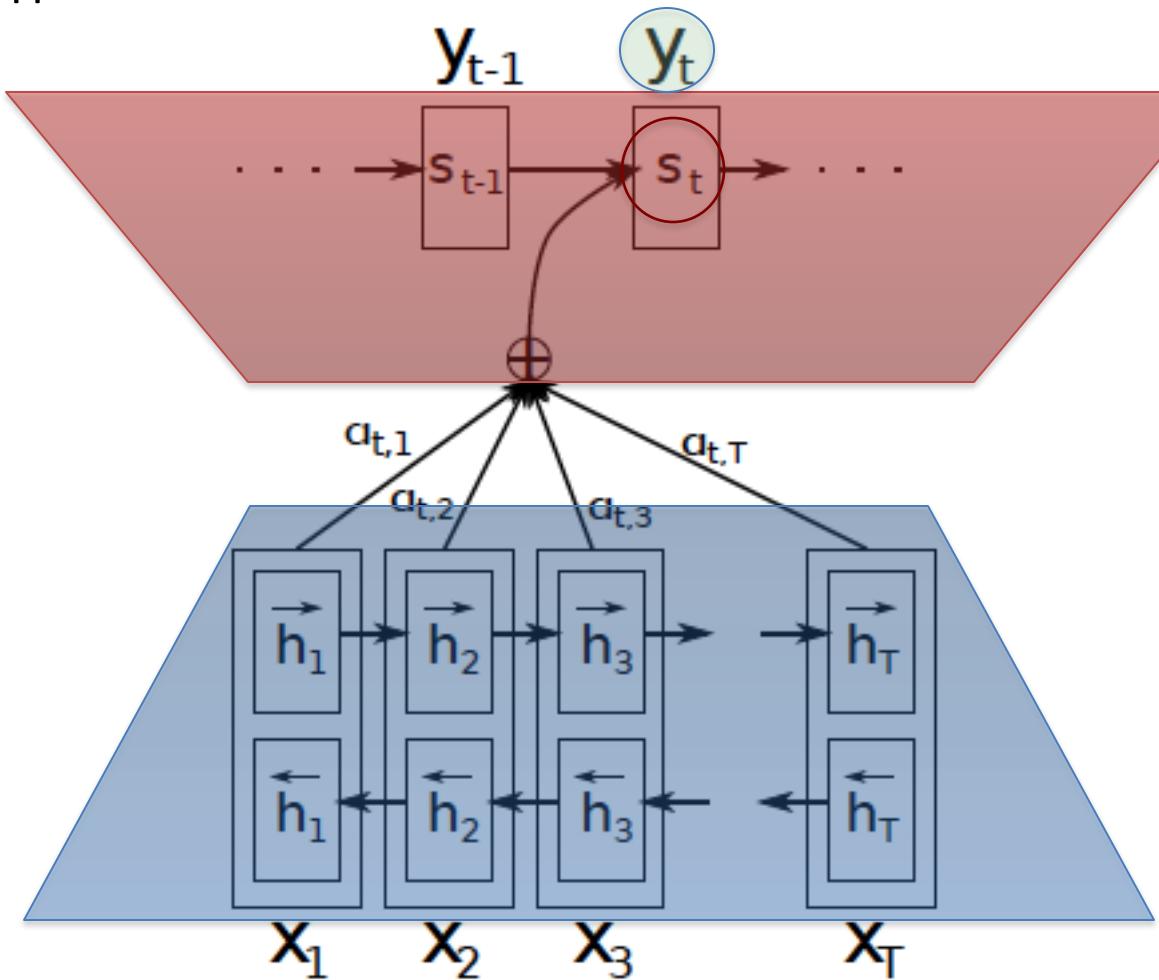
- **Attention** provides a solution to the bottleneck problem.
- Core idea: on each step of the decoder, *focus on a particular part* of the source sequence

Jointly Learning to Align and Translate

- Attention mechanism

$$p(y_t | y_{t-1}, \dots, y_1, x) = f(s_t)$$

$$s_t = g(s_{t-1}, y_{t-1}, c_t)$$



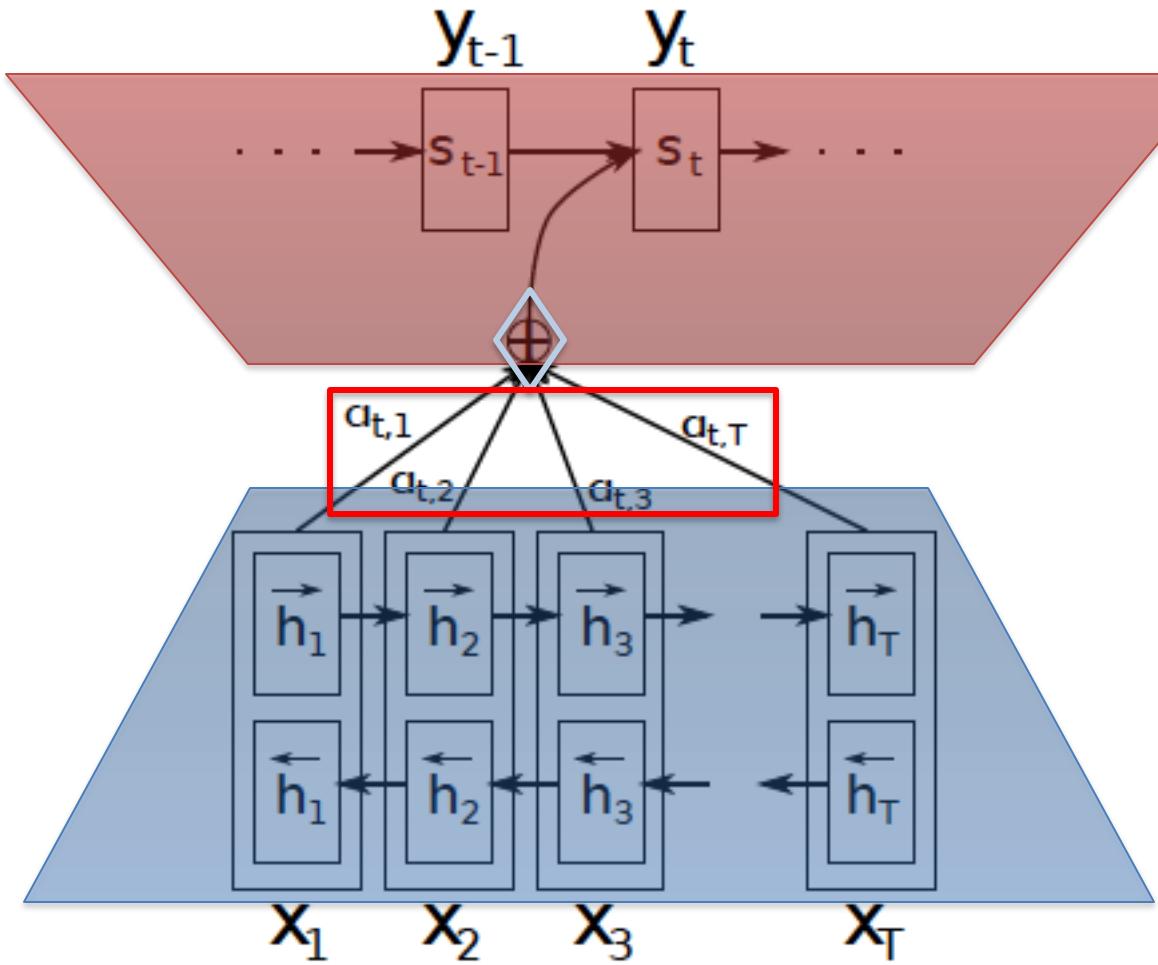
Jointly Learning to Align and Translate

- Attention mechanism

$$p(y_t | y_{t-1}, \dots, y_1, \mathbf{x}) = f(\mathbf{s}_t)$$

$$\mathbf{s}_t = g(\mathbf{s}_{t-1}, y_{t-1}, \mathbf{c}_t)$$

$$\mathbf{c}_t = \sum_{j=1}^T \alpha_{t,j} \mathbf{h}_j$$



Jointly Learning to Align and Translate

- Attention mechanism

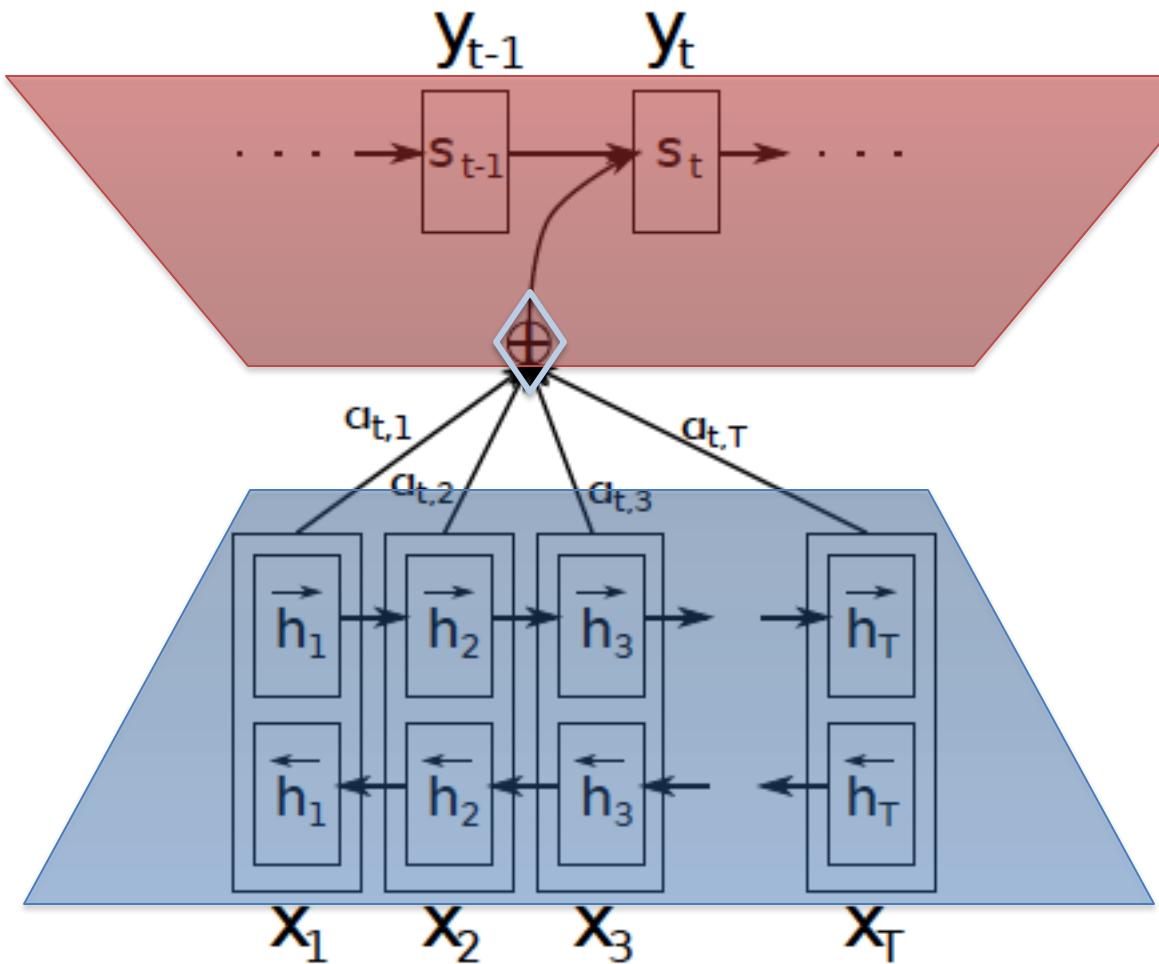
$$p(y_t | y_{t-1}, \dots, y_1, \mathbf{x}) = f(\mathbf{s}_t)$$

$$\mathbf{s}_t = g(\mathbf{s}_{t-1}, y_{t-1}, \mathbf{c}_t)$$

$$\mathbf{c}_t = \sum_{j=1}^T \alpha_{t,j} \mathbf{h}_j$$

$$\alpha_{t,j} = \frac{\exp(e_{t,j})}{\sum_{k=1}^T \exp(e_{t,k})}$$

$$e_{t,j} = a(\mathbf{s}_{t-1}, \mathbf{h}_j)$$



Jointly Learning to Align and Translate

- Attention mechanism

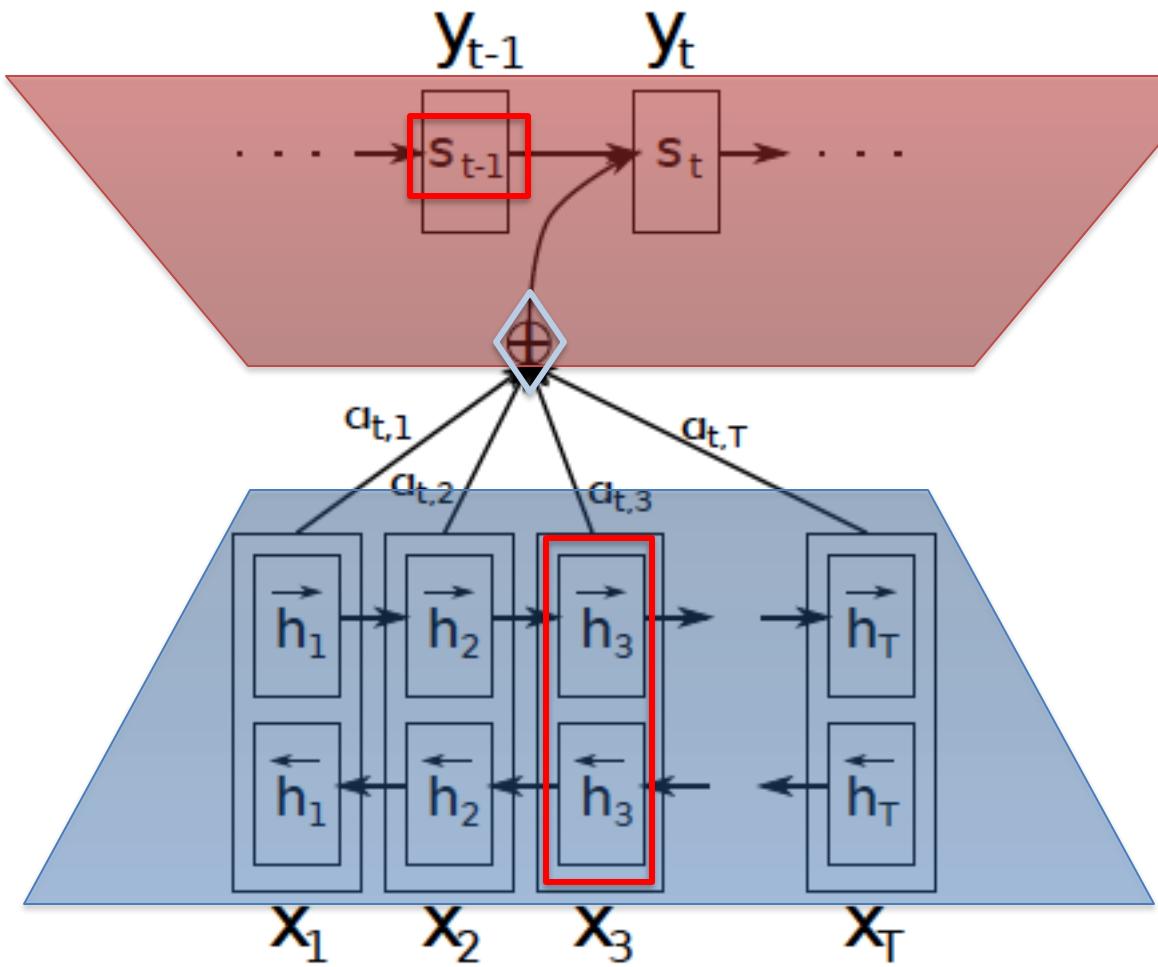
$$p(y_t | y_{t-1}, \dots, y_1, \mathbf{x}) = f(\mathbf{s}_t)$$

$$\mathbf{s}_t = g(\mathbf{s}_{t-1}, y_{t-1}, \mathbf{c}_t)$$

$$\mathbf{c}_t = \sum_{j=1}^T \alpha_{t,j} \mathbf{h}_j$$

$$\alpha_{t,j} = \frac{\exp(e_{t,j})}{\sum_{k=1}^T \exp(e_{t,k})}$$

$$e_{t,j} = a(\mathbf{s}_{t-1}, \mathbf{h}_j)$$



Attention variants

$$p(y_t | y_{t-1}, \dots, y_1, \mathbf{x}) = f(\mathbf{s}_t)$$

$$\mathbf{s}_t = g(\mathbf{s}_{t-1}, y_{t-1}, \mathbf{c}_t)$$

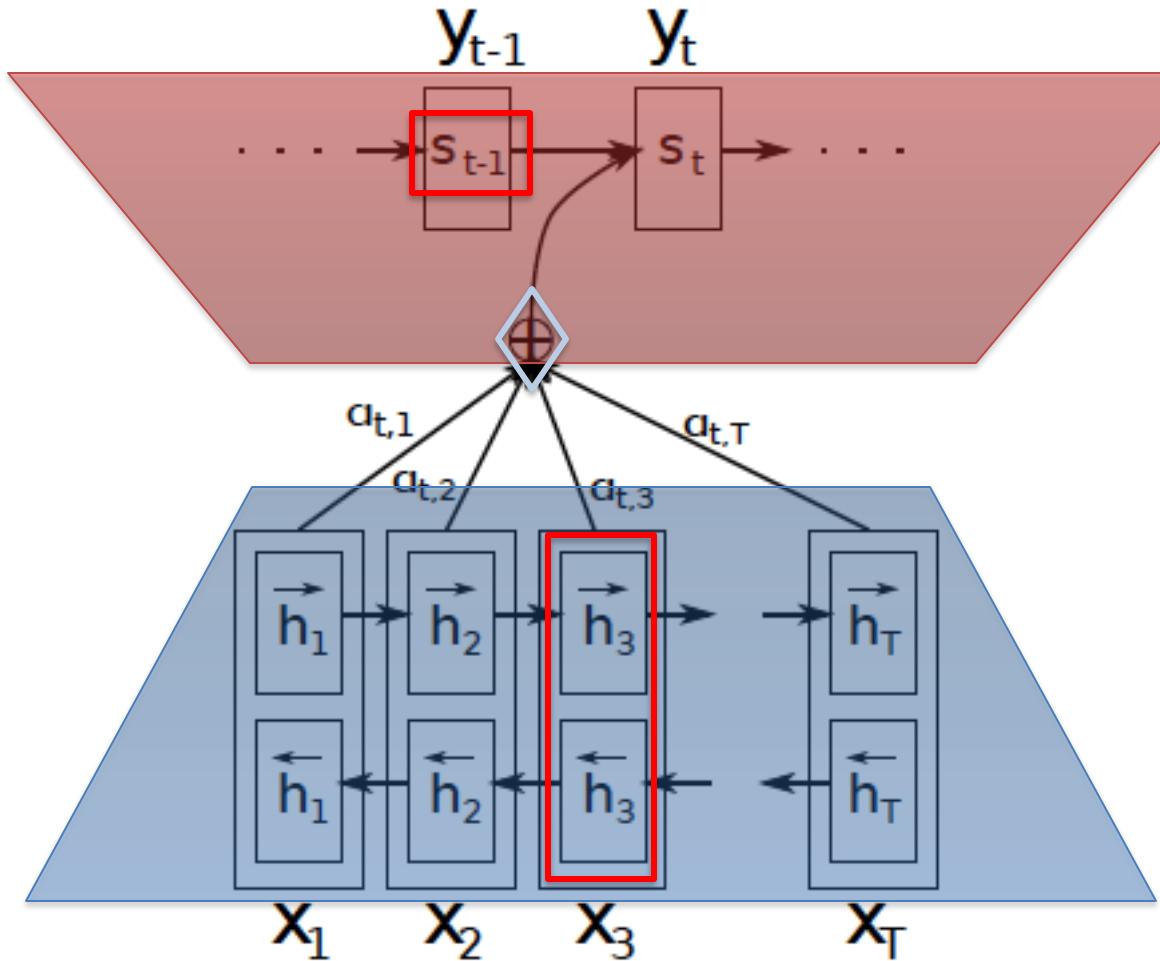
$$\mathbf{c}_t = \sum_{j=1}^T \alpha_{t,j} \mathbf{h}_j$$

$$\alpha_{t,j} = \frac{\exp(e_{t,j})}{\sum_{k=1}^T \exp(e_{t,k})}$$

$$e_{t,j} = a(\mathbf{s}_{t-1}, \mathbf{h}_j)$$

Basic dot-product attention

$$e_{tj} = \mathbf{s}_{t-1}^T \mathbf{h}_j$$



Attention variants

$$p(y_t | y_{t-1}, \dots, y_1, \mathbf{x}) = f(\mathbf{s}_t)$$

$$\mathbf{s}_t = g(\mathbf{s}_{t-1}, y_{t-1}, \mathbf{c}_t)$$

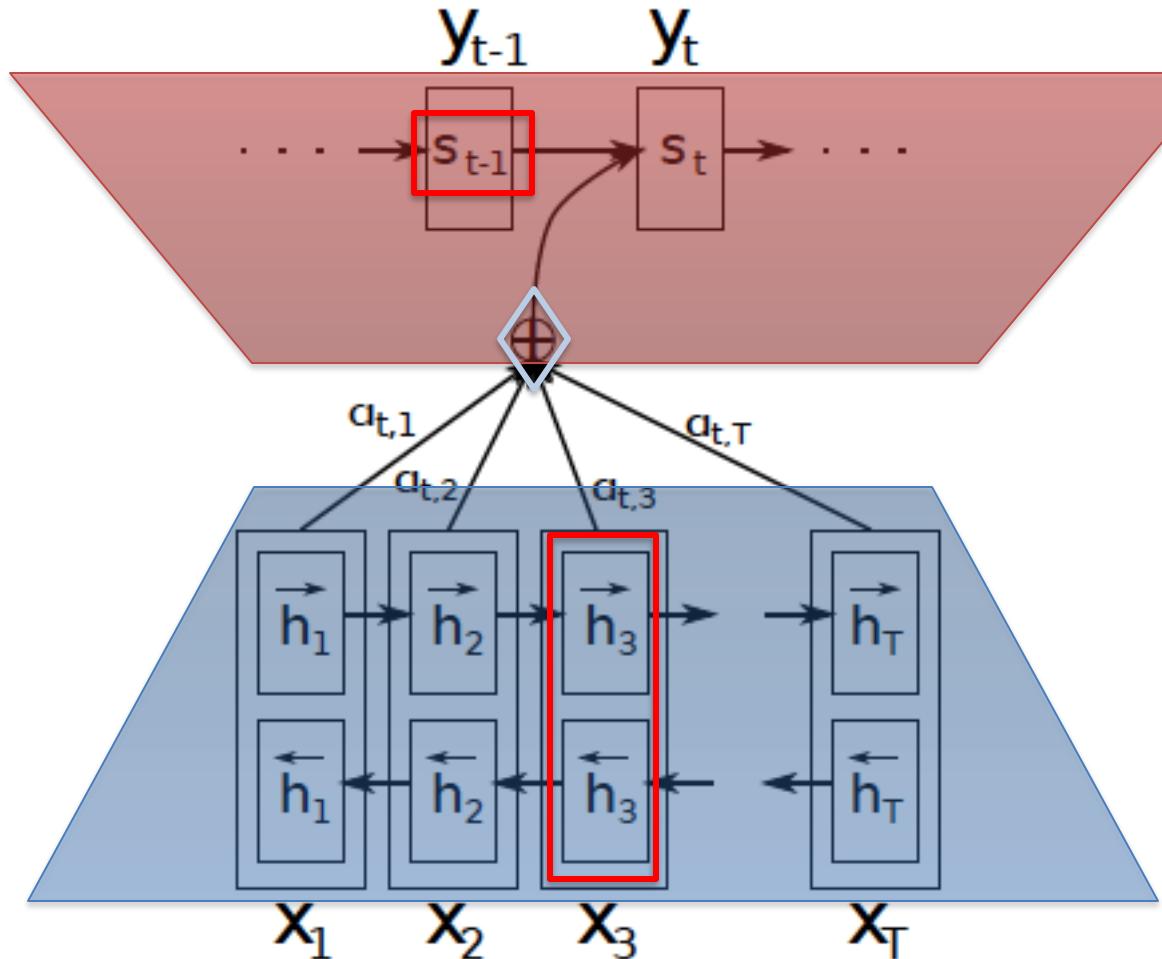
$$\mathbf{c}_t = \sum_{j=1}^T \alpha_{t,j} \mathbf{h}_j$$

$$\alpha_{t,j} = \frac{\exp(e_{t,j})}{\sum_{k=1}^T \exp(e_{t,k})}$$

$$e_{t,j} = a(\mathbf{s}_{t-1}, \mathbf{h}_j)$$

Multiplicative attention .
(bilinear mode)

$$e_{tj} = \mathbf{s}_{t-1}^T W \mathbf{h}_j$$



Attention variants

$$p(y_t | y_{t-1}, \dots, y_1, \mathbf{x}) = f(\mathbf{s}_t)$$

$$\mathbf{s}_t = g(\mathbf{s}_{t-1}, y_{t-1}, \mathbf{c}_t)$$

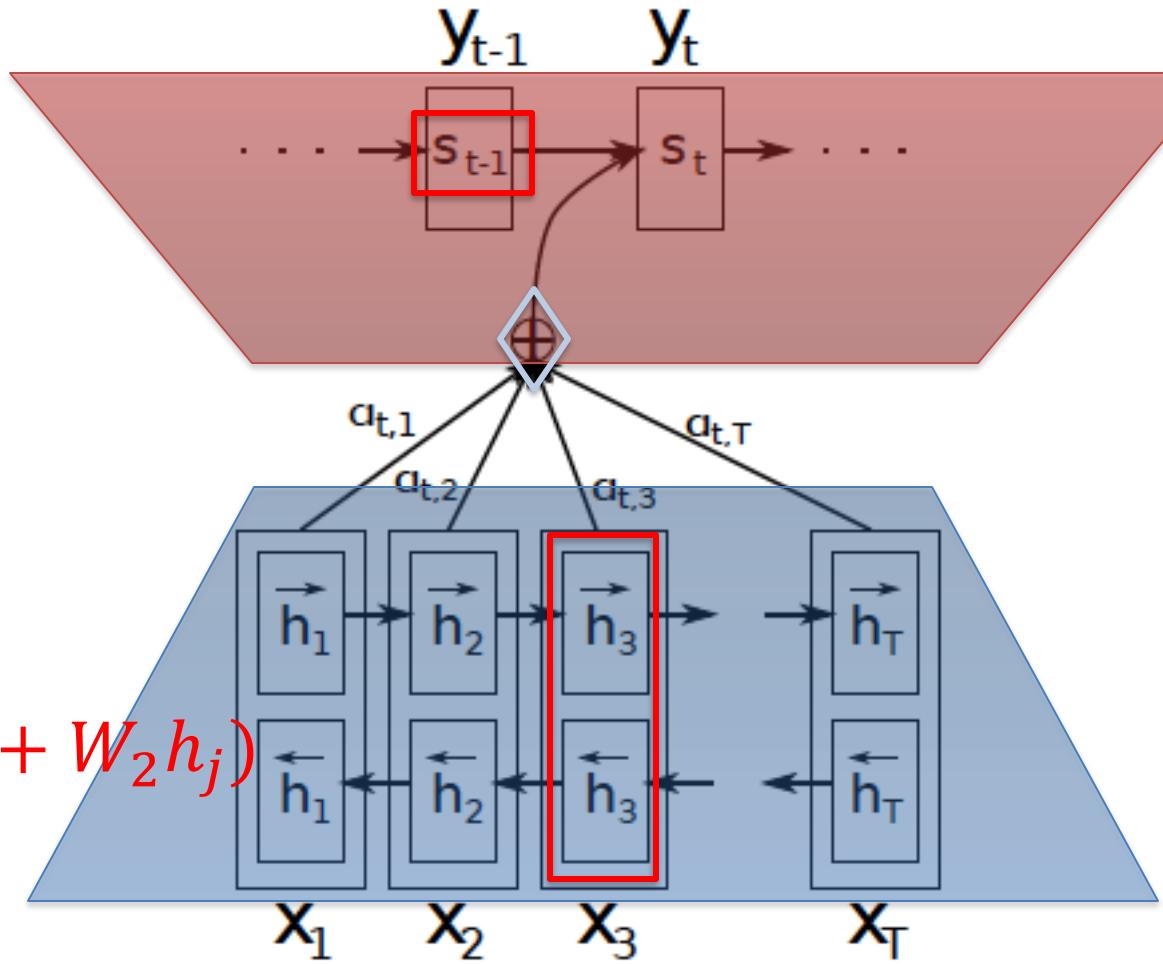
$$\mathbf{c}_t = \sum_{j=1}^T \alpha_{t,j} \mathbf{h}_j$$

$$\alpha_{t,j} = \frac{\exp(e_{t,j})}{\sum_{k=1}^T \exp(e_{t,k})}$$

$$e_{t,j} = a(\mathbf{s}_{t-1}, \mathbf{h}_j)$$

Additive attention

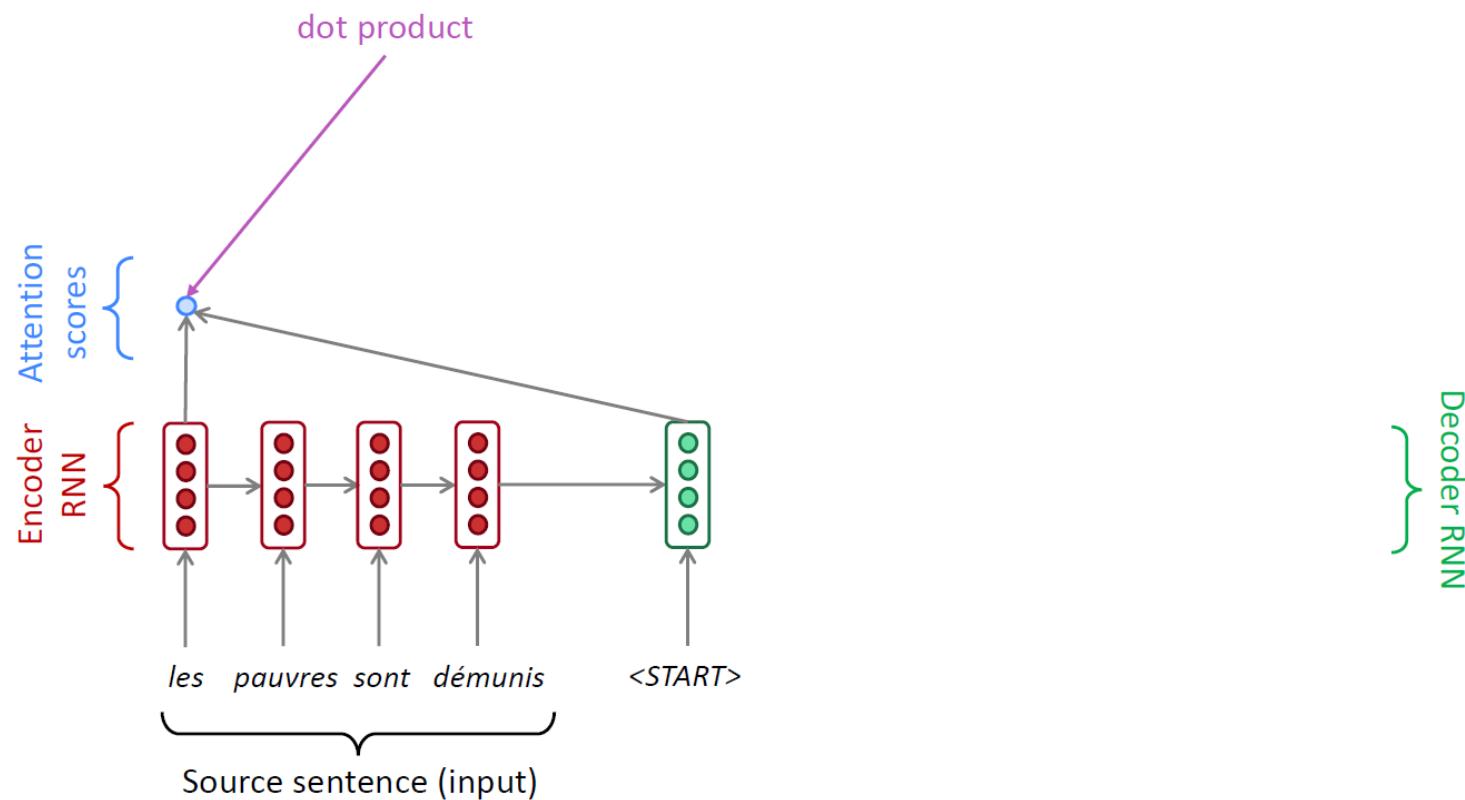
$$e_{t,j} = \mathbf{v}^T \tanh(W_1 \mathbf{s}_{t-1}^T + W_2 \mathbf{h}_j)$$



A Running Example (Slightly Different From Bengio's Original Method)

$$\alpha_{t,j} = \frac{\exp(e_{t,j})}{\sum_{k=1}^T \exp(e_{t,k})}$$

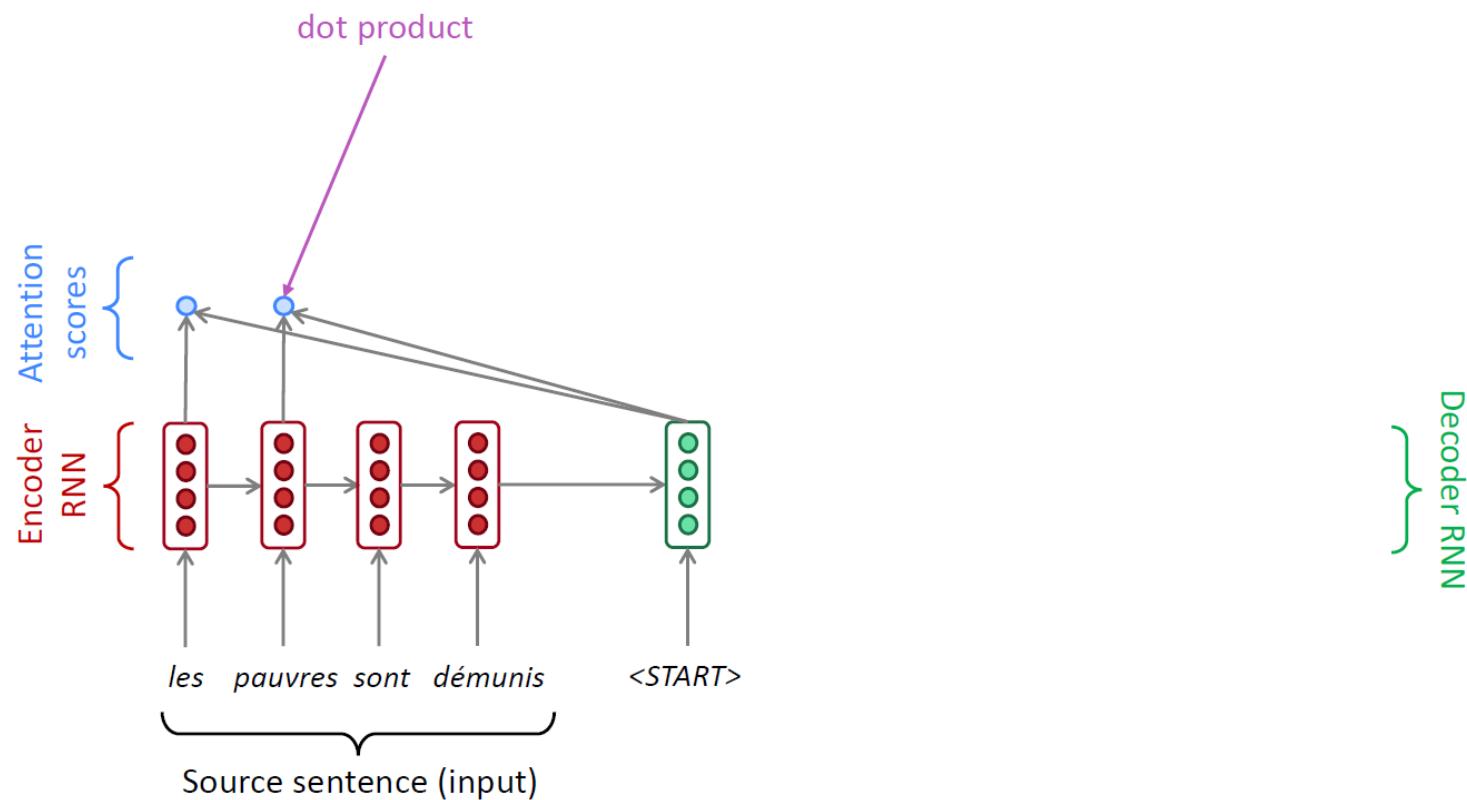
$$e_{t,j} = a(\mathbf{s}_t, \mathbf{h}_j)$$



A Running Example

$$\alpha_{t,j} = \frac{\exp(e_{t,j})}{\sum_{k=1}^T \exp(e_{t,k})}$$

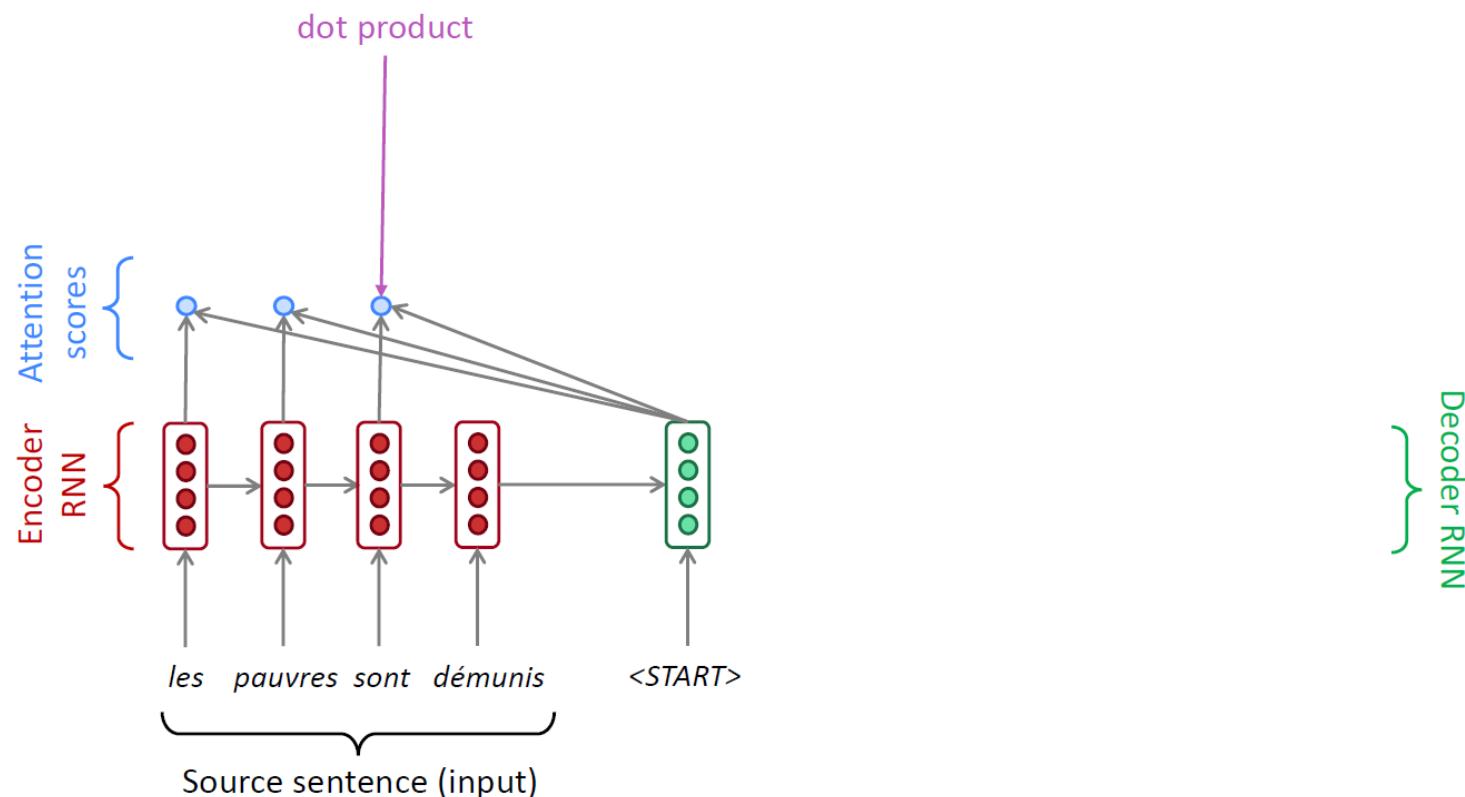
$$e_{t,j} = a(\mathbf{s}_t, \mathbf{h}_j)$$



A Running Example

$$\alpha_{t,j} = \frac{\exp(e_{t,j})}{\sum_{k=1}^T \exp(e_{t,k})}$$

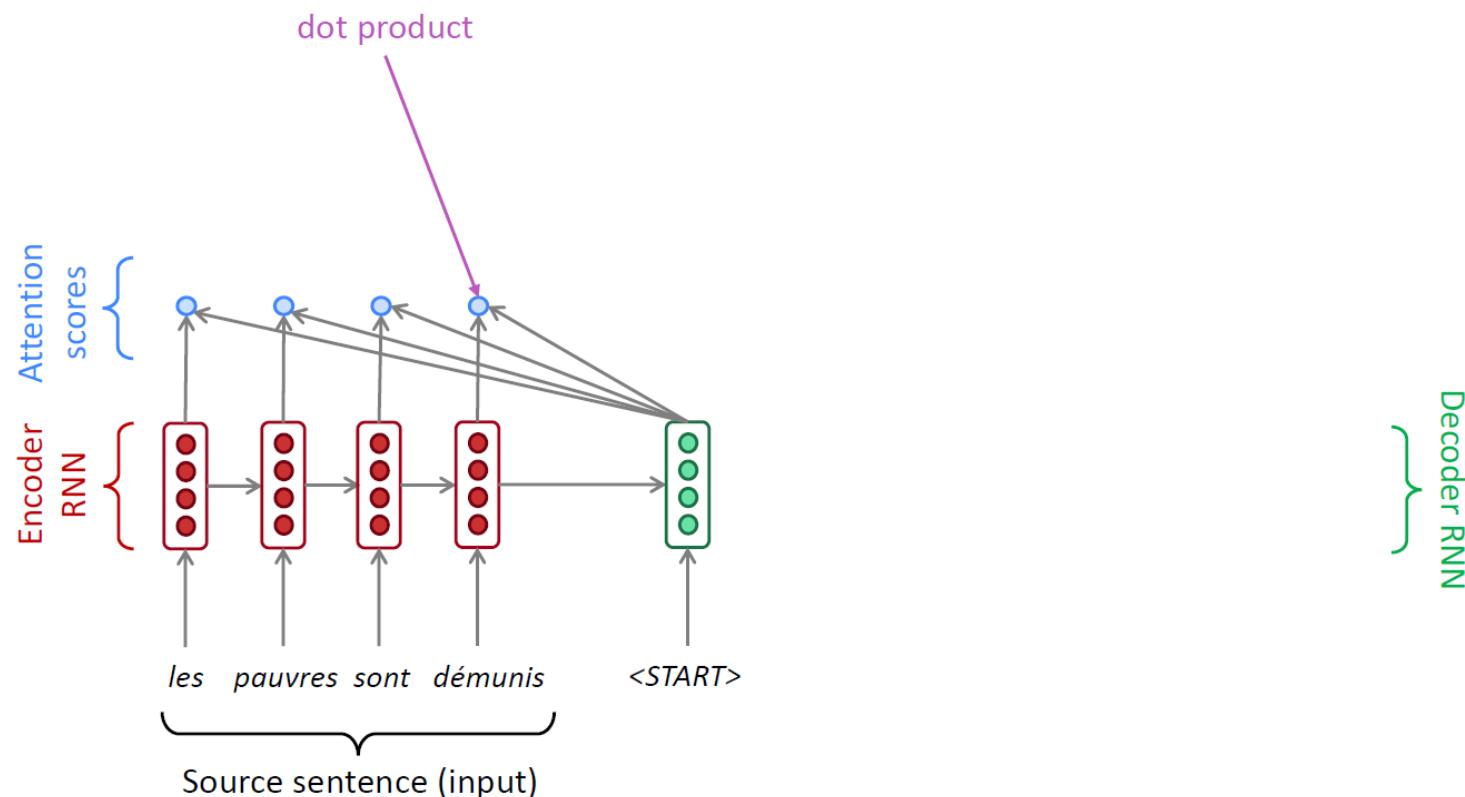
$$e_{t,j} = a(\mathbf{s}_t, \mathbf{h}_j)$$



A Running Example

$$\alpha_{t,j} = \frac{\exp(e_{t,j})}{\sum_{k=1}^T \exp(e_{t,k})}$$

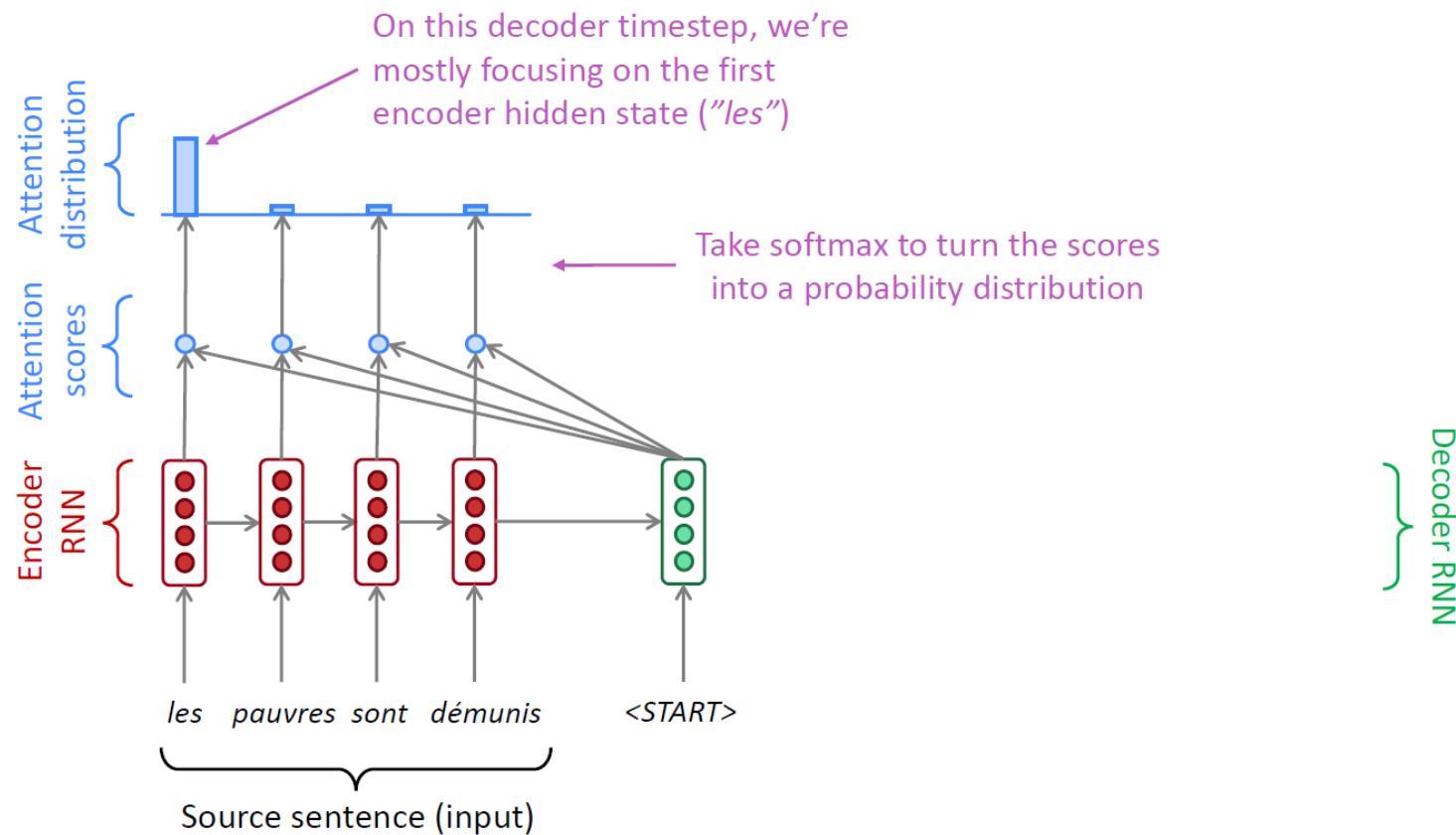
$$e_{t,j} = a(\mathbf{s}_t, \mathbf{h}_j)$$



A Running Example

$$\alpha_{t,j} = \frac{\exp(e_{t,j})}{\sum_{k=1}^T \exp(e_{t,k})}$$

$$e_{t,j} = a(\mathbf{s}_t, \mathbf{h}_j)$$

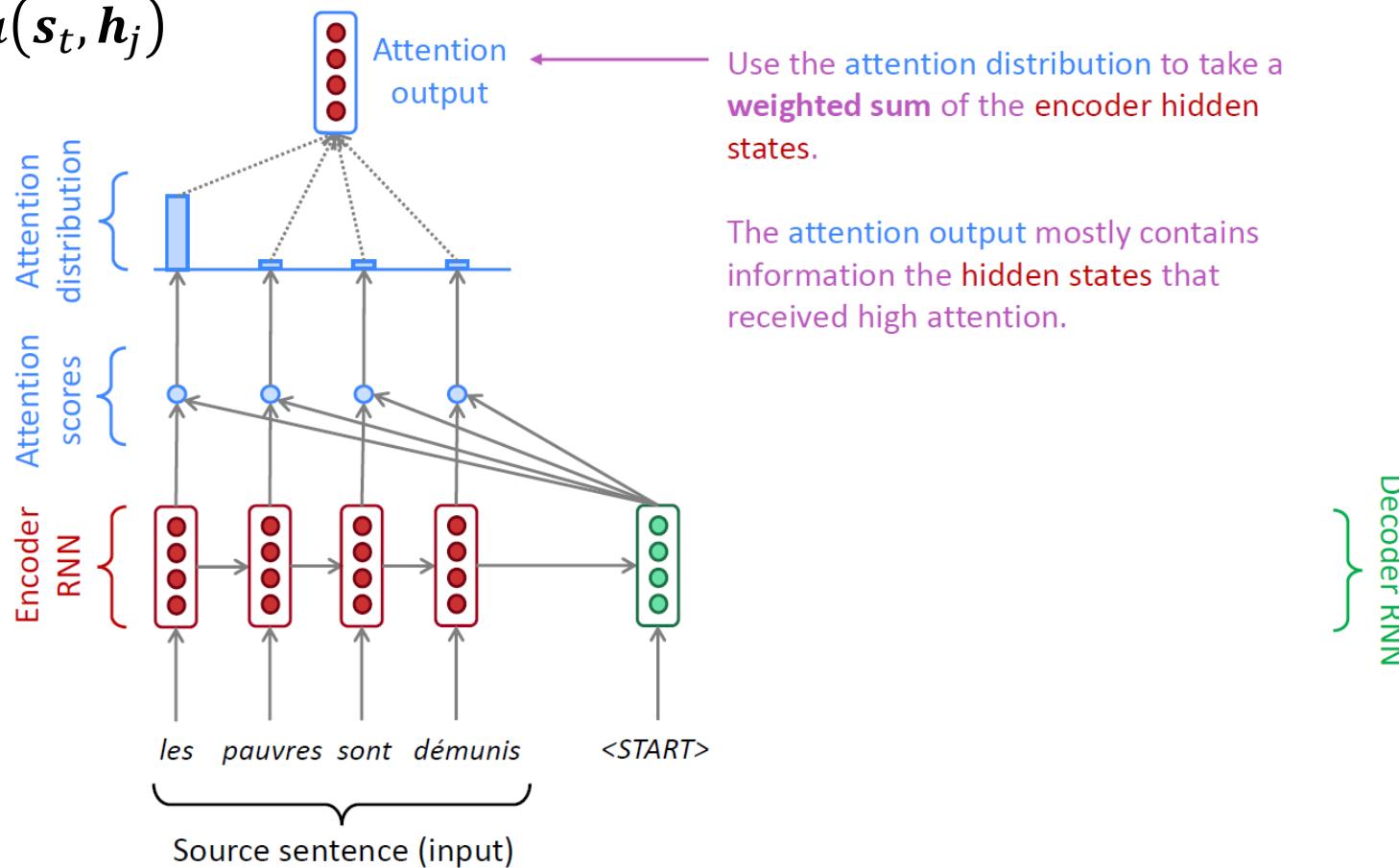


A Running Example

$$\alpha_{t,j} = \frac{\exp(e_{t,j})}{\sum_{k=1}^T \exp(e_{t,k})}$$

$$\mathbf{c}_t = \sum_{j=1}^T \alpha_{t,j} \mathbf{h}_j$$

$$e_{t,j} = a(\mathbf{s}_t, \mathbf{h}_j)$$



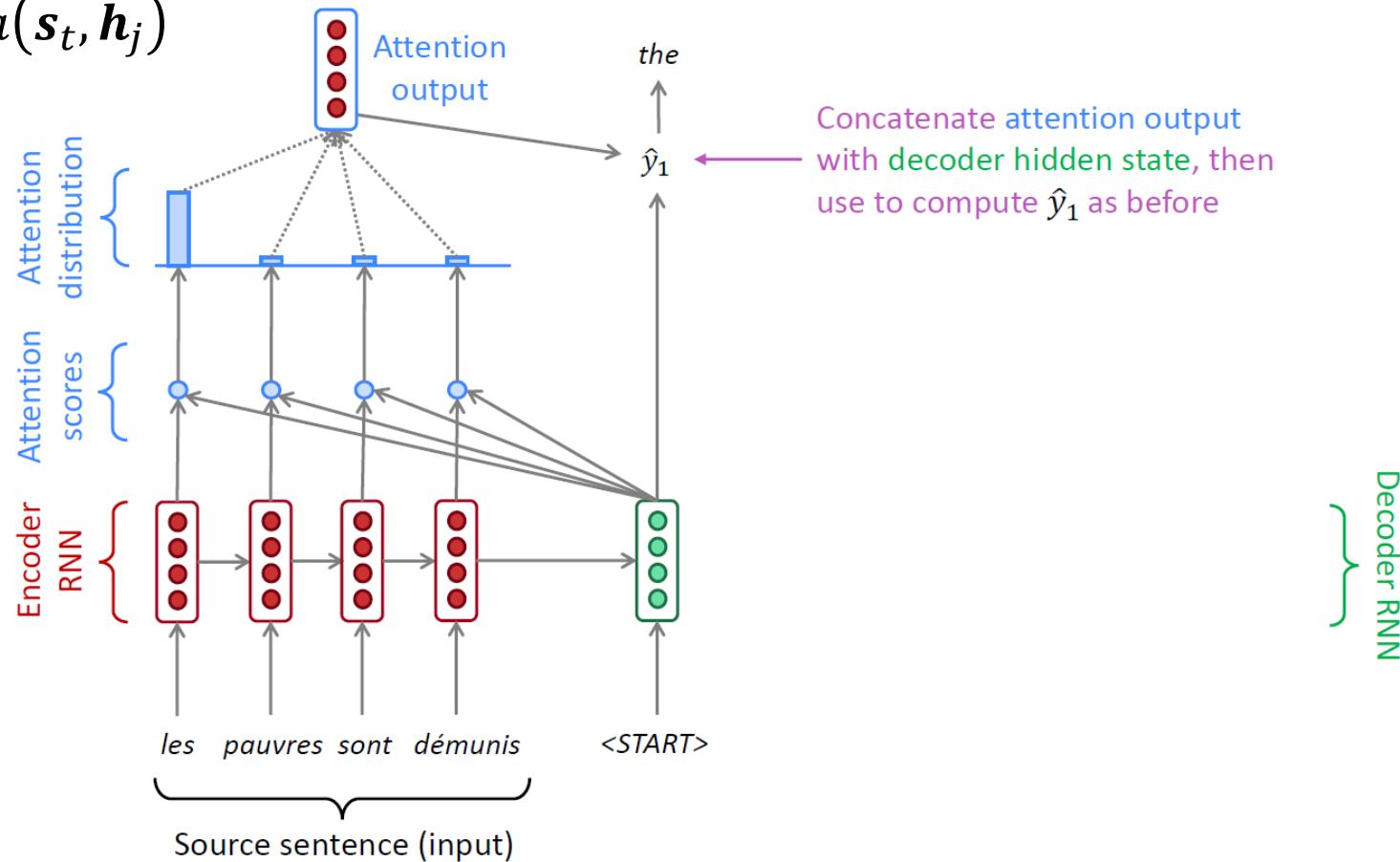
A Running Example

$$\alpha_{t,j} = \frac{\exp(e_{t,j})}{\sum_{k=1}^T \exp(e_{t,k})}$$

$$\mathbf{c}_t = \sum_{j=1}^T \alpha_{t,j} \mathbf{h}_j$$

$$p(y_t | y_{t-1}, \dots, y_1, \mathbf{x}) = f(y_{t-1}, \mathbf{s}_t, \mathbf{c}_t)$$

$$e_{t,j} = a(\mathbf{s}_t, \mathbf{h}_j)$$



A Running Example

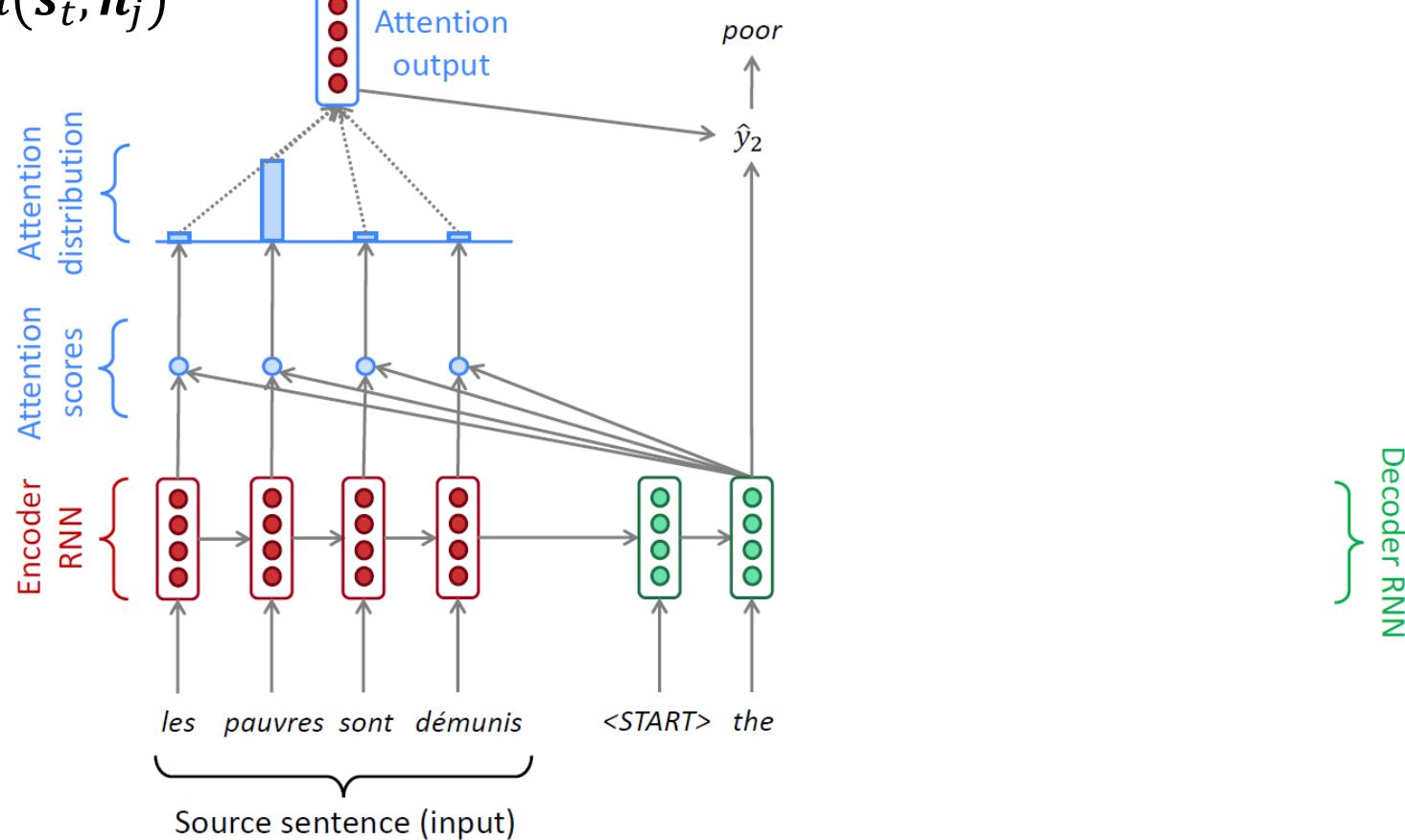
$$\alpha_{t,j} = \frac{\exp(e_{t,j})}{\sum_{k=1}^T \exp(e_{t,k})}$$

$$\mathbf{c}_t = \sum_{j=1}^T \alpha_{t,j} \mathbf{h}_j$$

$$p(y_t | y_{t-1}, \dots, y_1, \mathbf{x}) = f(y_{t-1}, \mathbf{s}_t, \mathbf{c}_t)$$

$$\mathbf{s}_{t+1} = g(\mathbf{s}_t, y_t)$$

$$e_{t,j} = a(\mathbf{s}_t, \mathbf{h}_j)$$



A Running Example

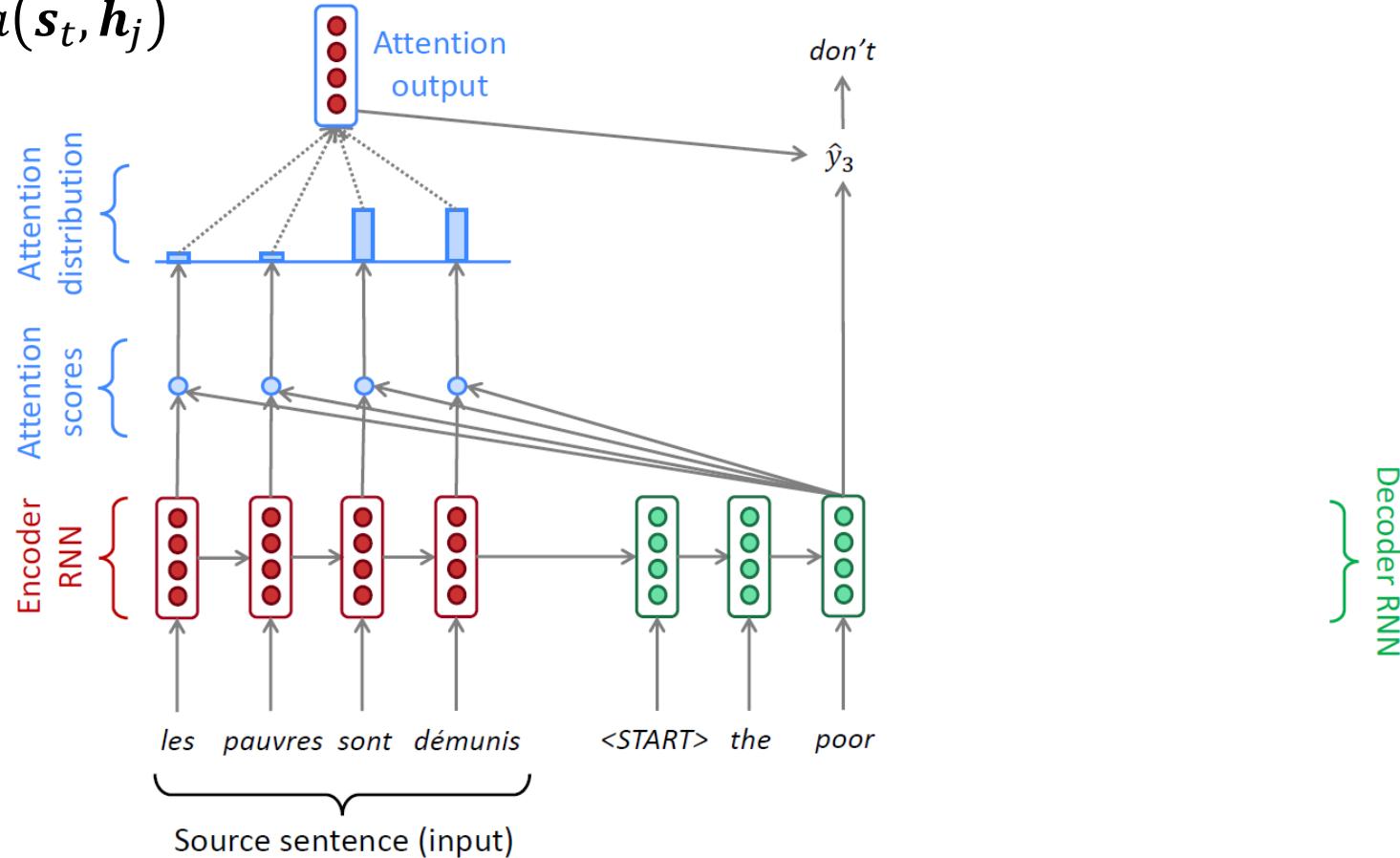
$$\alpha_{t,j} = \frac{\exp(e_{t,j})}{\sum_{k=1}^T \exp(e_{t,k})}$$

$$\mathbf{c}_t = \sum_{j=1}^T \alpha_{t,j} \mathbf{h}_j$$

$$p(y_t | y_{t-1}, \dots, y_1, \mathbf{x}) = f(y_{t-1}, \mathbf{s}_t, \mathbf{c}_t)$$

$$\mathbf{s}_{t+1} = g(\mathbf{s}_t, y_t)$$

$$e_{t,j} = a(\mathbf{s}_t, \mathbf{h}_j)$$



A Running Example

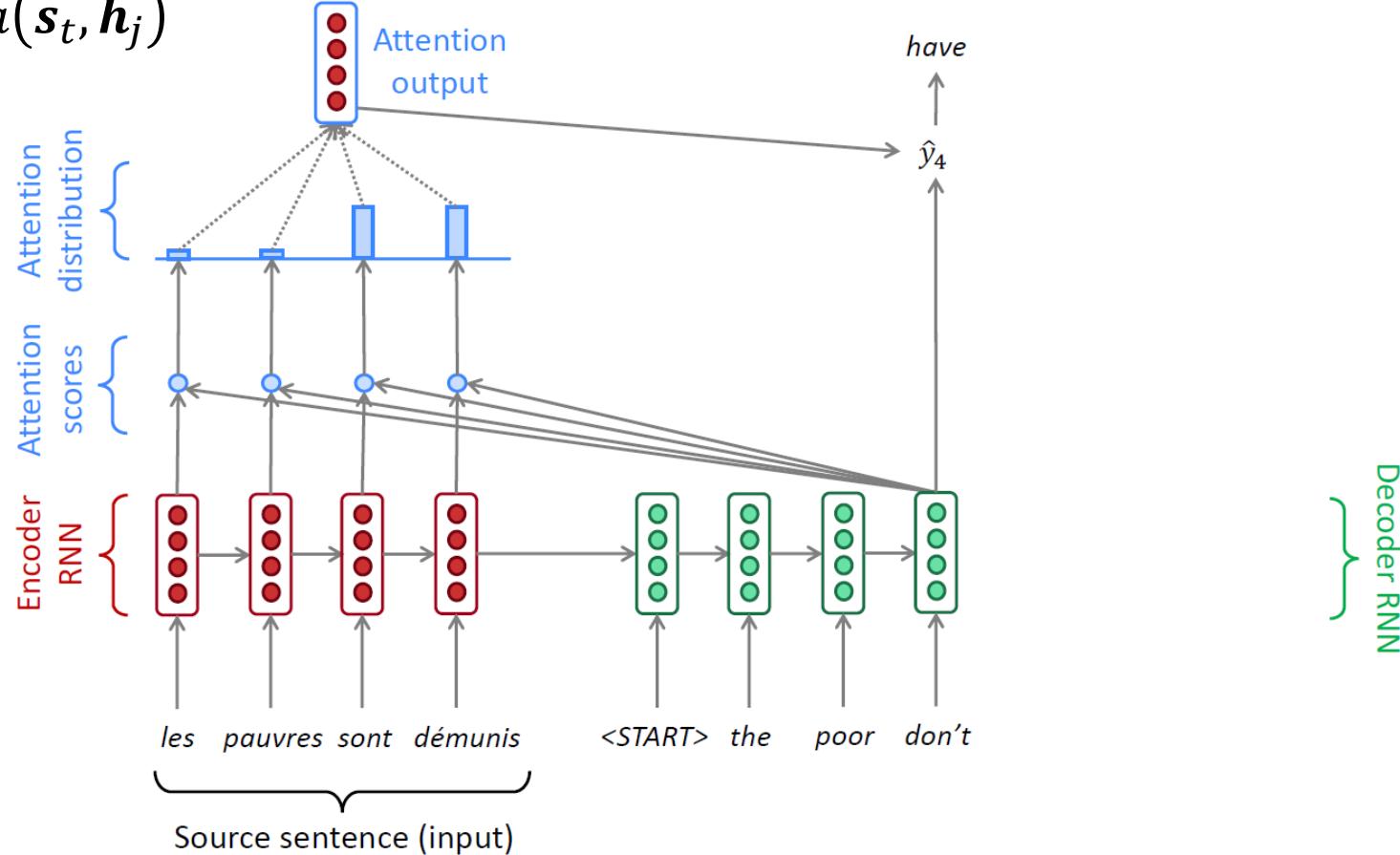
$$\alpha_{t,j} = \frac{\exp(e_{t,j})}{\sum_{k=1}^T \exp(e_{t,k})}$$

$$\mathbf{c}_t = \sum_{j=1}^T \alpha_{t,j} \mathbf{h}_j$$

$$p(y_t | y_{t-1}, \dots, y_1, \mathbf{x}) = f(y_{t-1}, \mathbf{s}_t, \mathbf{c}_t)$$

$$\mathbf{s}_{t+1} = g(\mathbf{s}_t, y_t)$$

$$e_{t,j} = a(\mathbf{s}_t, \mathbf{h}_j)$$



A Running Example

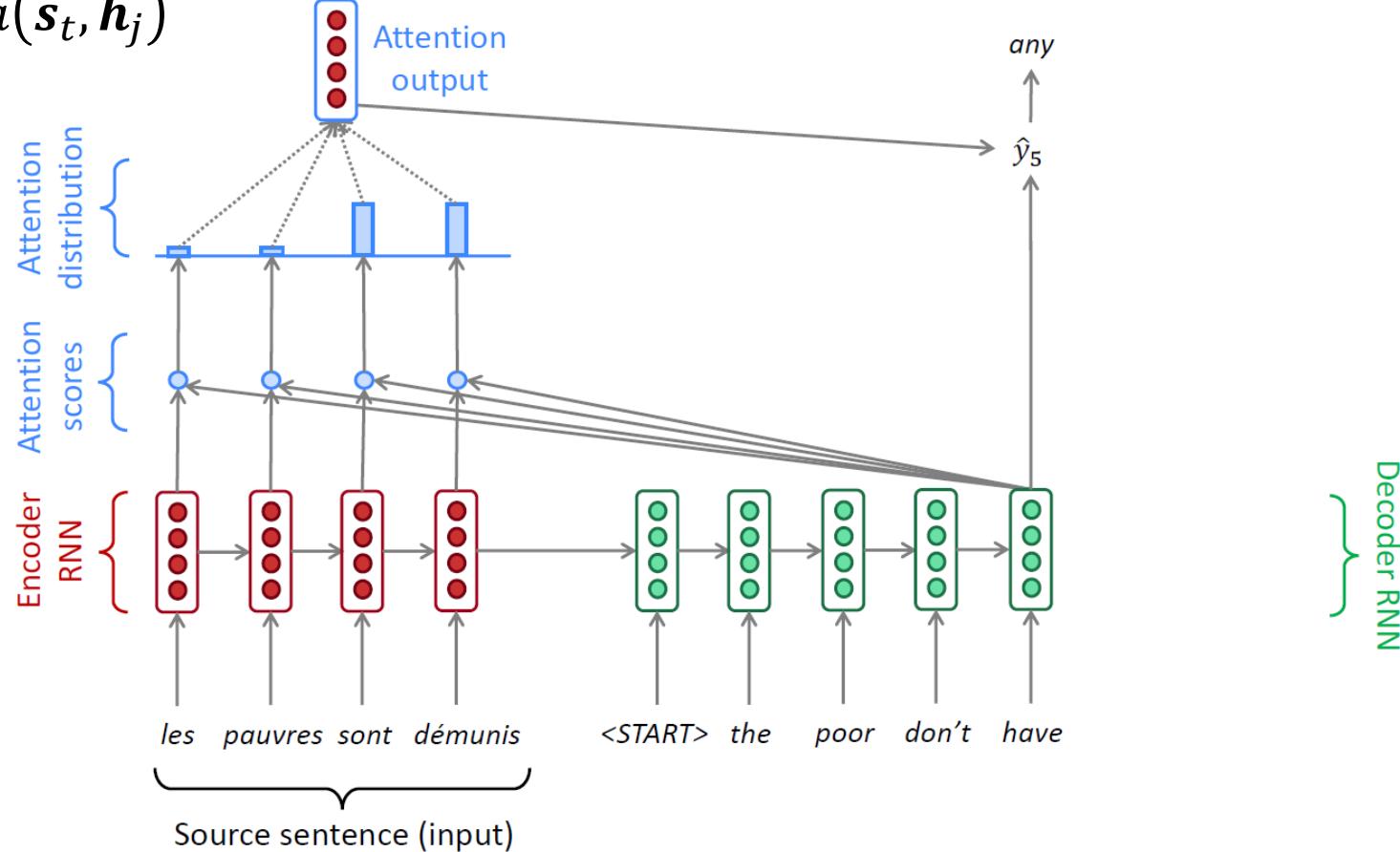
$$\alpha_{t,j} = \frac{\exp(e_{t,j})}{\sum_{k=1}^T \exp(e_{t,k})}$$

$$\mathbf{c}_t = \sum_{j=1}^T \alpha_{t,j} \mathbf{h}_j$$

$$p(y_t | y_{t-1}, \dots, y_1, \mathbf{x}) = f(y_{t-1}, \mathbf{s}_t, \mathbf{c}_t)$$

$$\mathbf{s}_{t+1} = g(\mathbf{s}_t, y_t)$$

$$e_{t,j} = a(\mathbf{s}_t, \mathbf{h}_j)$$



A Running Example

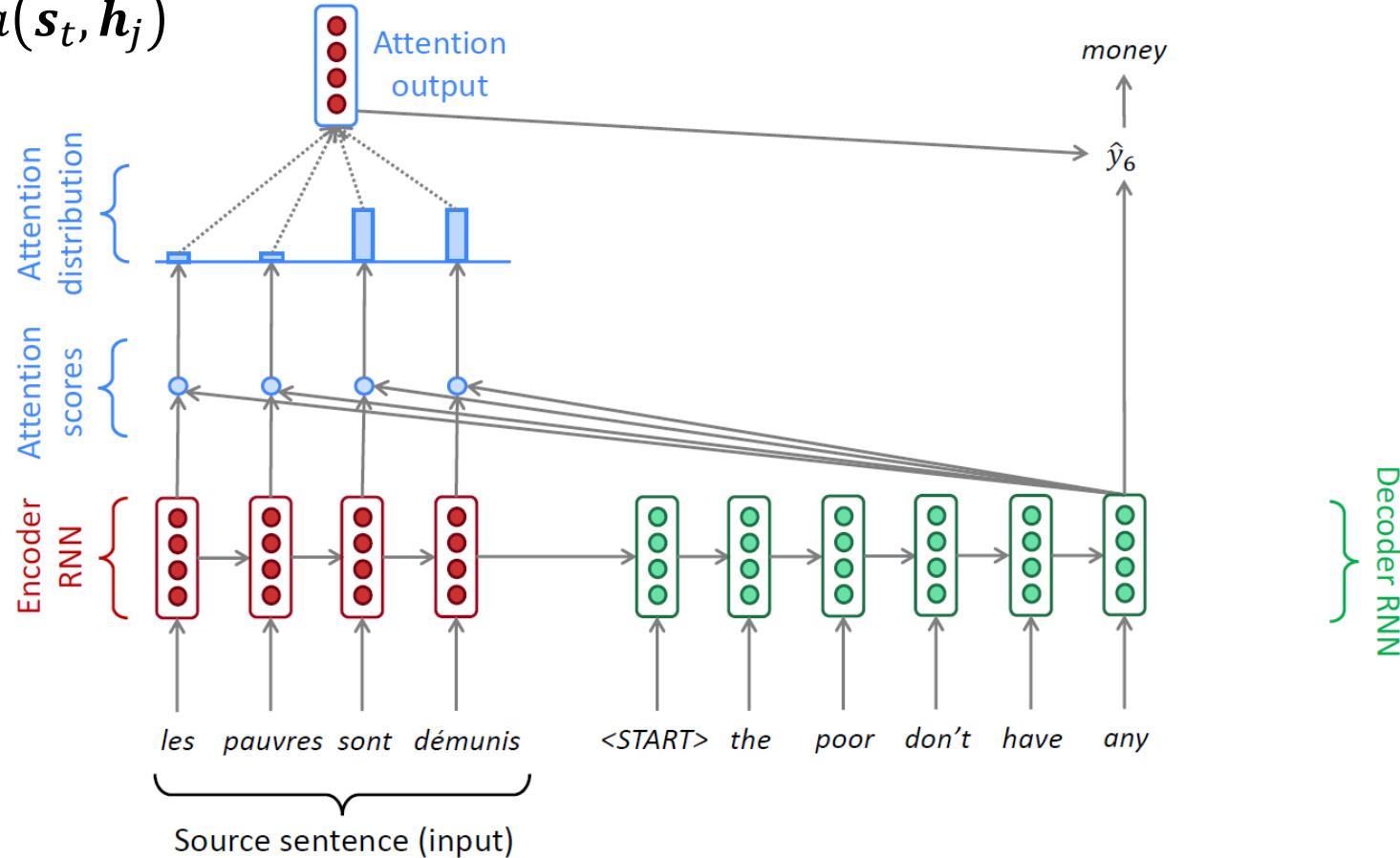
$$\alpha_{t,j} = \frac{\exp(e_{t,j})}{\sum_{k=1}^T \exp(e_{t,k})}$$

$$\mathbf{c}_t = \sum_{j=1}^T \alpha_{t,j} \mathbf{h}_j$$

$$p(y_t | y_{t-1}, \dots, y_1, \mathbf{x}) = f(y_{t-1}, \mathbf{s}_t, \mathbf{c}_t)$$

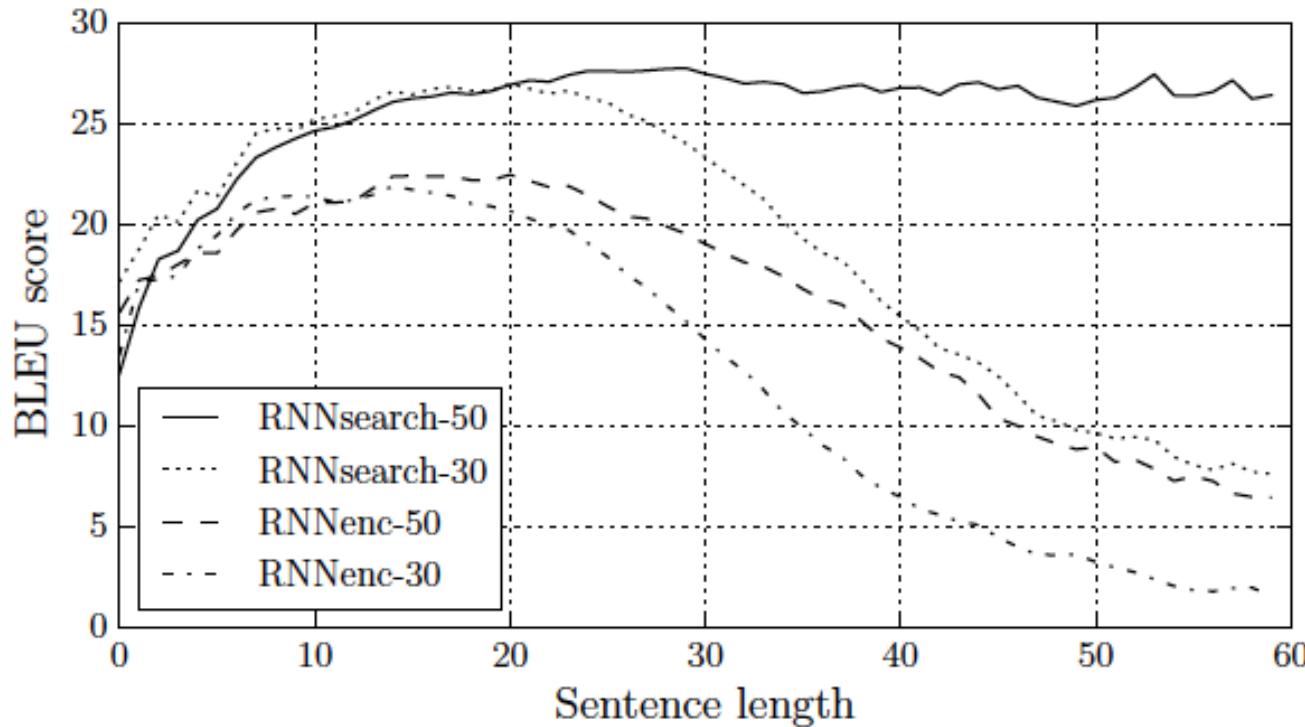
$$\mathbf{s}_{t+1} = g(\mathbf{s}_t, y_t)$$

$$e_{t,j} = a(\mathbf{s}_t, \mathbf{h}_j)$$



Jointly Learning to Align and Translate

- No performance drop for long sentences



"We train each model twice: first with the sentences of length up to 30 words (RNNencdec-30, RNNsearch-30) and then with the sentences of length up to 50 word (RNNencdec-50, RNNsearch-50)."

RNNencdec: Cho et al. 2014 <https://www.aclweb.org/anthology/D14-1179>

RNNsearch: Bahdanau et al, ICLR 2015, <https://arxiv.org/pdf/1409.0473.pdf>

OpenNMT

- <http://opennmt.net/>

Google Neural Machine Translation (Wu et al. (2016))

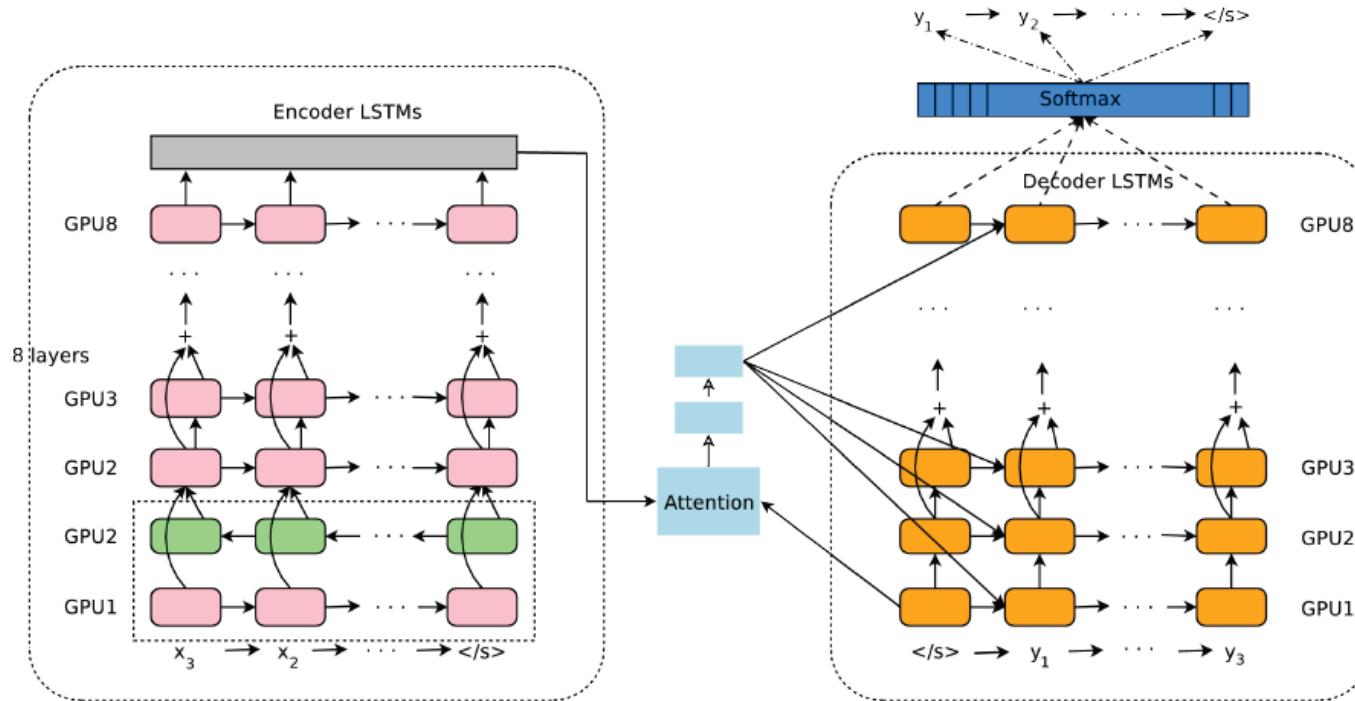
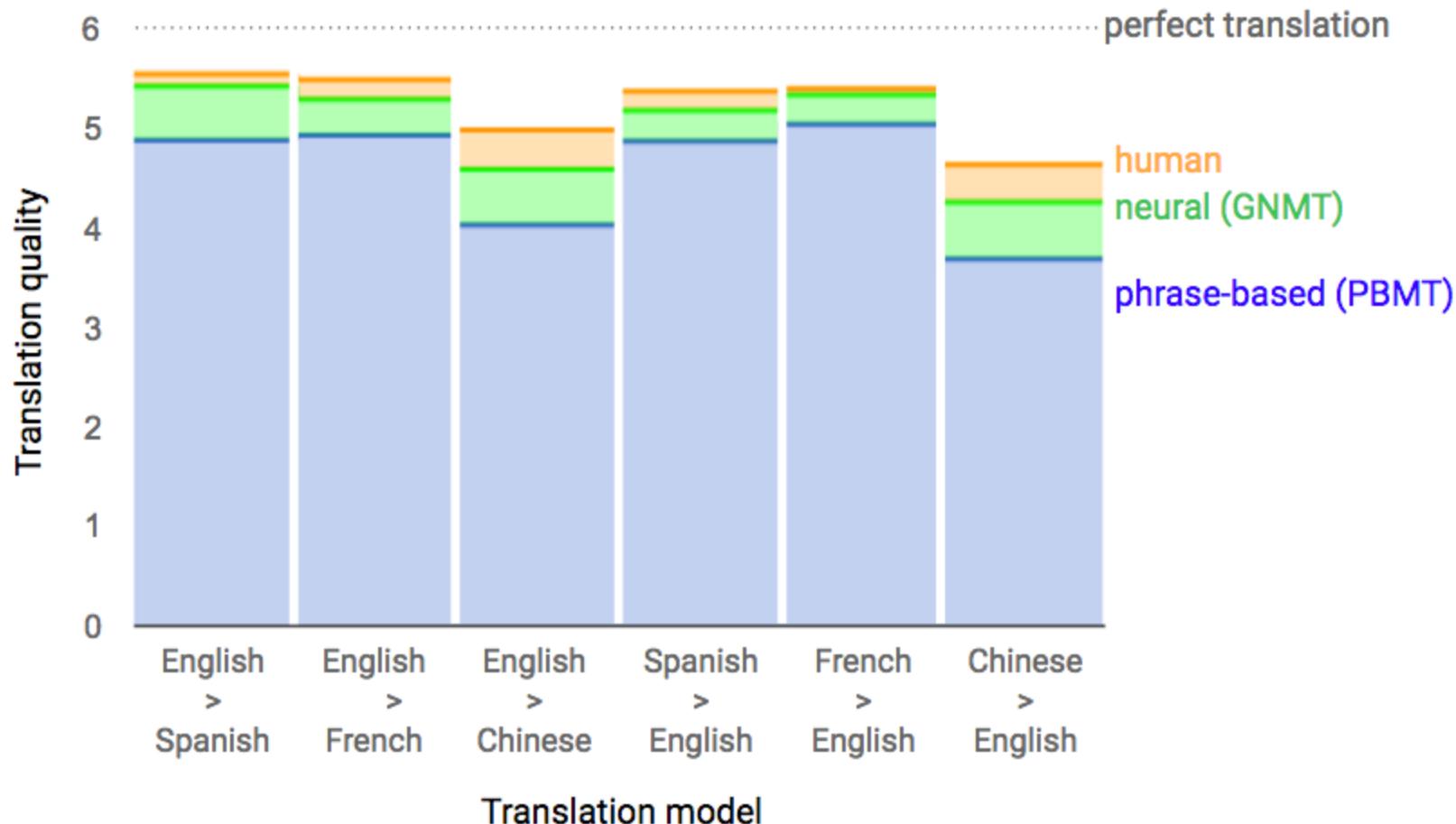
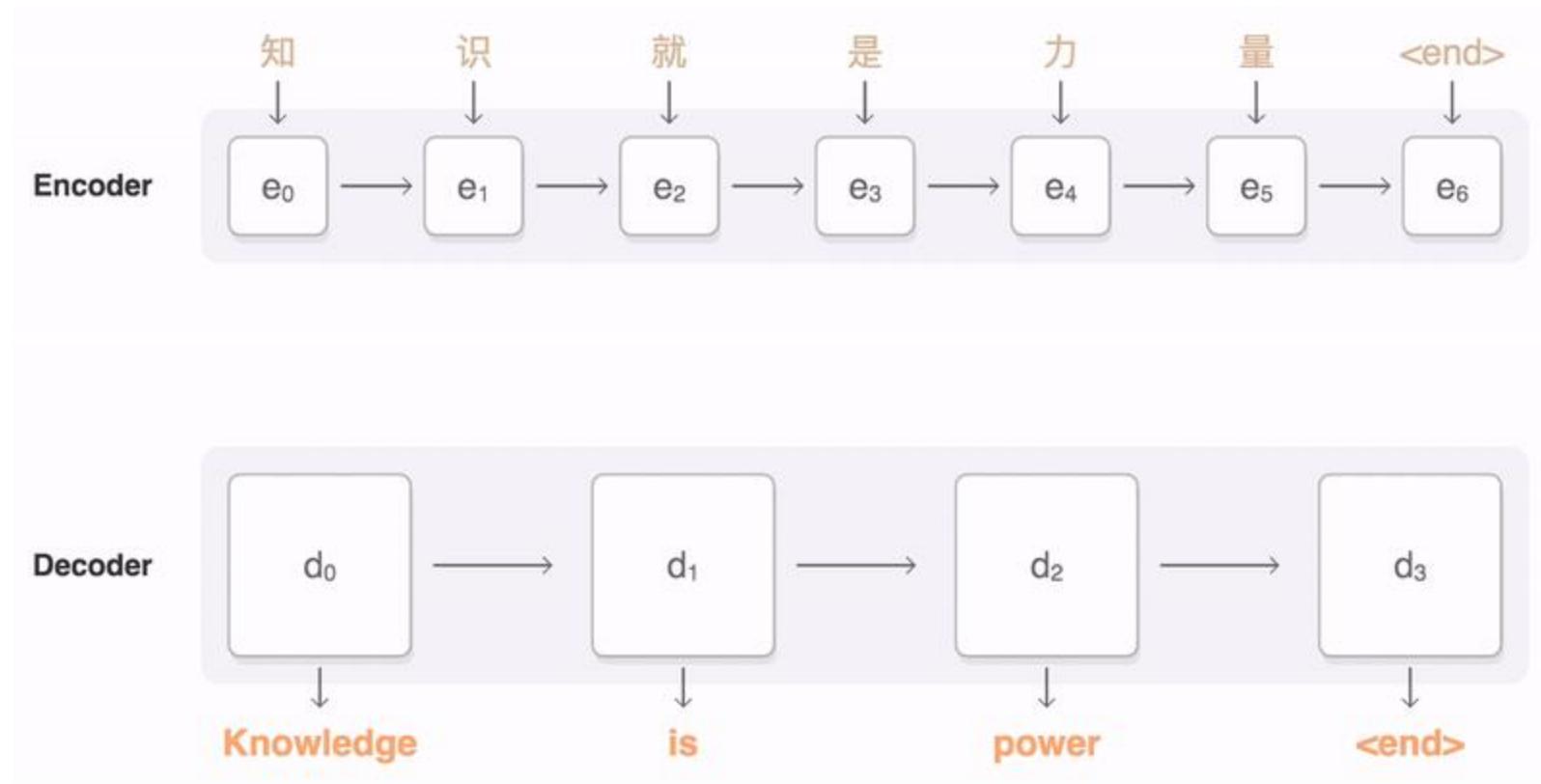


Figure 1: The model architecture of GNMT, Google’s Neural Machine Translation system. On the left is the encoder network, on the right is the decoder network, in the middle is the attention module. The bottom encoder layer is bi-directional: the pink nodes gather information from left to right while the green nodes gather information from right to left. The other layers of the encoder are uni-directional. Residual connections start from the layer third from the bottom in the encoder and decoder. The model is partitioned into multiple GPUs to speed up training. In our setup, we have 8 encoder LSTM layers (1 bi-directional layer and 7 uni-directional layers), and 8 decoder layers. With this setting, one model replica is partitioned 8-ways and is placed on 8 different GPUs typically belonging to one host machine. During training, the bottom bi-directional encoder layers compute in parallel first. Once both finish, the uni-directional encoder layers can start computing, each on a separate GPU. To retain as much parallelism as possible during running the decoder layers, we use the bottom decoder layer output only for obtaining recurrent attention context, which is sent directly to all the remaining decoder layers. The softmax layer is also partitioned and placed on multiple GPUs. Depending on the output vocabulary size we either have them run on the same GPUs as the encoder and decoder networks, or have them run on a separate set of dedicated GPUs.

Google Neural Machine Translation (Wu et al. (2016))

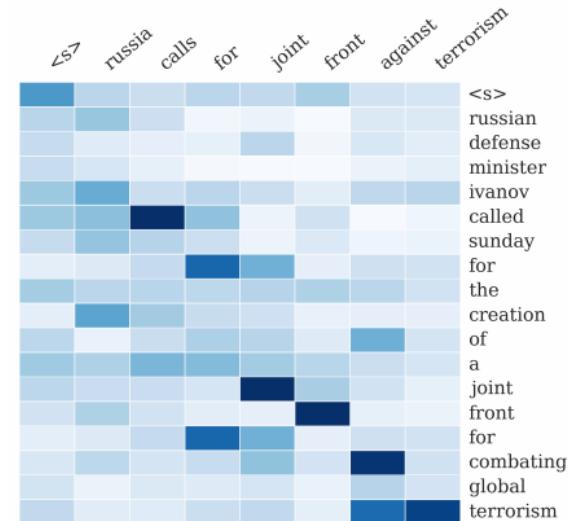


Google Neural Machine Translation (Wu et al. (2016))



Summary

- Attention significantly improves NMT performance
 - It's very useful to allow decoder to focus on certain parts of the source
- Attention solves the bottleneck problem
 - Attention allows decoder to look directly at source; bypass bottleneck
- Attention helps with vanishing gradient problem
 - Provides shortcut to faraway states
- Attention provides some interpretability
 - By inspecting attention distribution, we can see what the decoder was focusing on
 - We get alignment for free!
 - This is cool because we never explicitly trained an alignment system
 - The network just learned alignment by itself



Sequence-to-sequence is versatile!

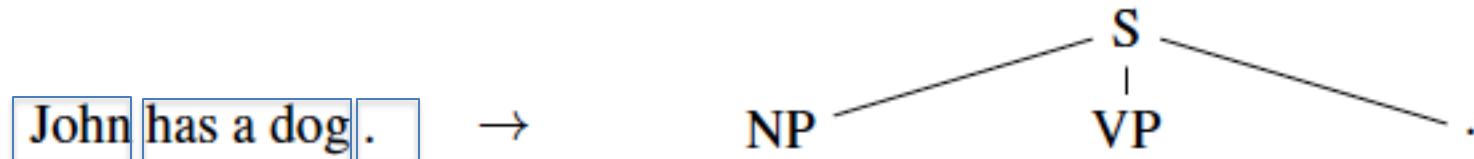
- Sequence-to-sequence is useful for *more than just MT*
- Many NLP tasks can be phrased as sequence-to-sequence:
 - Summarization (long text → short text)
 - Dialogue (previous utterances → next utterance)
 - Parsing (input text → output parse as sequence)
 - Code generation (natural language → Python code)

Grammar as a Foreign Language

- Vinyals et al., 2015

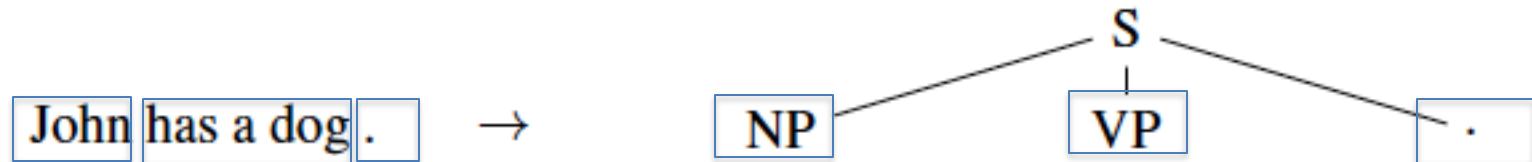
Grammar as a Foreign Language

Parsing tree



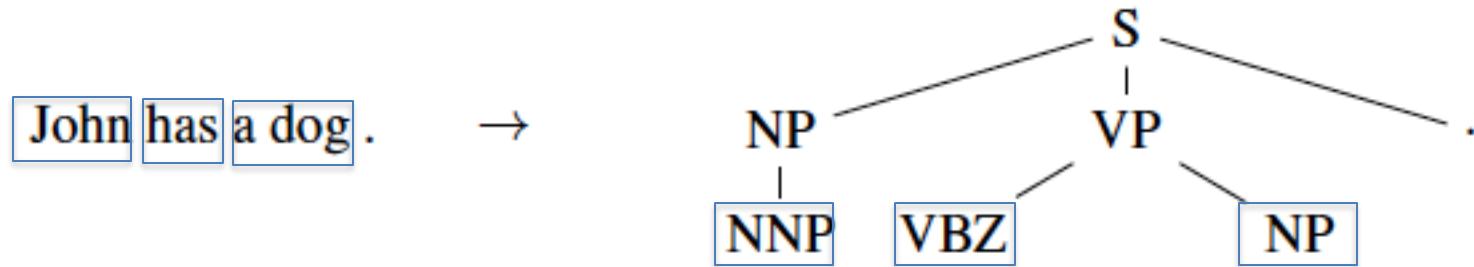
Grammar as a Foreign Language

Parsing tree



Grammar as a Foreign Language

Parsing tree

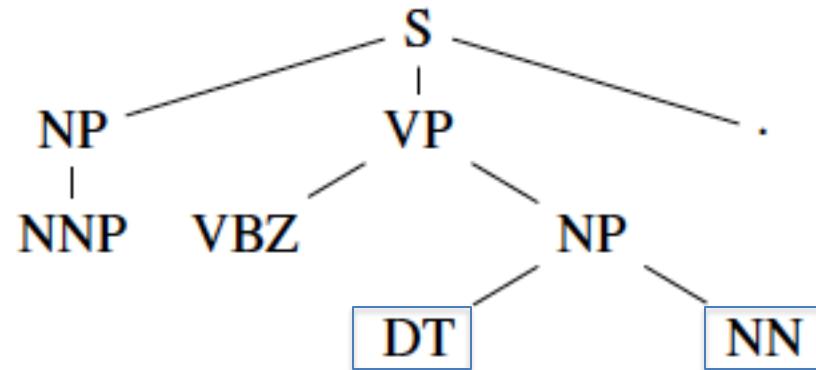


Grammar as a Foreign Language

Parsing tree

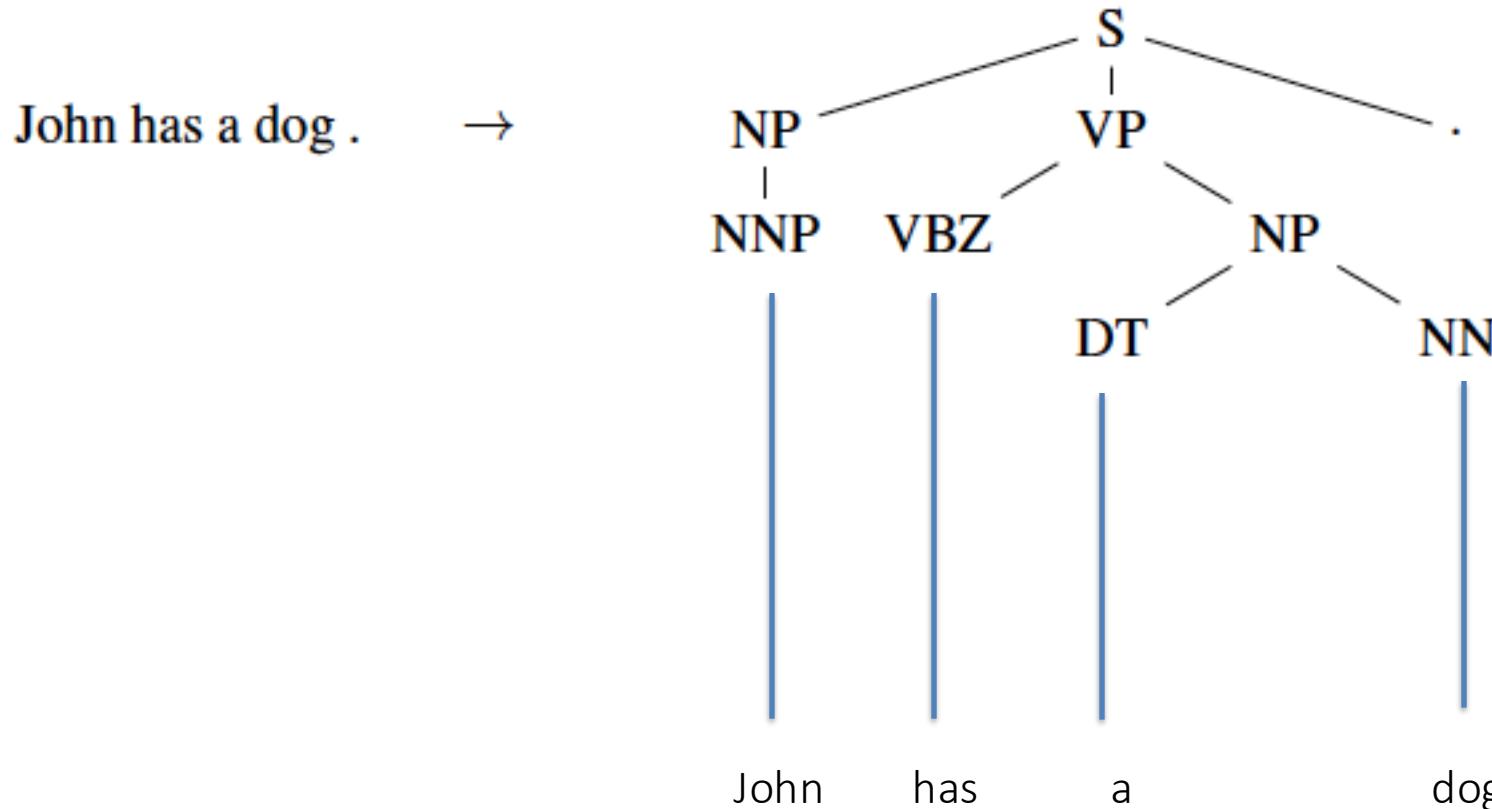
John has **a** **dog**.

→



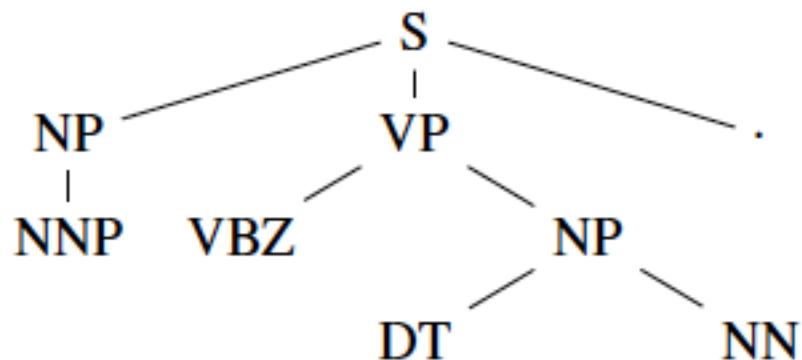
Grammar as a Foreign Language

Parsing tree



Grammar as a Foreign Language

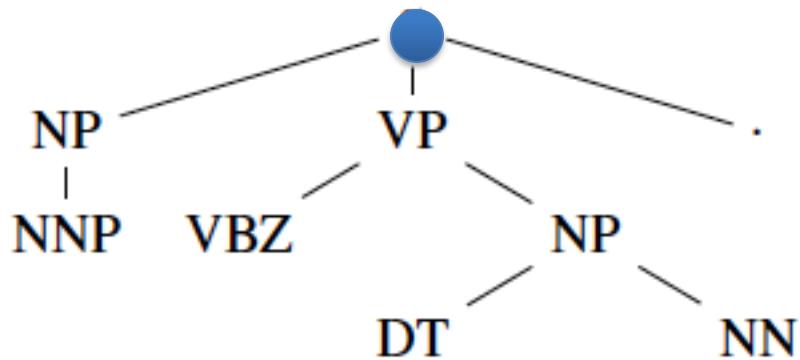
Converting tree to sequence



(S (NP NNP)_{NP} (VP VBZ (NP DT NN)_{NP})_{VP} .)_S

Grammar as a Foreign Language

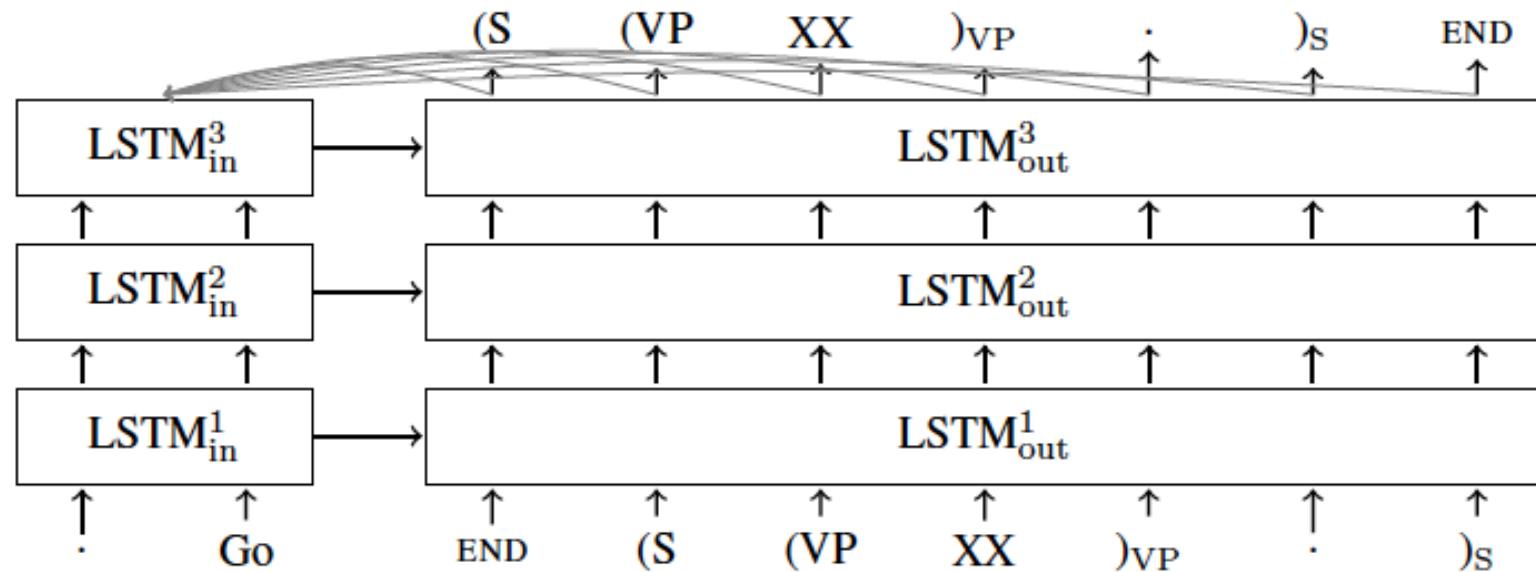
Converting tree to sequence



● (NP NNP)_{NP} (VP VBZ (NP DT NN)_{NP})_{VP} .)_S

Grammar as a Foreign Language

Model



Grammar as a Foreign Language

Results

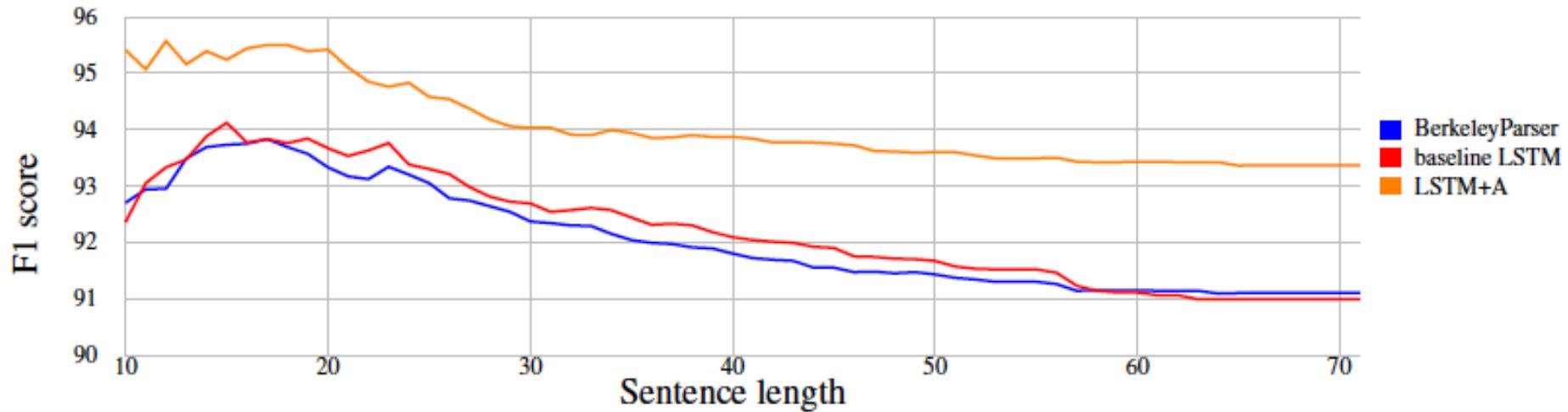


Image Captioning

- Neural image captioning (Donahue et al. (2014); Vinyals et al. (2015))
 - Encode image
 - Decode text



a cat is sitting on a toilet seat
logprob: -7.79



a display case filled with lots of different types of donuts
logprob: -7.78



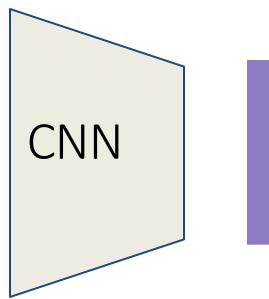
a group of people sitting at a table with wine glasses
logprob: -6.71



Image:
H x W x 3



Image:
 $H \times W \times 3$



Features:
D

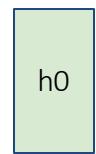
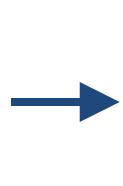
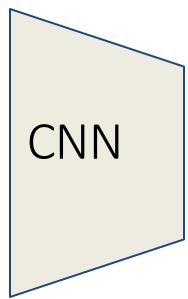


Image:
 $H \times W \times 3$

Features:
 D

Hidden
state: H

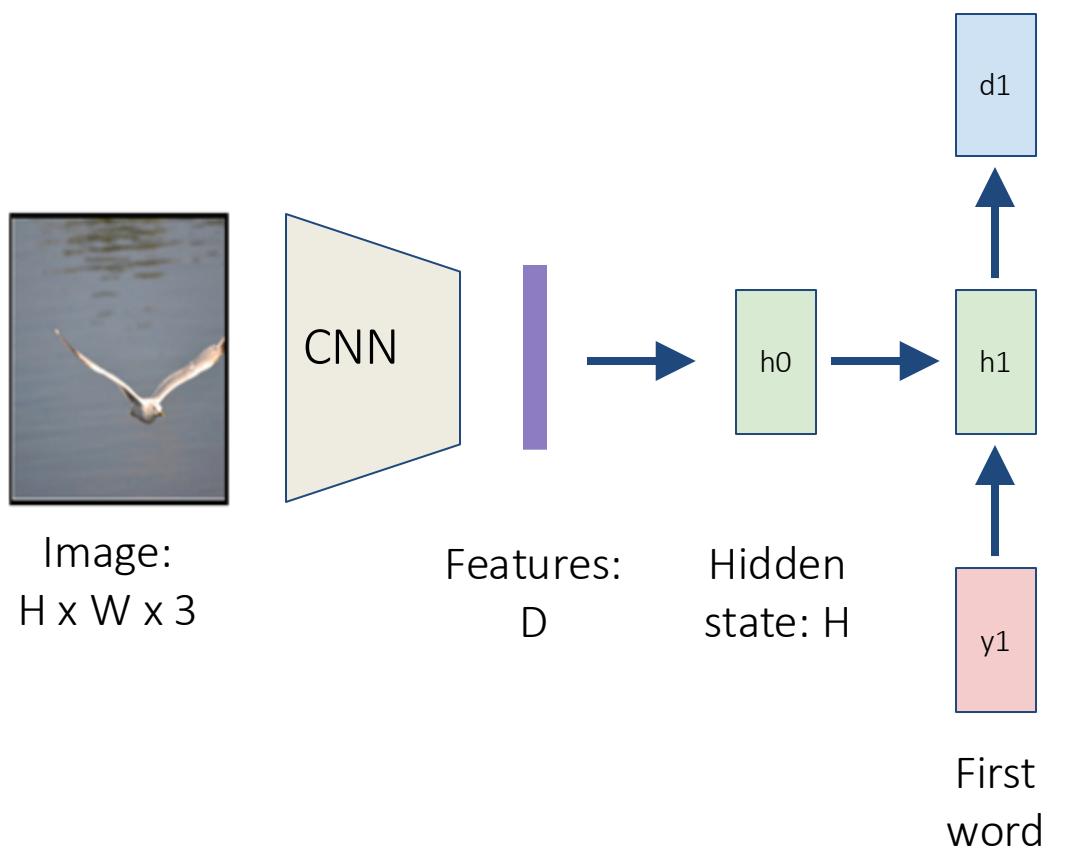




Image:
 $H \times W \times 3$

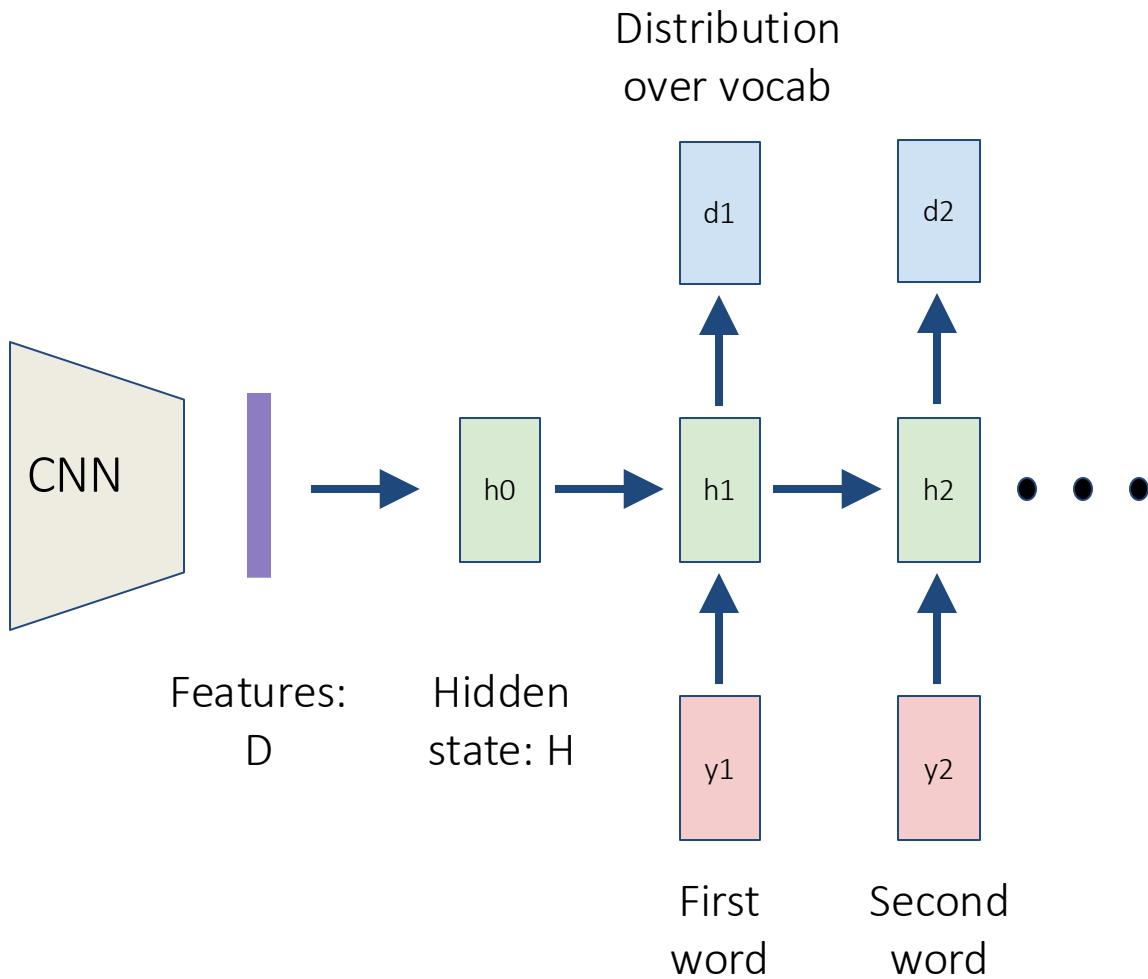
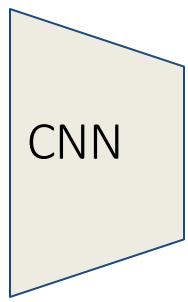




Image:
 $H \times W \times 3$



Features:
 D

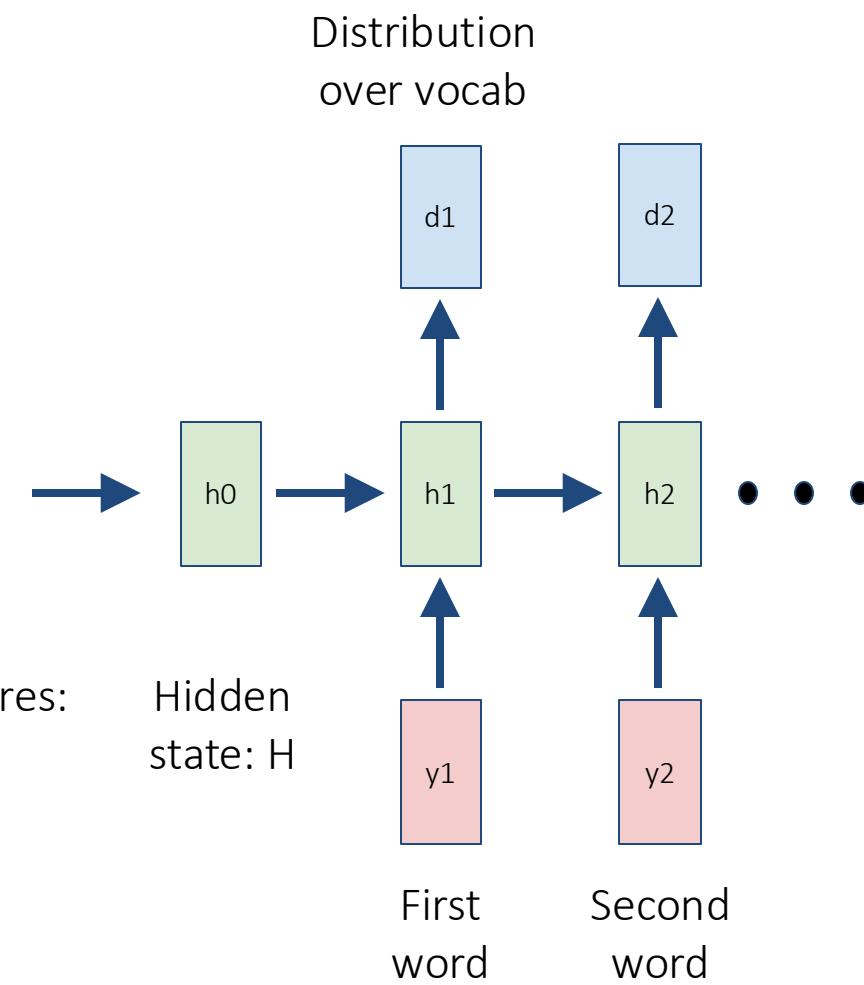
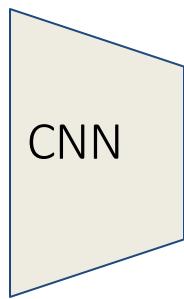




Image:
 $H \times W \times 3$

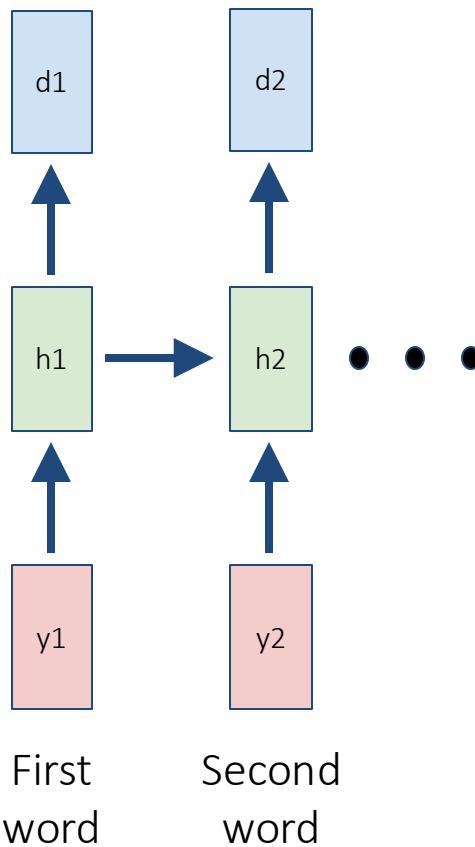


Features:
 D

→

Hidden
state: H

Distribution
over vocab

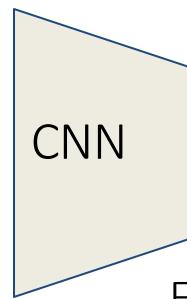


RNN only looks
at whole image,
once

What if the RNN
looks at different
parts of the
image at each
timestep?



Image:
 $H \times W \times 3$

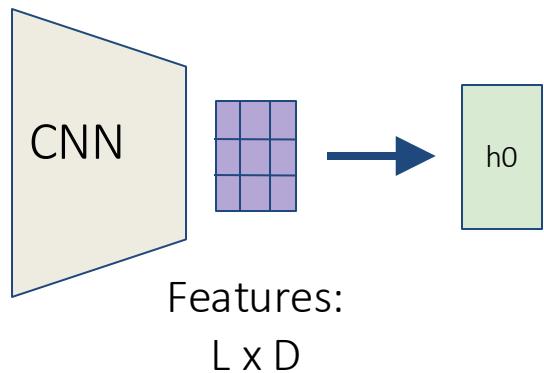


Features:
 $L \times D$

Xu et al, "Show, Attend and Tell:
Neural Image Caption Generation with
Visual Attention", ICML 2015

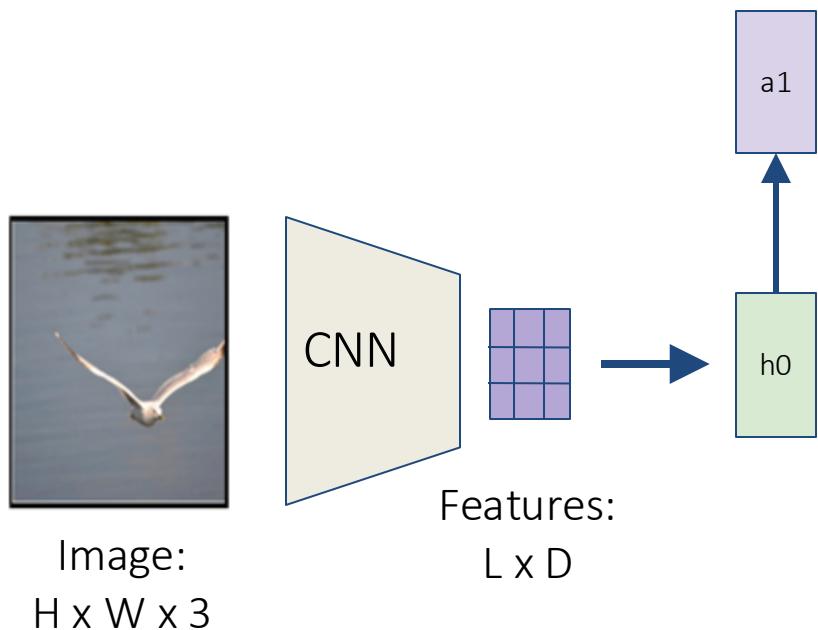


Image:
 $H \times W \times 3$



Xu et al, "Show, Attend and Tell:
Neural Image Caption Generation with
Visual Attention", ICML 2015

Distribution over L
locations

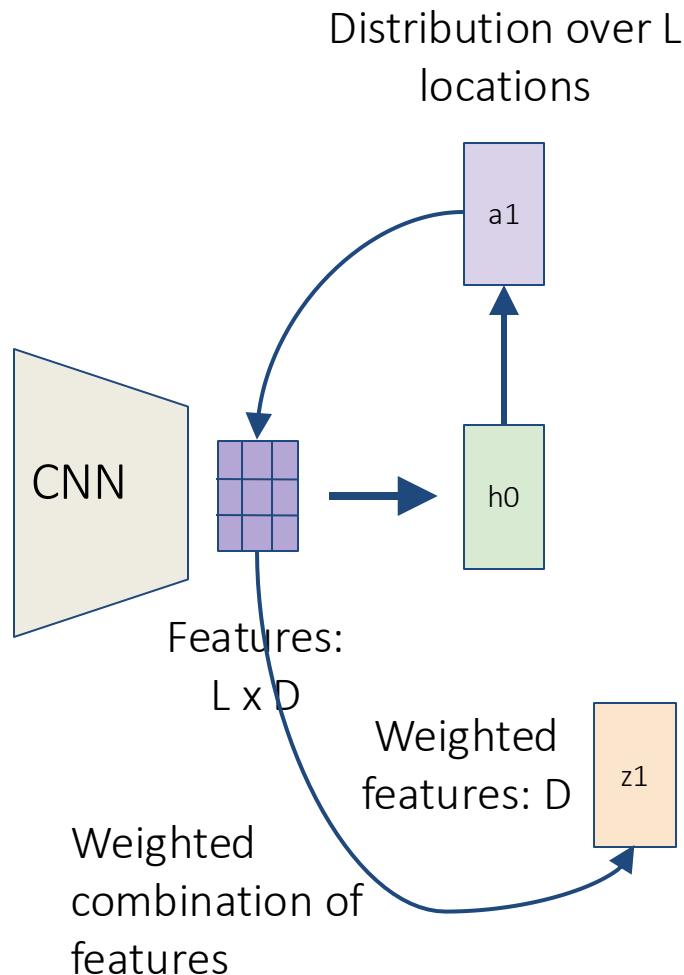


Xu et al, "Show, Attend and Tell:
Neural Image Caption Generation with
Visual Attention", ICML 2015

Image



Image:
 $H \times W \times 3$



Distribution over
L locations

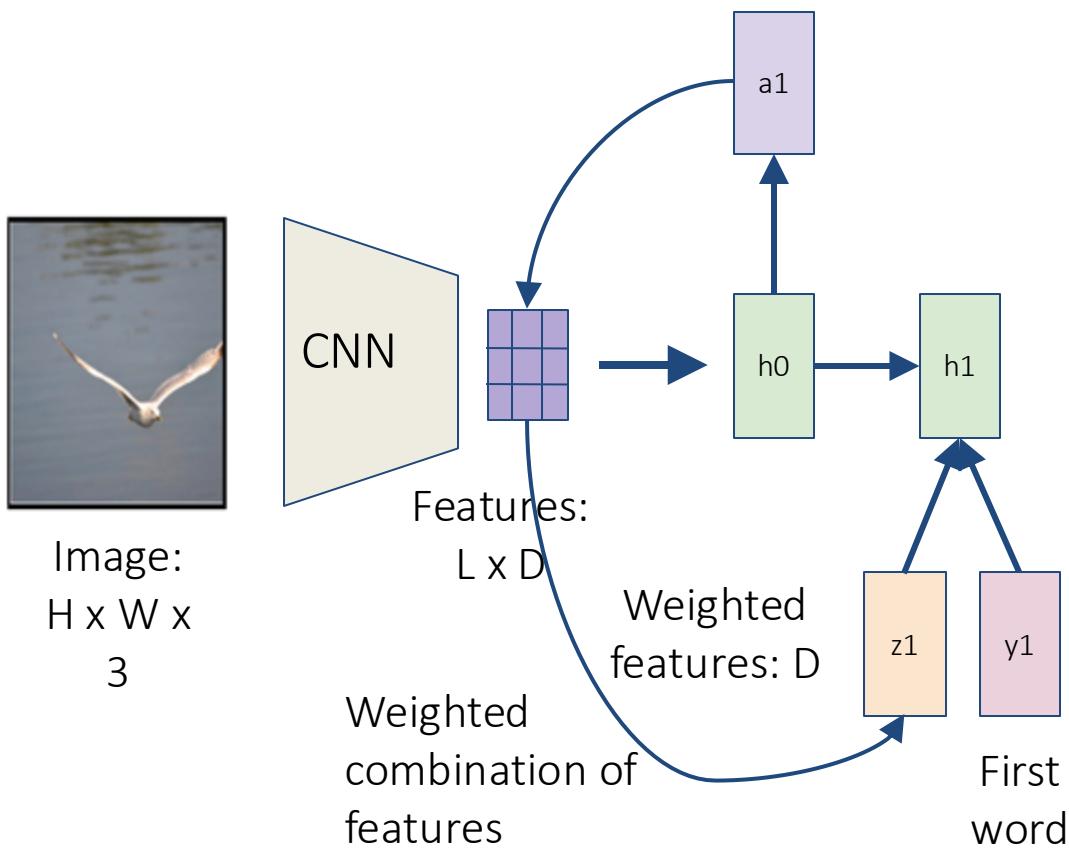




Image:
 $H \times W \times 3$

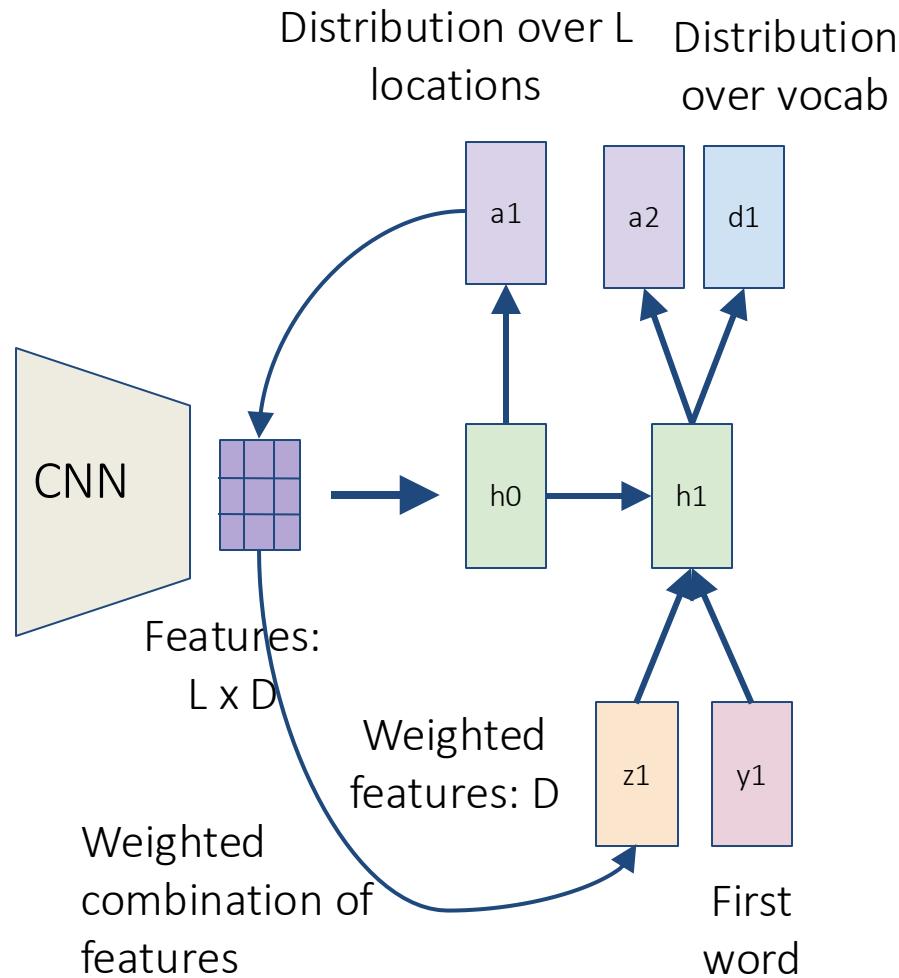
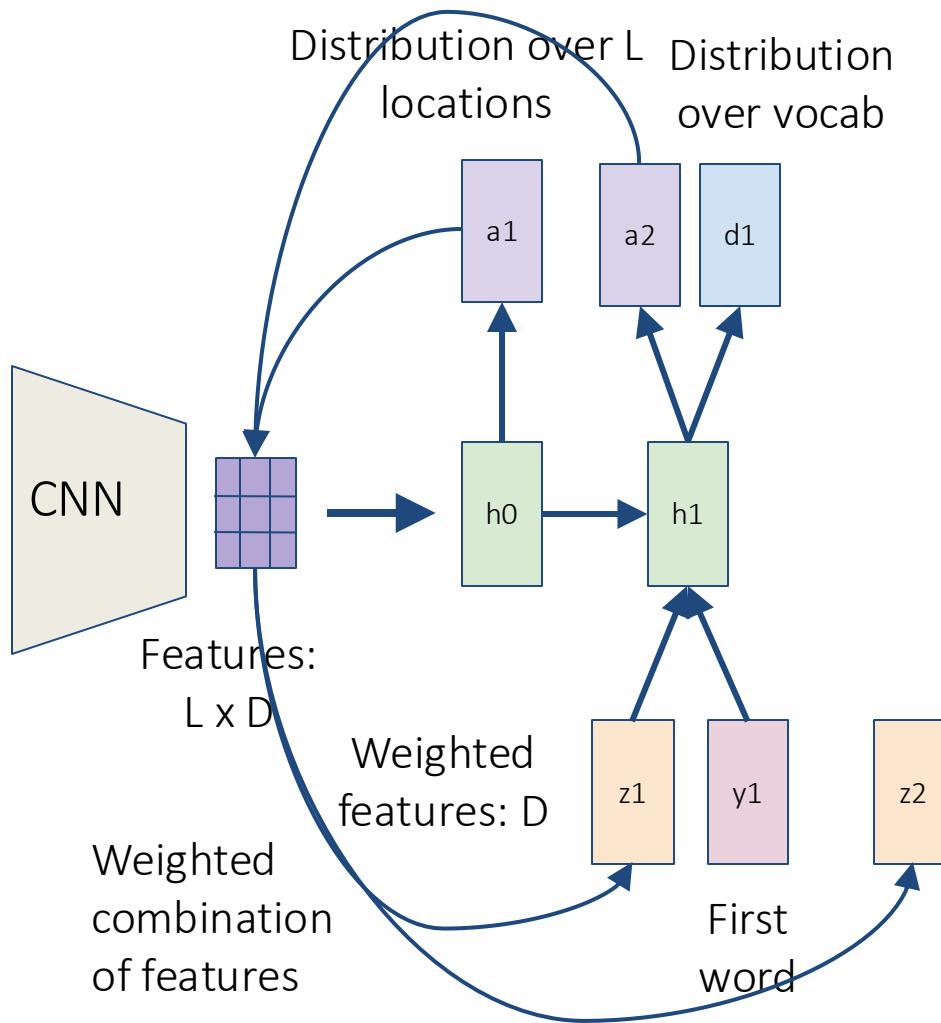




Image:
 $H \times W \times 3$



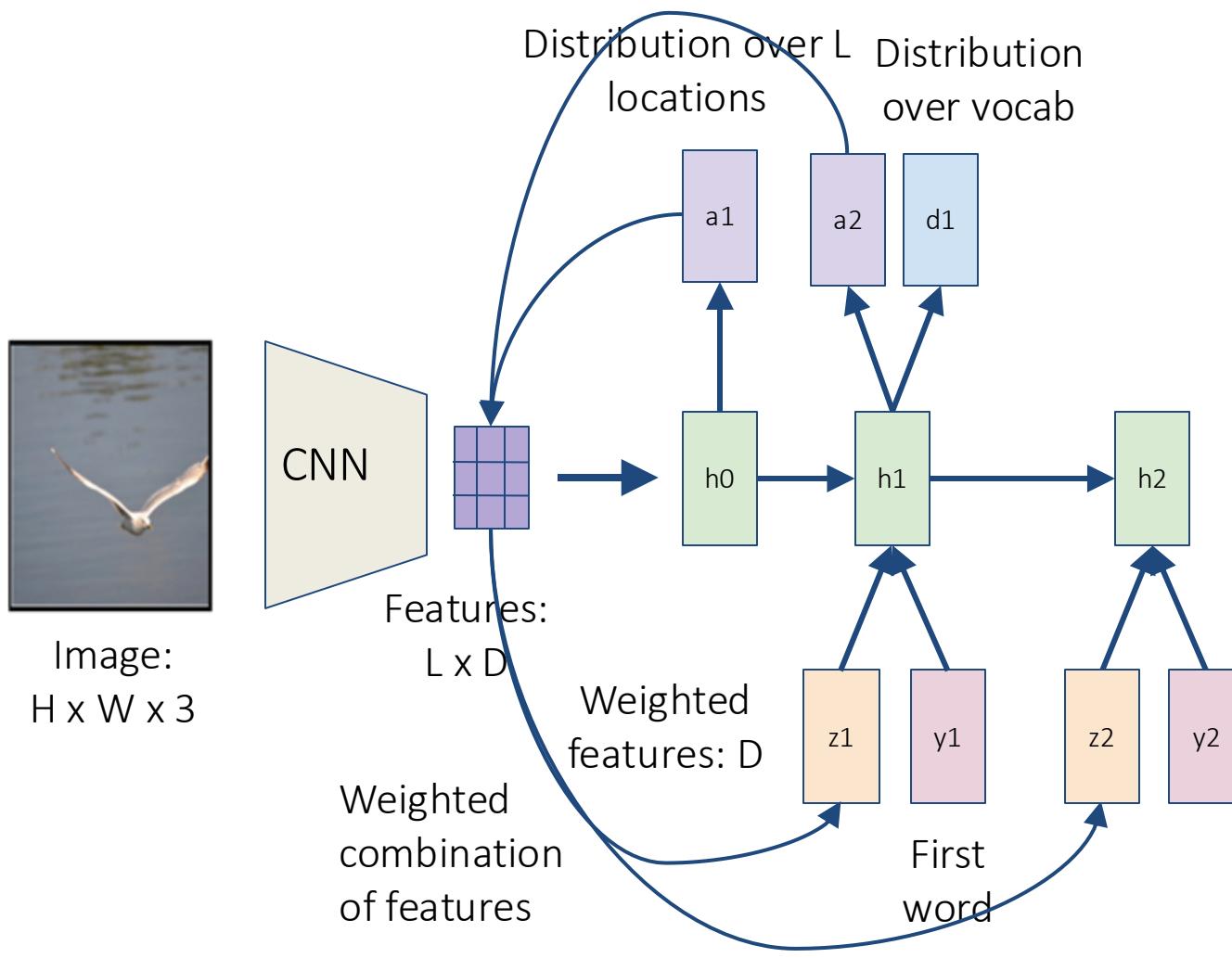
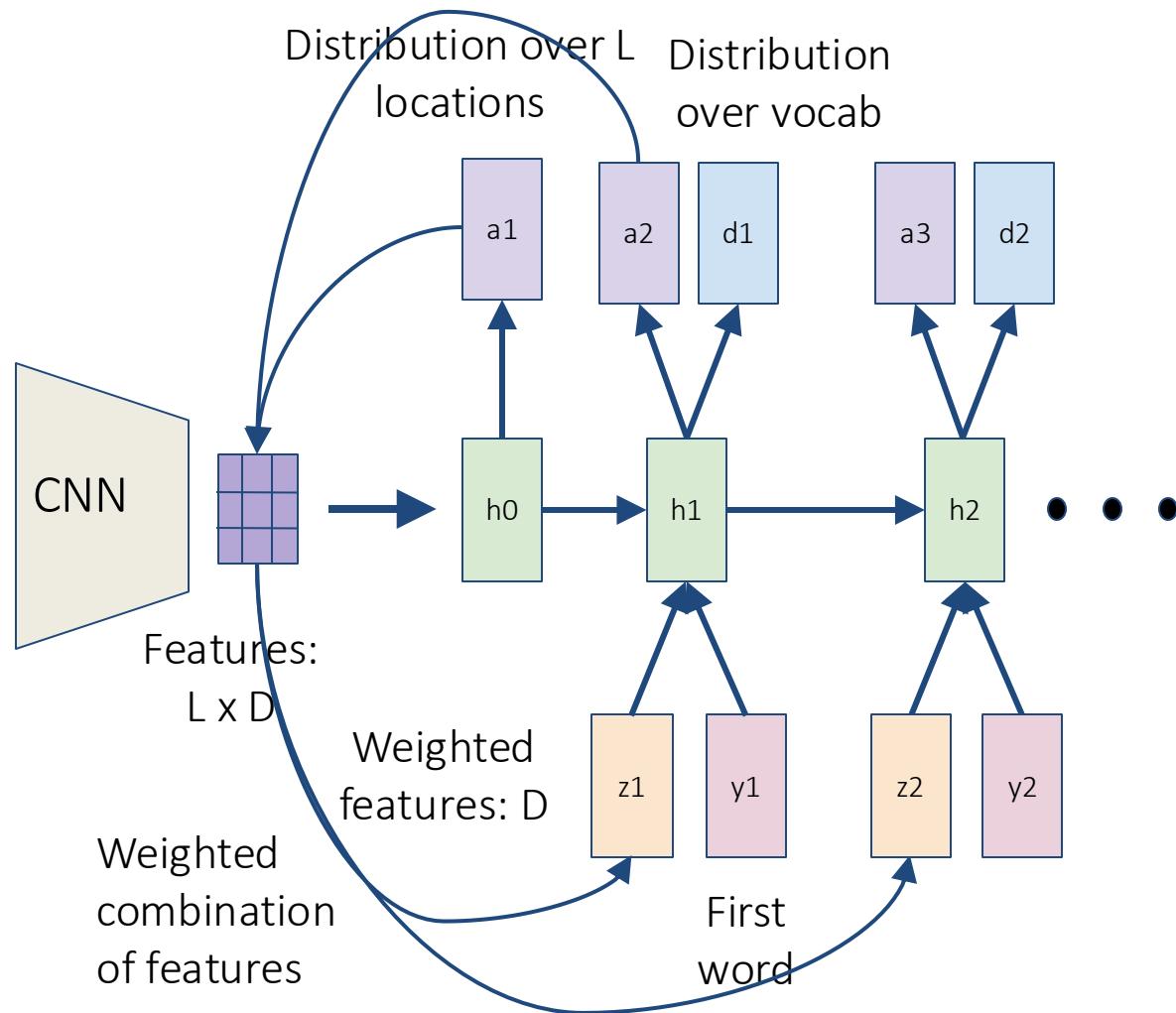




Image:
 $H \times W \times 3$



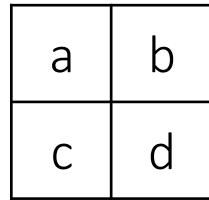
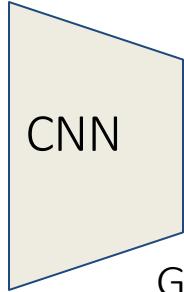
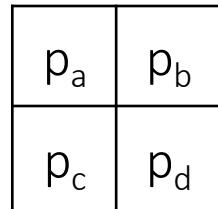


Image:
 $H \times W \times 3$

Grid of features
(Each D-dimensional)

From
RNN:

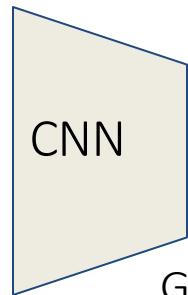


Xu et al, "Show, Attend and Tell:
Neural Image Caption Generation with
Visual Attention", ICML 2015

Distribution over grid
locations
 $p_a + p_b + p_c + p_d = 1$



Image:
 $H \times W \times 3$



Grid of features
(Each D-dimensional)

a	b
c	d

Context vector z
(D-dimensional)

From
RNN:

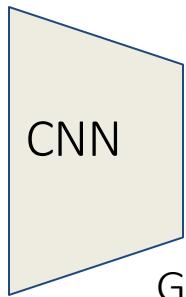
p_a	p_b
p_c	p_d

Xu et al, "Show, Attend and Tell:
Neural Image Caption Generation with
Visual Attention", ICML 2015

Distribution over grid
locations
 $p_a + p_b + p_c + p_d = 1$



Image:
 $H \times W \times 3$



Grid of features
(Each D-dimensional)

a	b
c	d

Soft attention:
Summarize ALL locations
$$z = p_a a + p_b b + p_c c + p_d d$$

Derivative dz/dp is nice!
Train with gradient descent

Context vector z
(D-dimensional)

From
RNN:

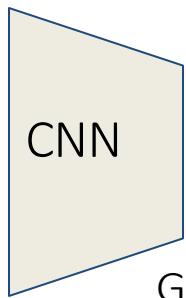
p_a	p_b
p_c	p_d

Xu et al, "Show, Attend and Tell:
Neural Image Caption Generation with
Visual Attention", ICML 2015

Distribution over grid
locations
$$p_a + p_b + p_c + p_d = 1$$



Image:
 $H \times W \times 3$



Grid of features
(Each D-dimensional)

a	b
c	d

From
RNN:

p_a	p_b
p_c	p_d

Distribution over grid
locations
 $p_a + p_b + p_c + p_d = 1$

Xu et al, "Show, Attend and Tell:
Neural Image Caption Generation with
Visual Attention", ICML 2015

Context vector z
(D-dimensional)

Soft attention:
Summarize ALL locations
$$z = p_a a + p_b b + p_c c + p_d d$$

Derivative dz/dp is nice!
Train with gradient descent

Hard attention:
Sample ONE location
according to p , $z =$ that vector

With argmax , dz/dp is zero
almost everywhere ...
Can't use gradient descent;
need reinforcement learning

Soft
attention



Hard
attention



A

bird

flying

over

a

body

of

water

.



A woman is throwing a frisbee in a park.



A dog is standing on a hardwood floor.



A stop sign is on a road with a mountain in the background.



A little girl sitting on a bed with a teddy bear.



A group of people sitting on a boat in the water.



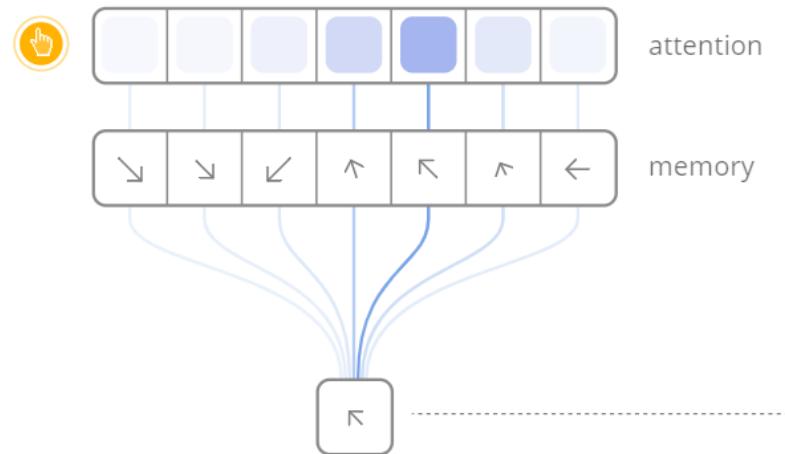
A giraffe standing in a forest with trees in the background.

What did we learn today?

- Attention can be quite useful in a variety of applications
- Attention = a form of memory access and can help with capturing long-term dependencies!

Attention is a *general* deep Learning technique

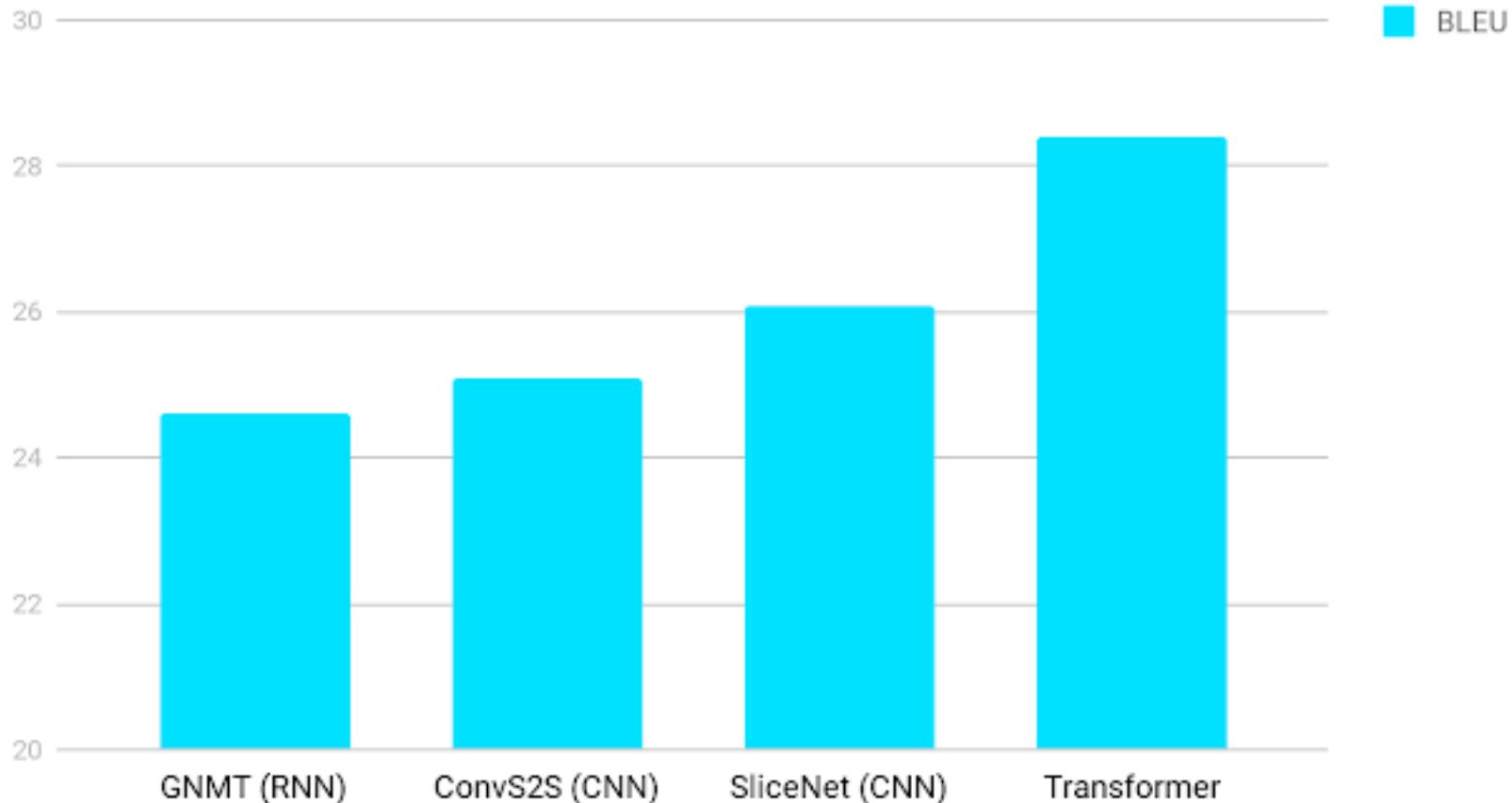
- More general definition of attention
 - Given a set of vector *values*, and a vector *query*, **attention** is a technique to compute a **weighted sum of the values**, dependent on the query.
 - We sometimes say that the query *attends to* the values
 - For example, in the seq2seq + attention model, each decoder hidden state *attends to* the encoder hidden states



<https://distill.pub/2016/augmented-rnns/>

Attention is all you need? – Transformer Networks

English German Translation quality



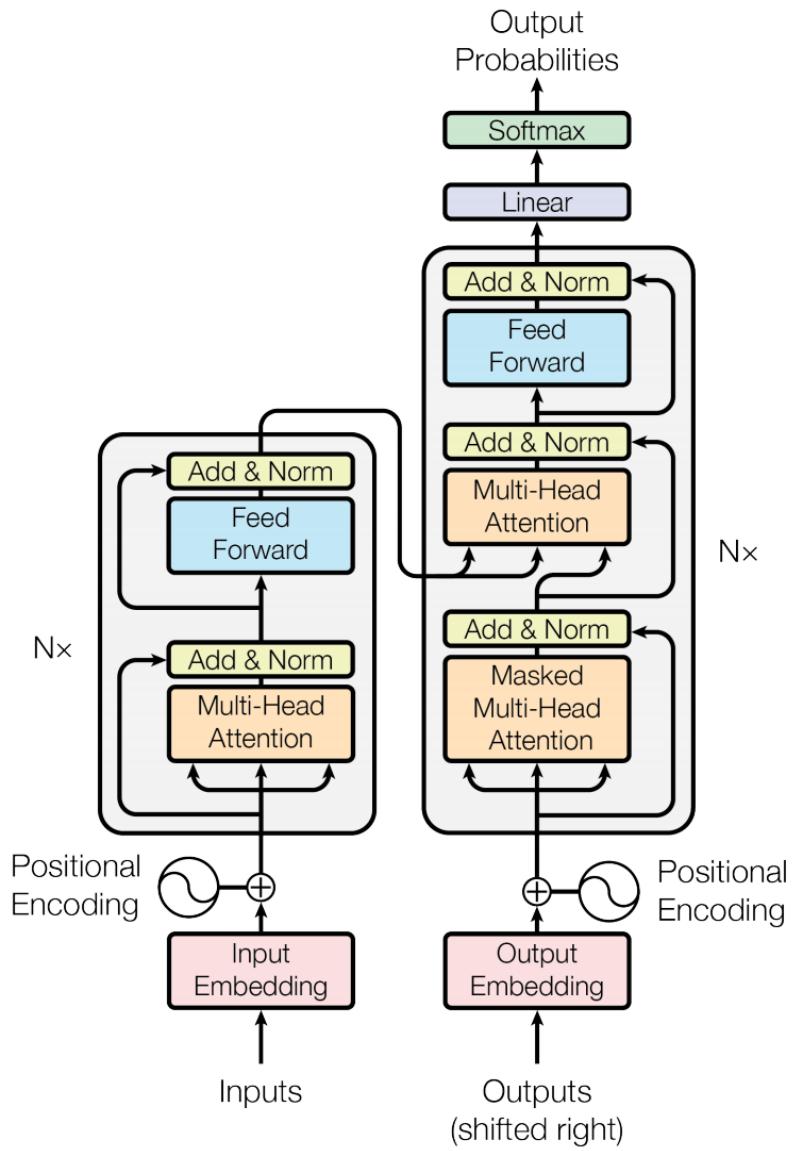
<https://ai.googleblog.com/2017/08/transformer-novel-neural-network.html>

Problems with RNNs = Motivation for Transformers

- Recurrent models typically factor computation along the symbol positions of the input and output sequences
 - Sequential computation **prevents parallelization**
 - Critical at **longer sequence lengths**, as memory constraints limit batching across examples
- Despite advanced RNNs like LSTMs, RNNs still need **attention mechanism** to deal with **long range dependencies** – path length for codependent computation between states grows with sequence
- But if attention gives us access to any state... maybe we don't need the RNN?

Transformer Overview

- Sequence-to-sequence
- Encoder-Decoder
- Task: machine translation with parallel corpus
- Predict each translated word
- Final cost/error function is standard cross-entropy error on top of a softmax classifier



Transformer Basics

- Let's define the basic building blocks of transformer networks first: new attention layers!

Dot-Product Attention (Extending our previous def.)

- Inputs: a query q and a set of key-value (k-v) pairs to an output
 - Query, keys, values, and output are all vectors
- Output is weighted sum of values, where
 - Weight of each value is computed by an inner product of query and corresponding key
 - Queries and keys have same dimensionality

$$A(q, K, V) = \sum_i \frac{e^{q \cdot k_i}}{\sum_j e^{q \cdot k_j}} v_i$$

Dot-Product Attention – Matrix notation

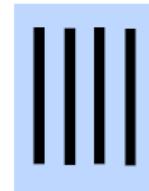
- When we have multiple queries q , we stack them in a matrix Q :

$$A(q, K, V) = \sum_i \frac{e^{q \cdot k_i}}{\sum_j e^{q \cdot k_j}} v_i$$

- Becomes: $A(Q, K, V) = \text{softmax}(QK^T)V$

$$[|Q| \times d_k] \times [d_k \times |K|] \times [|K| \times d_v]$$

softmax
row-wise

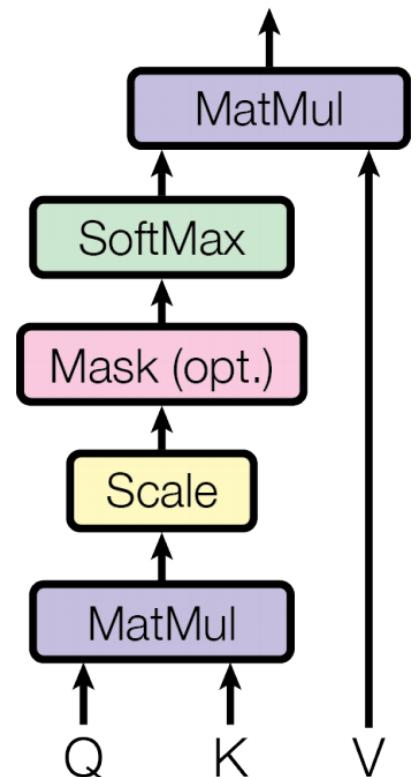


$$= [|Q| \times d_v]$$

Scaled Dot-Product Attention

- Problem: As d_k gets large, the variance of $q^T k$ increases → some values inside the softmax get large → the softmax gets very peaked → hence its gradient gets smaller.
- Solution: Scale by length of query/key vectors:

$$A(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

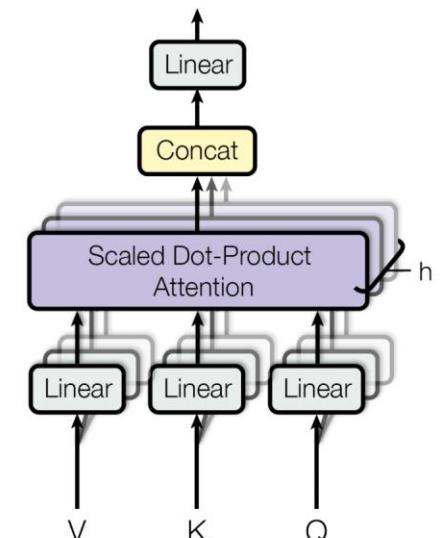


Self-attention and Multi-head attention

- The input word vectors could be the queries, keys and values
 - In other words: the word vectors themselves select each other
 - Word vector stack = $Q = K = V$
- Problem: Only one way for words to interact with one-another
- Solution: Multi-head attention
 - First map Q, K, V into h many lower dimensional spaces via W matrices
 - Then apply attention, then concatenate outputs and pipe through linear layer

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$

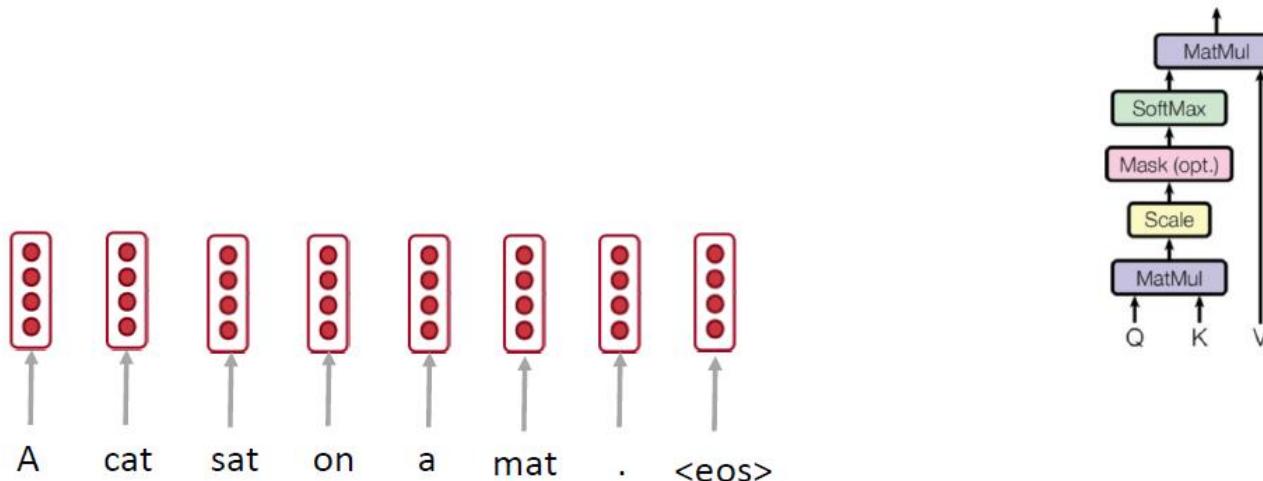
$$\text{where } \text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$



Self-attention: A Running Example

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

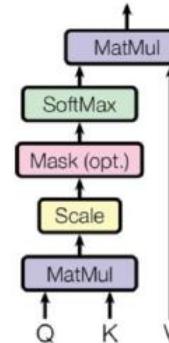
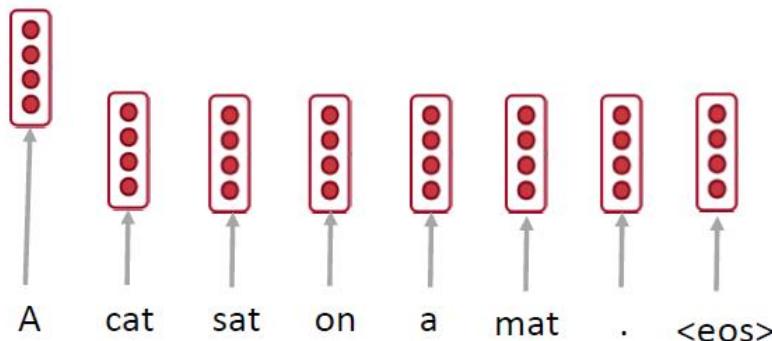
Scaled Dot-Product Attention



Self-attention: A Running Example

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

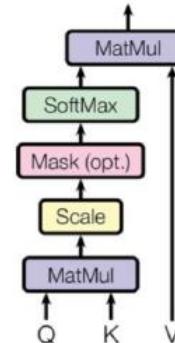
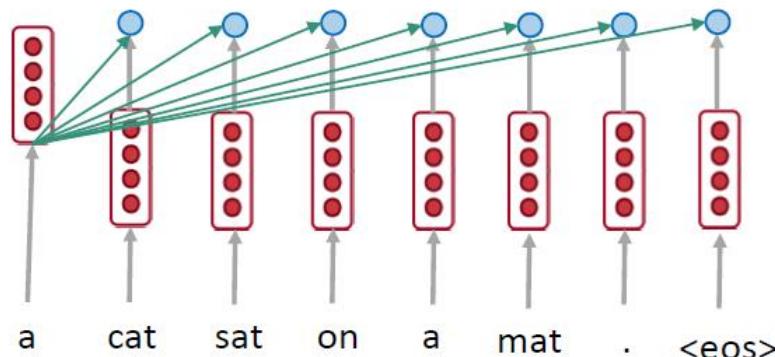
Scaled Dot-Product Attention



Self-attention: A Running Example

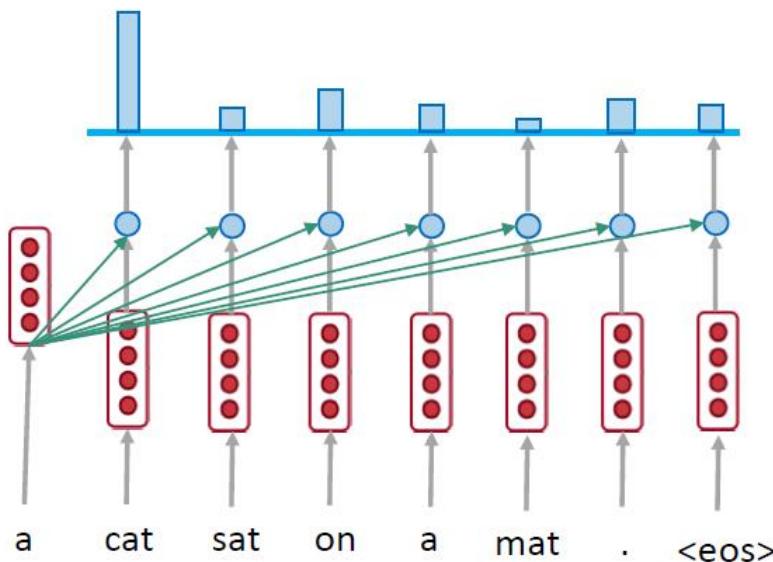
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Scaled Dot-Product Attention

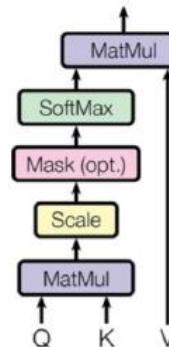


Self-attention: A Running Example

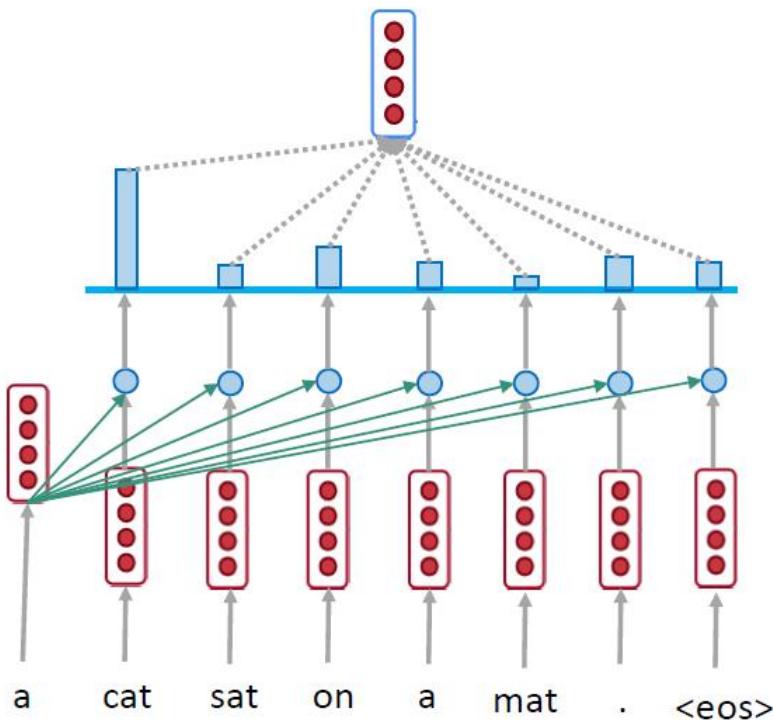
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$



Scaled Dot-Product Attention

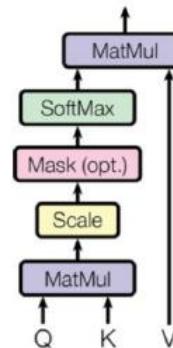


Self-attention: A Running Example

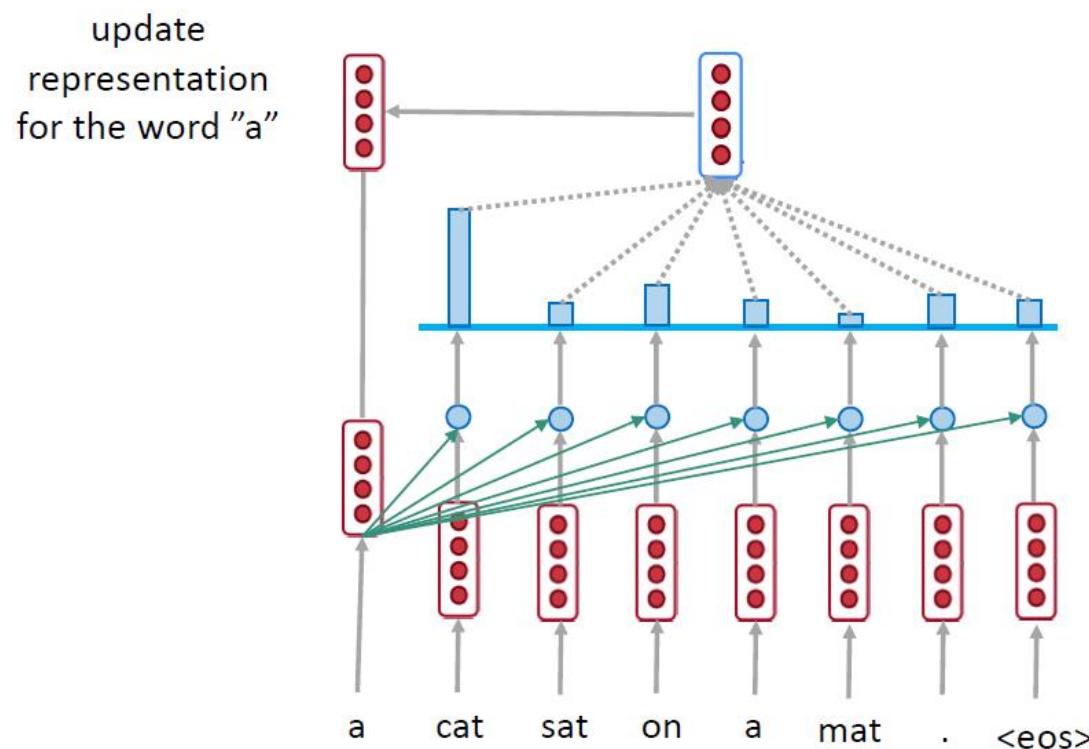


$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Scaled Dot-Product Attention

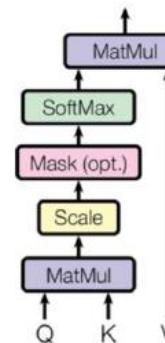


Self-attention: A Running Example



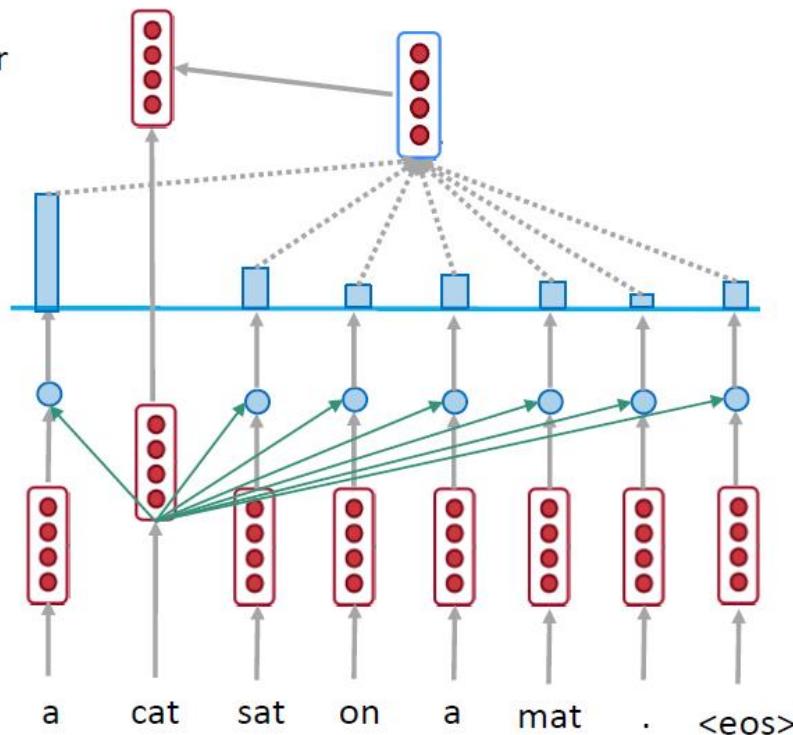
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Scaled Dot-Product Attention



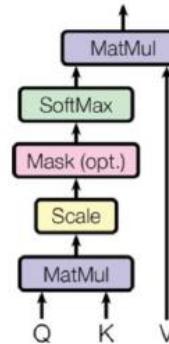
Self-attention: A Running Example

update
representation for
the word "cat"



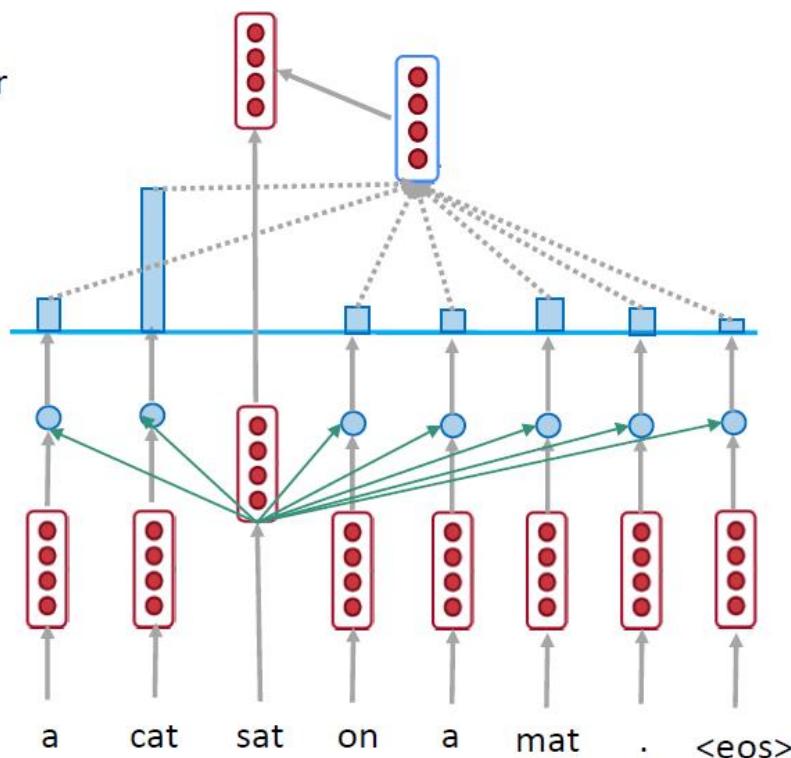
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Scaled Dot-Product Attention



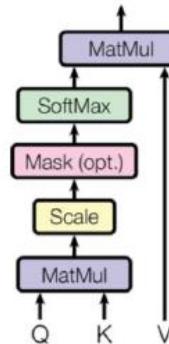
Self-attention: A Running Example

update
representation for
the word "sat"



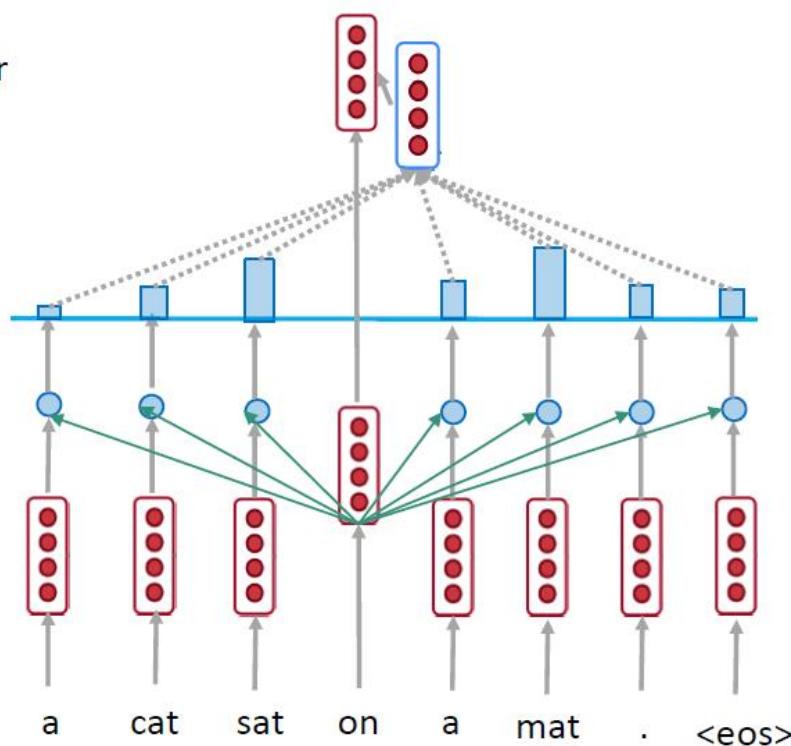
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Scaled Dot-Product Attention



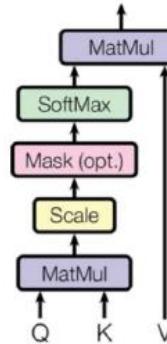
Self-attention: A Running Example

update
representation for
the word "on"



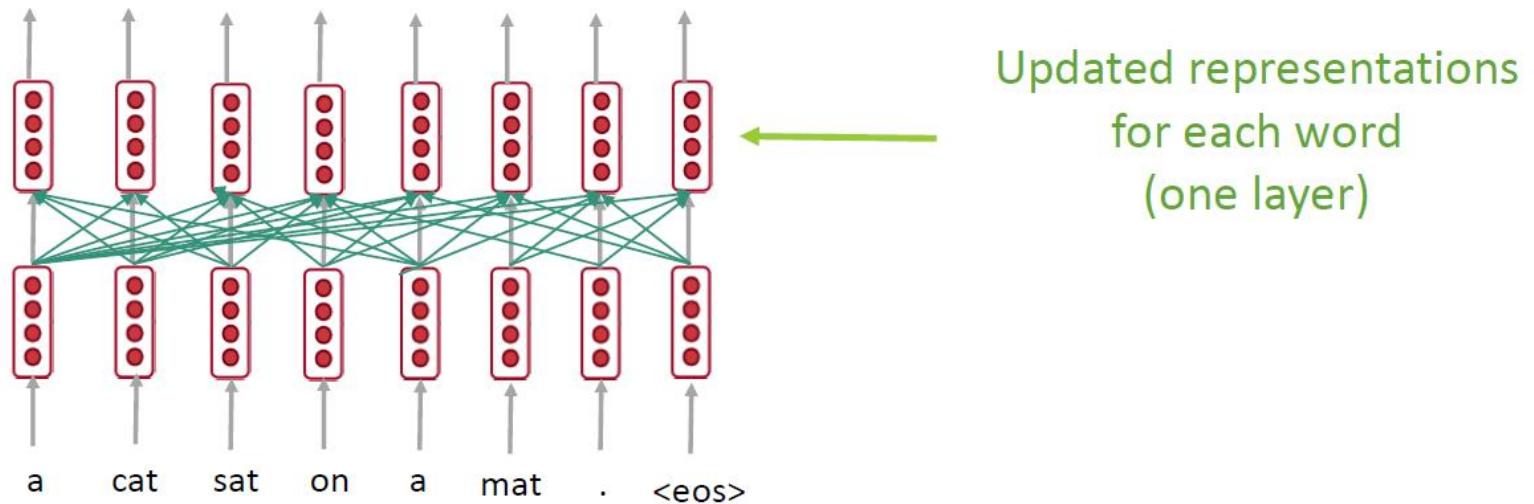
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Scaled Dot-Product Attention



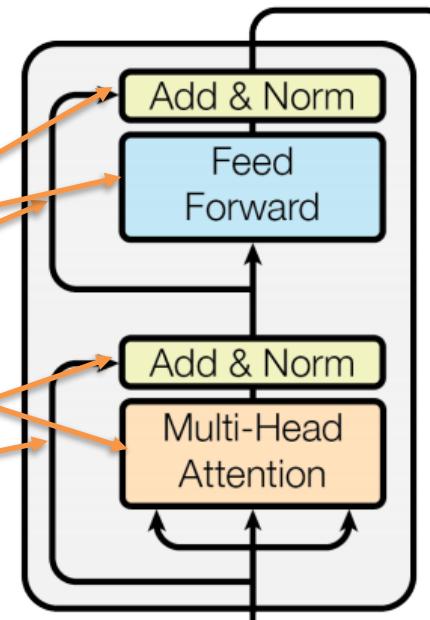
Self-attention: A Running Example

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$



Complete transformer block

- Each block has two “sublayers”
 - 1. Multihead attention
 - 2. 2 layer feed-forward Nnet (with relu)
- Each of these two steps also has:
 - Residual (short-circuit) connection and LayerNorm:
 - $\text{LayerNorm}(x + \text{Sublayer}(x))$
 - Layernorm changes input to have mean 0 and variance 1, per layer and per training point (and adds two more parameters)



$$\mu^l = \frac{1}{H} \sum_{i=1}^H a_i^l \quad \sigma^l = \sqrt{\frac{1}{H} \sum_{i=1}^H (a_i^l - \mu^l)^2} \quad h_i = f\left(\frac{g_i}{\sigma_i} (a_i - \mu_i) + b_i\right)$$

New Fun Fact (Found in May 2023)

Post-LN Transformer

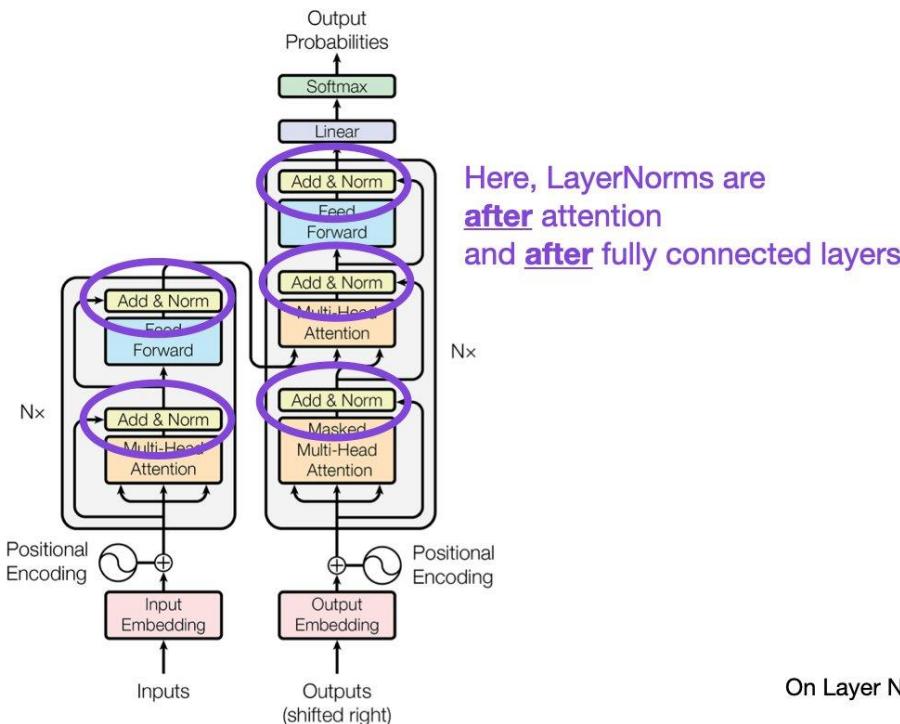
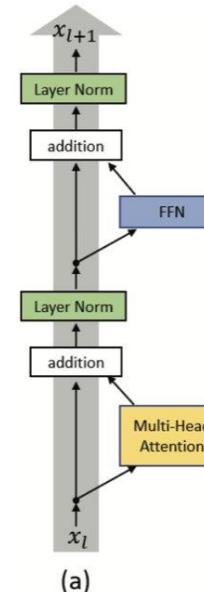


Figure 1: The Transformer - model architecture.

Attention Is All You Need, <https://arxiv.org/abs/1706.03762>

Places the layer normalization between the residual blocks: the expected gradients of the parameters near the output layer are large

Paper: addition + layer norm
Code: dropout + addition + layer norm



On Layer Normalization in the Transformer Architecture,
<https://arxiv.org/abs/2002.04745>

 **Sebastian Raschka**  [@rasbt](#) · May 8

Or maybe to be fair, the "original original" code did match the figure, but then they submitted this PR in 2017 to change it, while not updating the figure. So yeah, it's utterly confusing.

1 23 4,473

New Fun Fact (Found in May 2023)

Better gradients if LayerNorm is placed inside residual connection, before the attention and fully connected layers

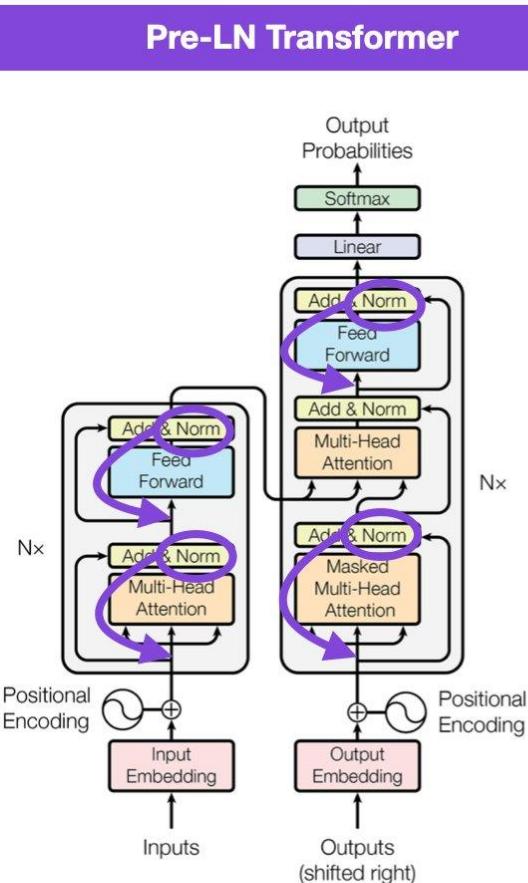
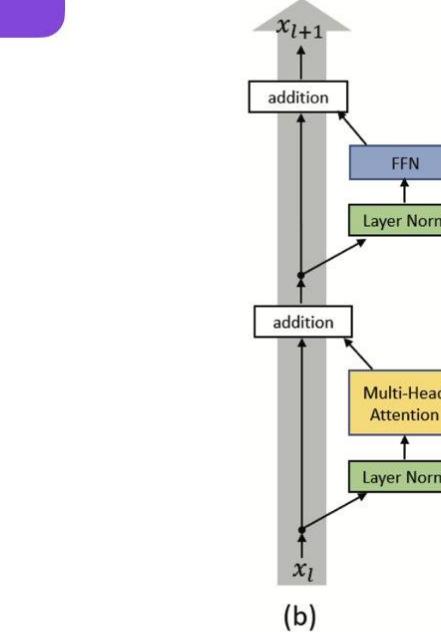


Figure 1: The Transformer - model architecture.

Attention Is All You Need, <https://arxiv.org/abs/1706.03762>



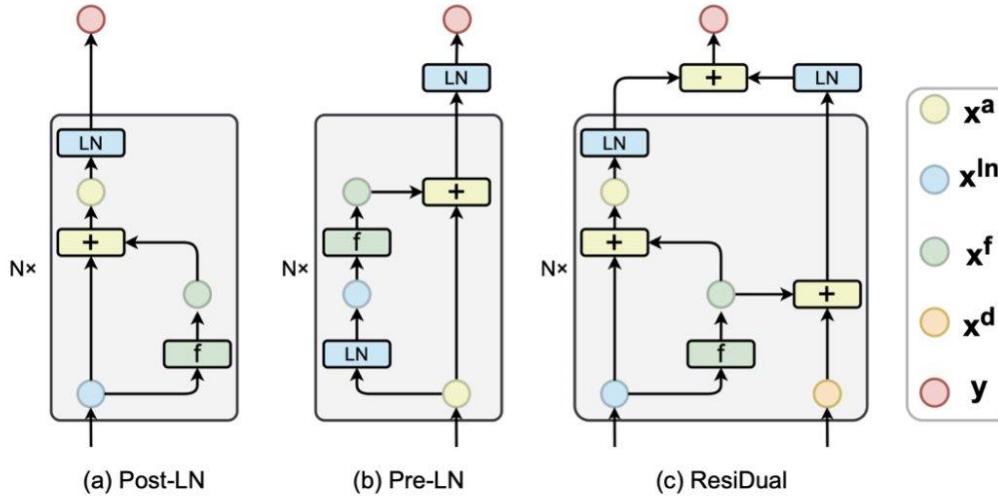
On Layer Normalization in the Transformer Architecture,
<https://arxiv.org/abs/2002.04745>

Method	Gradient Vanishing	Representation Collapse
Post-LN	:(sad face)	:(smiling face)
Pre-LN	:(smiling face)	:(sad face)

ResiDual: Transformer with Dual Residual Connections,
<https://arxiv.org/abs/2304.14802>

New Fun Fact (Found in May 2023)

Dual-Residual Transformer



No representation collapse

No vanishing gradient

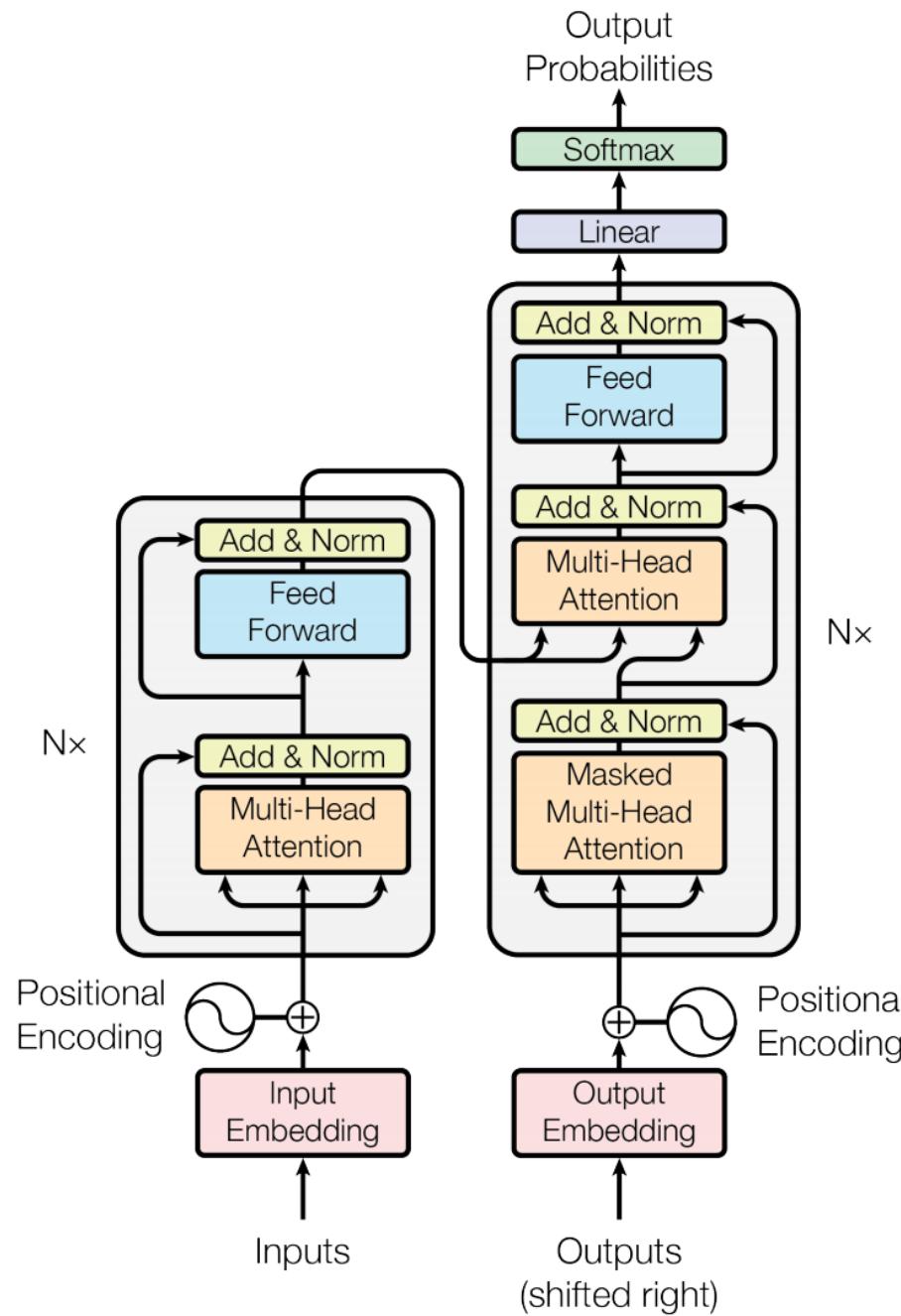
Best of both worlds

Method	Gradient Vanishing	Representation Collapse
Post-LN	:(sad face)	:(smiling face)
Pre-LN	:(smiling face)	:(sad face)
ResiDual	:(smiling face)	:(smiling face)

ResiDual: Transformer with Dual Residual Connections,
<https://arxiv.org/abs/2304.14802>

<https://twitter.com/rasbt/status/1655575611979489282>

Blocks are repeated
N=6 times



Encoder Input

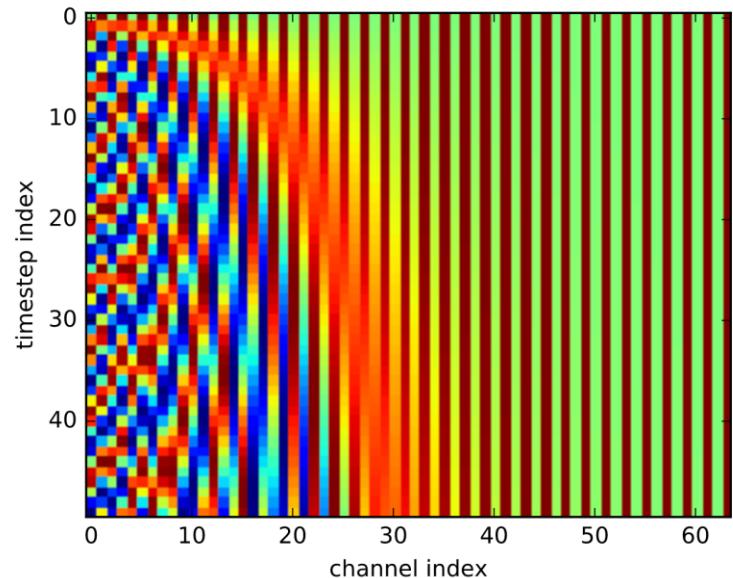
- Added is a positional encoding so same words at different locations have different overall representations:

$$PE_{(pos,2i)} = \sin(pos/10000^{2i/d_{\text{model}}})$$

$$PE_{(pos,2i+1)} = \cos(pos/10000^{2i/d_{\text{model}}})$$

pos is the position of a word

i is the dimension index



https://kazemnejad.com/blog/transformer_architecture_positional_encoding/

Encoder Input

- Actual word representations are byte-pair encodings
 - Rico Sennrich, Barry Haddow, and Alexandra Birch. Neural Machine Translation of Rare Words with Subword Units. ACL 2016.
- Need to handle **large, open vocabulary**
 - Rich morphology:
 - nejneobhospodařovávatelnějšímu - Czech = “to the worst farmable one”
 - Donaudampfschiffahrtsgesellschaftskapitän – German = Danube steamship company captain
 - Informal spelling: goooooood morning !!!!!
- Need to be able to operate at sub-word levels!

Original Byte Pair Encoding for Data Compression

- Suppose the data to be encoded is
 - aaabdaaabac
- The byte pair "aa" occurs most often, so it will be replaced by a byte that is not used in the data, "Z". Now there is the following data and replacement table:
 - ZabdZabac
 - Z=aa
- Then the process is repeated with byte pair "ab", replacing it with Y:
 - ZYdZYac
 - Y=ab
 - Z=aa
- The only literal byte pair left occurs only once, and the encoding might stop here. Or the process could continue with recursive byte pair encoding, replacing "ZY" with "X":
 - XdXac
 - X=ZY
 - Y=ab
 - Z=aa
- This data cannot be compressed further by byte pair encoding because there are no pairs of bytes that occur more than once.

Byte Pair Encoding used in NLP

- A compression algorithm:
 - Most frequent byte pair \mapsto a new byte.
 - Replace bytes with character ngrams
- A word segmentation algorithm:
 - Start with a vocabulary of characters.
 - Most frequent ngram pairs \mapsto a new ngram.

Dictionary

5 low
2 lower
6 newest
3 widest

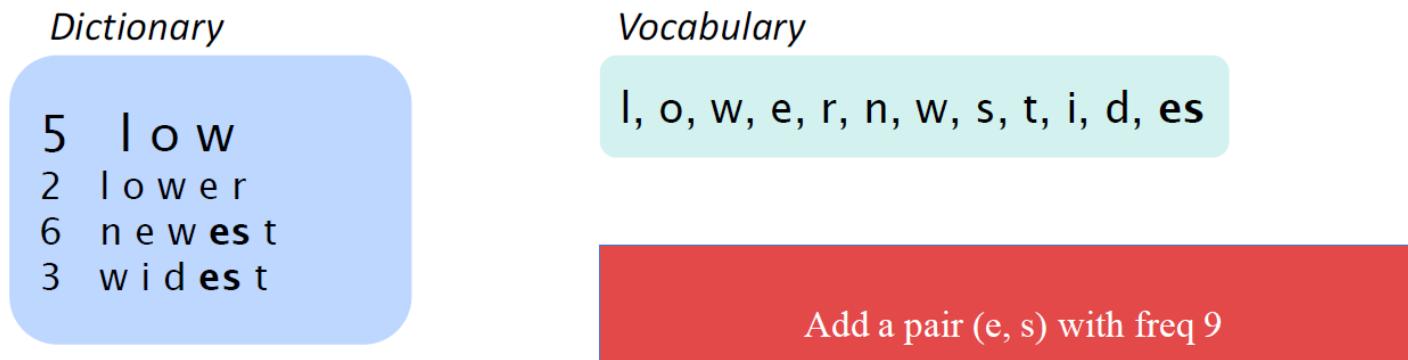
Vocabulary

I, o, w, e, r, n, w, s, t, i, d

Start with all characters in vocab

Byte Pair Encoding used in NLP

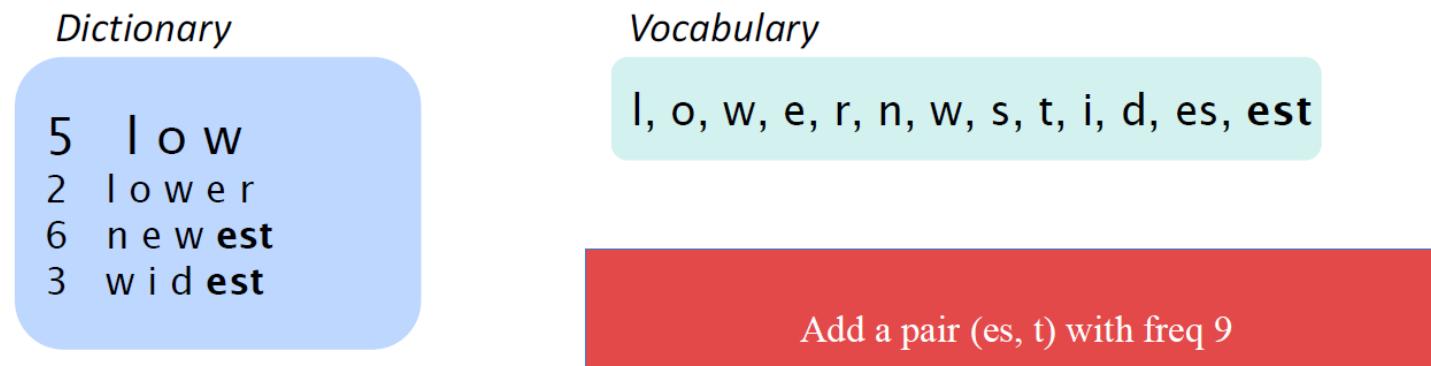
- A compression algorithm:
 - Most frequent byte pair \mapsto a new byte.
 - Replace bytes with character ngrams
- A word segmentation algorithm:
 - Start with a vocabulary of characters.
 - Most frequent ngram pairs \mapsto a new ngram.



(Example from Sennrich)

Byte Pair Encoding used in NLP

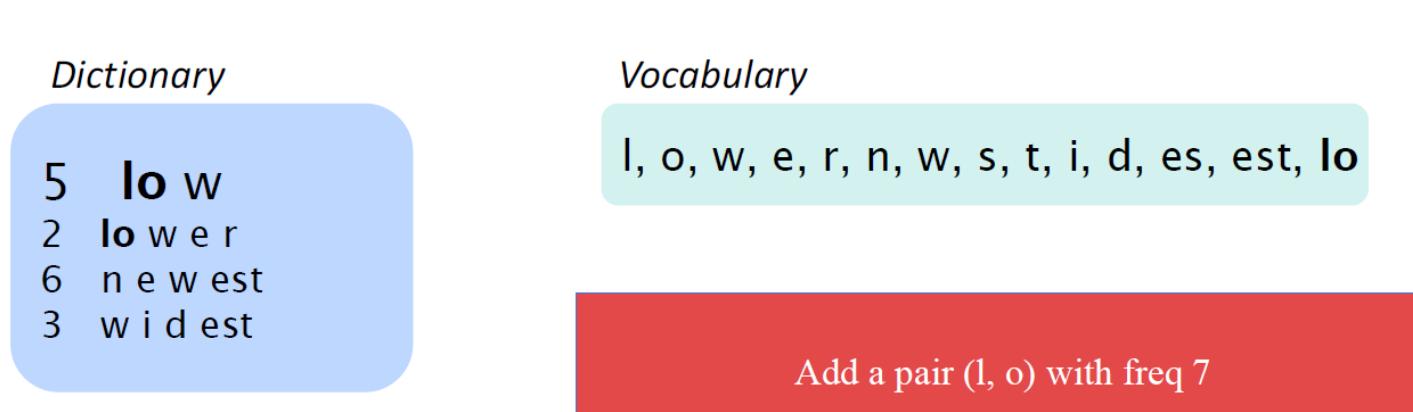
- A compression algorithm:
 - Most frequent byte pair \mapsto a new byte.
 - Replace bytes with character ngrams
- A word segmentation algorithm:
 - Start with a vocabulary of characters.
 - Most frequent ngram pairs \mapsto a new ngram.



(Example from Sennrich)

Byte Pair Encoding used in NLP

- A compression algorithm:
 - Most frequent byte pair \mapsto a new byte.
 - Replace bytes with character ngrams
- A word segmentation algorithm:
 - Start with a vocabulary of characters.
 - Most frequent ngram pairs \mapsto a new ngram.



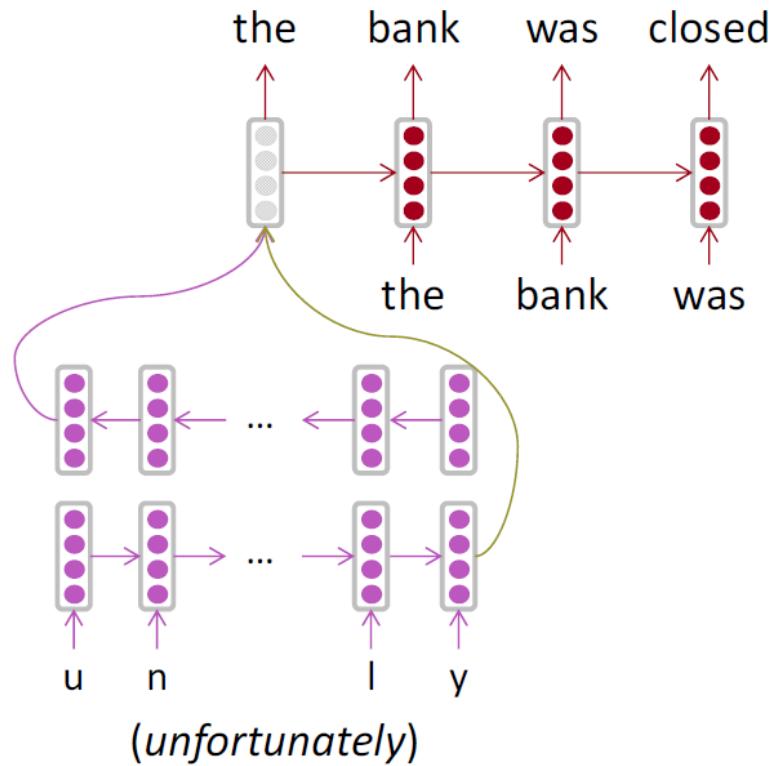
(Example from Sennrich)

Byte Pair Encoding used in NLP

- Automatically decide vocabs for NMT
 - Top places in WMT 2016!
 - <https://github.com/rsennrich/nematus>
- Google’s WordPiece Embeddings (Google Neural Machine Translation System (Wu et al. (2016)))
 - **Word:** Jet makers feud over seat width with big orders at stake
 - **wordpieces:** _J et _makers _fe ud _over _seat _width _with _big _orders _at _stake
 - The word “Jet” is broken into two wordpieces “_J” and “et”, and the word “feud” is broken into two wordpieces “_fe” and “ud”. The other words remain as single wordpieces. “_” is a special character added to mark the beginning of a word.
 - The wordpiece model is generated using a data-driven approach to maximize the language-model likelihood of the training data, given an evolving word definition.

Many ways of using subwords

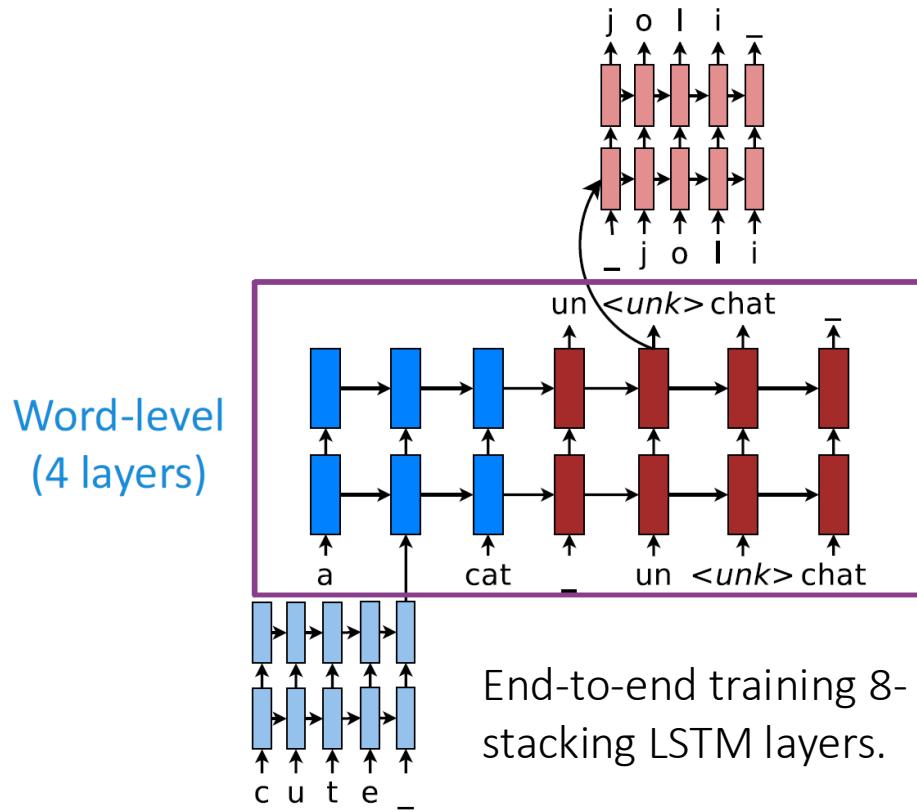
- Character-based LSTM



Many ways of using subwords

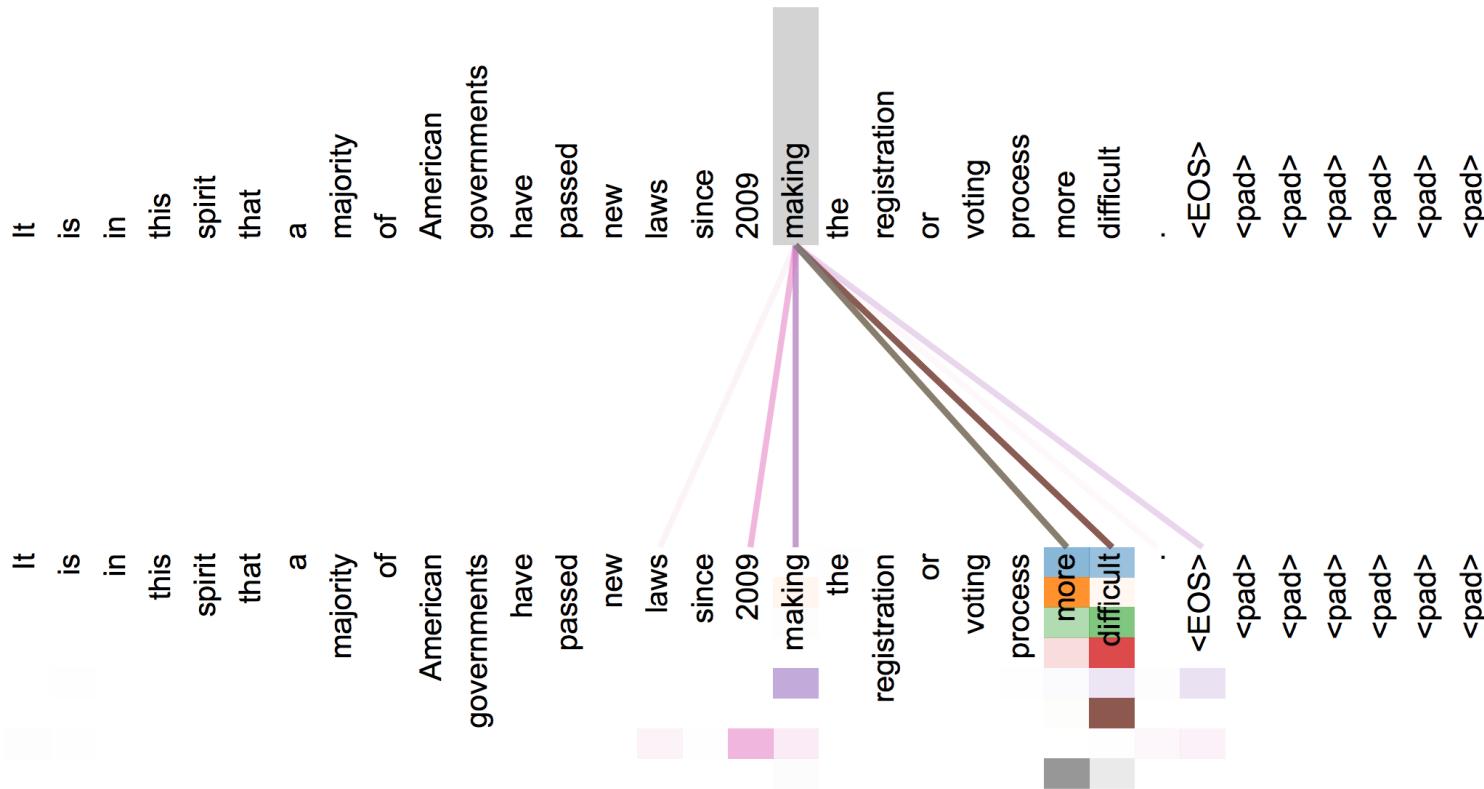
- **Hybrid NMT**

- A *best-of-both-worlds* architecture:
 - Translate mostly at the word level
 - Only go to the character level when needed.
- Word-level beam search
- Char-level beam search for <unk>.
- More than 2 BLEU improvement over a copy mechanism.



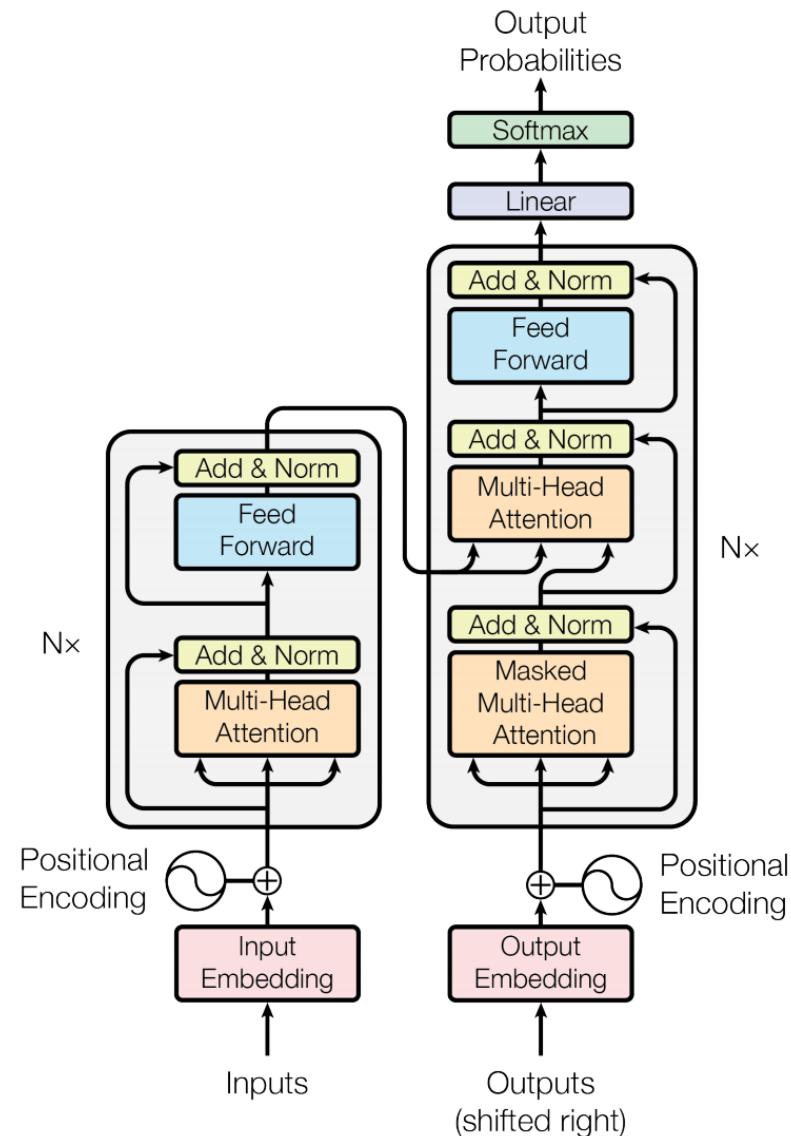
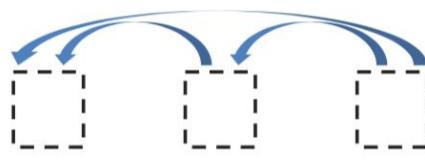
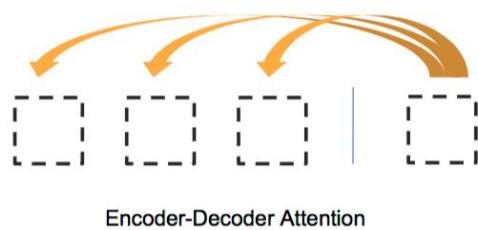
Self Attention Visualization in Layer 5

- Words start to pay attention to other words in sensible ways



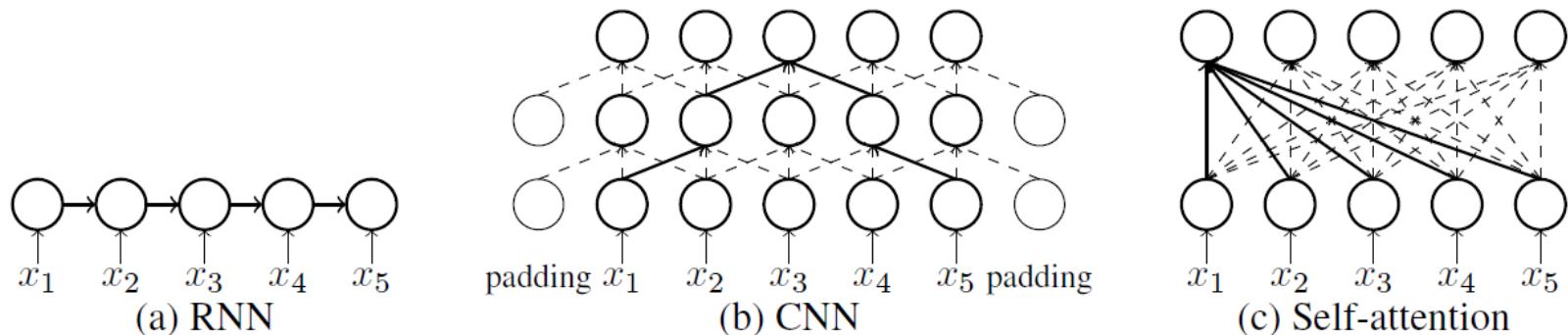
Transformer Decoder

- 2 sublayer changes in decoder
- Masked decoder self-attention
 - Only depends on previous words
- Encoder-Decoder Attention
 - Queries come from previous decoder layer and keys and values come from output of encoder



Advantages

- No recurrence: parallel encoding
- Fast training: both encoder and decoder are parallel
- No long range problem: $O(1)$ for all tokens direct connections
- Three attentions: the model does not have to remember too much
- Multi-head attention allows to pay attention to different aspects
- Comparison of self-attention and CNN/RNN



A comparison of RNN, CNN, and self-attention
<http://aclweb.org/anthology/D18-1458>

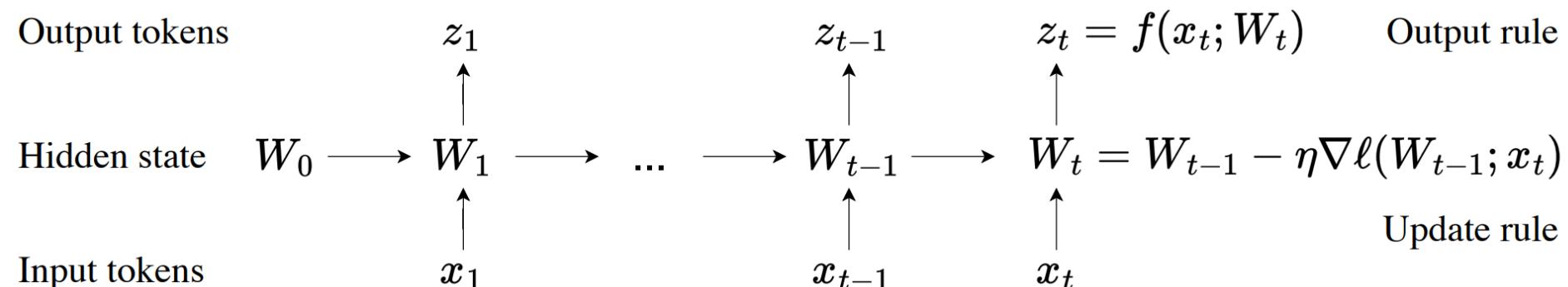
Transformer Networks Visualization

<https://ai.googleblog.com/2017/08/transformer-novel-neural-network.html>

Test Time Training (TTT)

- The hidden state is updated by training even on test sequences
- Loss: denoising autoencoders
 - process x into a corrupted input (low-rank projection)

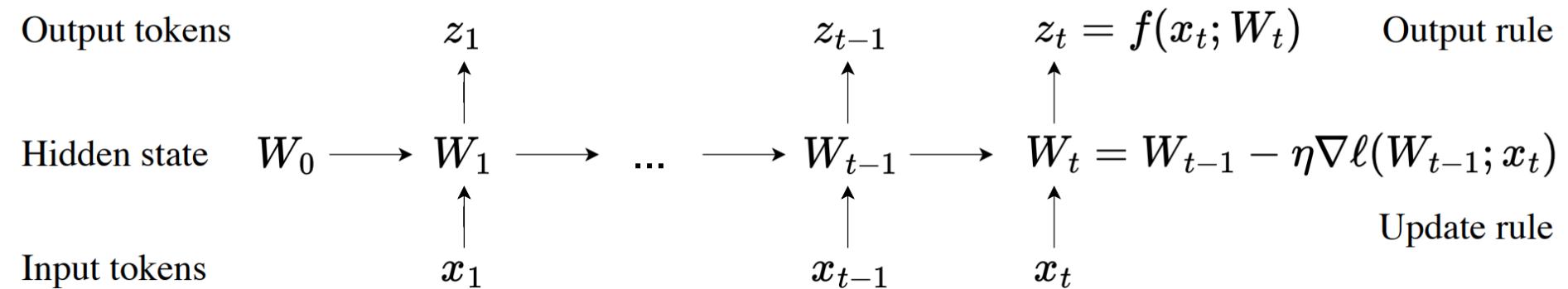
$$\ell(W; x_t) = \|f(\tilde{x}_t; W) - x_t\|^2.$$



Test Time Training (TTT)

- Input corruption by projection using θ_V
- Reconstruction label by another low-rank projection θ_K
- Loss: $\ell(W; x_t) = \|f(\theta_K x_t; W) - \theta_V x_t\|^2.$
- Output decoding by another projection: θ_Q

$$z_t = f(\theta_Q x_t; W_t).$$



Test Time Training (TTT)

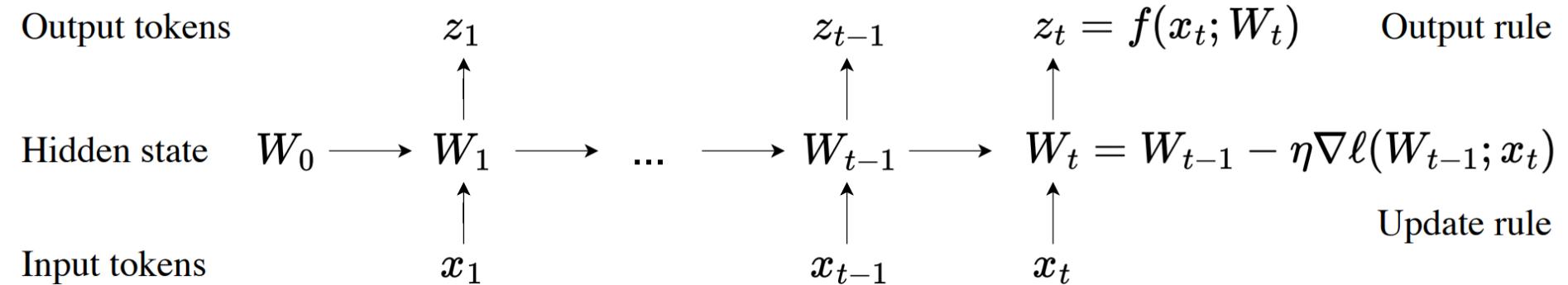
- Loss: $\ell(W; x_t) = \|f(\theta_K x_t; W) - \theta_V x_t\|^2.$
- Updating rule:

$$\nabla \ell(W_0; x_t) = -2(\theta_V x_t)(\theta_K x_t)^T$$

$$W_t = W_{t-1} - \eta \nabla \ell(W_0; x_t) = W_0 - \eta \sum_{s=1}^t \nabla \ell(W_0; x_s) = \sum_{s=1}^t (\theta_V x_s)(\theta_K x_s)^T.$$

$$z_t = f(\theta_Q x_t; W_t) = \sum_{s=1}^t (\theta_V x_s)(\theta_K x_s)^T (\theta_Q x_t)$$

Exactly linear transformer



Linear Attention

- Original Transformer

$$\mathbf{Q} = x\mathbf{W}_Q, \quad \mathbf{K} = x\mathbf{W}_K, \quad \mathbf{V} = x\mathbf{W}_V,$$

$$\mathbf{y}_i = \sum_{j=1}^i \frac{\exp\left(\mathbf{Q}_i^\top \mathbf{K}_j / \sqrt{d_{\text{in}}}\right) \mathbf{V}_j}{\sum_{\ell=1}^i \exp\left(\mathbf{Q}_i^\top \mathbf{K}_\ell / \sqrt{d_{\text{in}}}\right)},$$

- Linearization

$$\mathbf{y}_i = \sum_{j=1}^i \frac{\phi(Q_i^\top K_j)}{\sum_{\ell=1}^i \phi(Q_i^\top K_\ell)} V_j = \sum_{j=1}^i \frac{\phi(Q_i)^\top \phi(K_j)}{\sum_{\ell=1}^i \phi(Q_i)^\top \phi(K_\ell)} V_j = \frac{\phi(Q_i)^\top \sum_{j=1}^i \phi(K_j) V_j}{\phi(Q_i)^\top \sum_{\ell=1}^i \phi(K_\ell)},$$

- Memory matrix

$$\mathcal{M}_t = \mathcal{M}_{t-1} + K_t^\top V_t,$$

$$\mathbf{y}_t = Q_t \mathcal{M}_t,$$

Further Nonlinearization

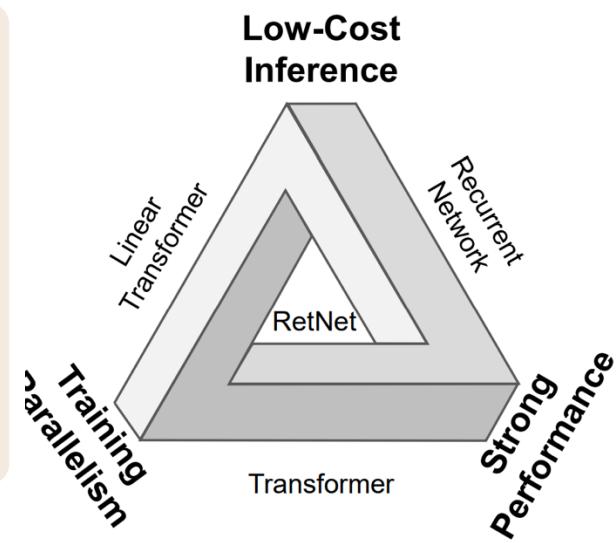
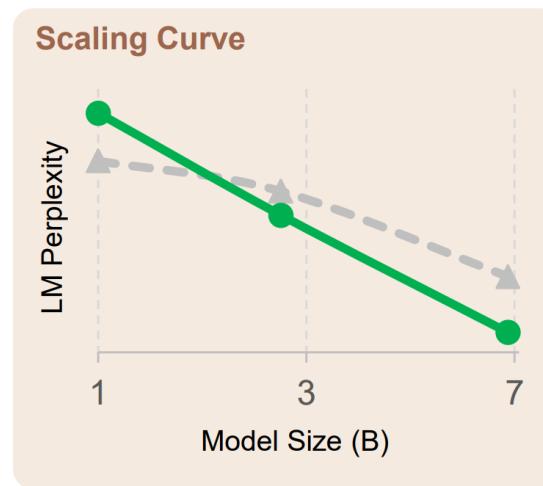
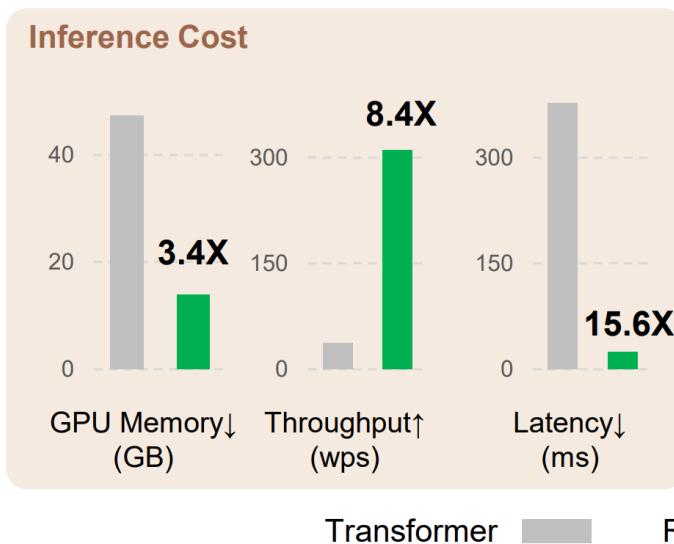
- Memory matrix $\mathcal{M}_t = \mathcal{M}_{t-1} + K_t^\top V_t ,$
 $\mathbf{y}_t = Q_t \mathcal{M}_t ,$
- Generalization
 $\mathcal{M}_t = f(\mathcal{M}_{t-1}, x_t),$ Write Operation
 $\mathbf{y}_t = g(\mathcal{M}_t, x_t),$ Read Operation

- We can use any kernel function here

$$\kappa(x, x'; \theta_K, \theta_Q) \propto e^{(\theta_K x)^T \theta_Q x'}$$

The “Impossible Triangle” of Scaling

- Achieving low-cost inference, efficient long sequence modeling



Test Time Scaling

- Still a very hot topic
 - LLM architecture design, e.g., TTT
 - Reasoning models, e.g., R1, o3
- No SOTA LLM is using TTT so far
 - Possibly because of the sunk cost
 - But it's a promising direction, as the inference cost has become a bottleneck for LLM service